# Context-Aware Semantic Inpainting

Haofeng Li, Guanbin Li, Liang Lin, and Yizhou Yu

*Abstract*—Recently image inpainting has witnessed rapid progress due to generative adversarial networks (GAN) that are able to synthesize realistic contents. However, most existing GAN-based methods for semantic inpainting apply an auto-encoder architecture with a fully connected layer, which cannot accurately maintain spatial information. In addition, the discriminator in existing GANs struggle to understand high-level semantics within the image context and yield semantically consistent content. Existing evaluation criteria are biased towards blurry results and cannot well characterize edge preservation and visual authenticity in the inpainting results. In this paper, we propose an improved generative adversarial network to overcome the aforementioned limitations. Our proposed GAN-based framework consists of a fully convolutional design for the generator which helps to better preserve spatial structures and a joint loss function with a revised perceptual loss to capture high-level semantics in the context. Furthermore, we also introduce two novel measures to better assess the quality of image inpainting results. Experimental results demonstrate that our method outperforms the state of the art under a wide range of criteria.

*Index Terms*—Image Completion, Image Inpainting, Convolutional Neural Network.

## I. INTRODUCTION

IMAGE inpainting aims at synthesizing the missing or damaged parts of an image. It is a fundamental problem in low-level vision and has attracted widespread interest in the computer vision and graphics communities as it can be used for filling occluded image regions or repairing damaged photos. Due to the inherent ambiguity of this problem and the complexity of natural images, synthesizing content with reasonable details for arbitrary natural images still remains a challenging task.

High-quality inpainted result should be not only realistic but also semantically consistent with the image context surrounding the missing or damaged region at different scales. First, colorization should be reasonable and spatially coherent. Second, structural features such as salient contours and edges should be connected inside the missing region or across its boundary. Third, texture generated within the missing region should be consistent with the image context and contains high-frequency details. In addition, missing object parts need to be recovered correctly, which is challenging and requires capturing high-level semantics.

Deep convolutional neural networks are capable of learning powerful image representations and have been applied to inpainting [3], [4] with varying degrees of success. Recently semantic image inpainting has been formulated as an image

H. Li and Y. Yu are with the Department of Computer Science, The University of Hong Kong, HK

G. Li and L. Lin are with the school of Data and Computer Science, Sun Yat-sen University, Guangzhou 510006, China.



(a) Input　　　　　　(b) CASI
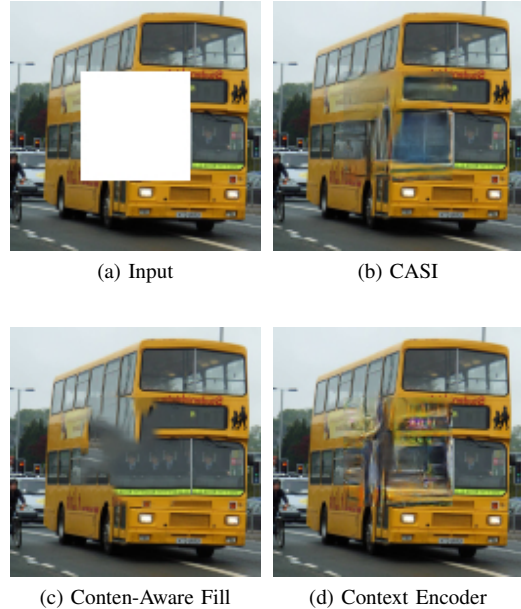
(c) Conten-Aware Fill　　(d) Context Encoder

Fig. 1: Our proposed CASI with perceptual loss synthesizes content with a more reasonable colorization and structure than Content-Aware Fill [1] and Context Encoder [2].

generation problem and solved within the framework of generative adversarial networks (GAN) [5]. GAN trains a generator against a discriminator and successfully generates plausible visual content with sharp details. State-of-the-art results [2], [6], [7] have been achieved.

However, all existing GAN-based solutions to inpainting share common limitations. First of all, they utilize an encoder-decoder architecture with fully connected layers as the *bottleneck* structure in the middle of the network. The bottleneck structure contains two fully connected (fc) layers. The first fc layer converts convolutional features with spatial dimensions to a single feature vector and another fc layer maps the feature vector backward to features with spatial information. The first fully connected layer collapses the spatial structure of the input image so that location related information cannot be accurately recovered during the decoding process. Second, the discriminator only takes a synthesized region without its image context as the input. Thus neither structural continuity nor texture consistency can be guaranteed between the synthesized region and its image context. Moveover, existing GANs struggle to understand high-level semantics within the image context and yield semantically consistent content.

To overcome the aforementioned limitations, we conceive a novel fully convolutional generative network for semantic inpainting. First, we adopt a fully convolutional design without

the bottleneck structure to preserve more spatial information. Second, we composite the synthesized region and its image context together as a whole, and measures the similarity between this composite image and the ground truth. To increase such similarity, a perceptual loss is computed for the composite image. This perceptual loss defined in terms of high-level deep features is promising in capturing the semantics of the image context.

Furthermore, noticing that the $L2$ loss and PSNR are unable to rate blurry results accurately and quantitative measures do not exist for assessing how well the intended semantics have been restored, we define a local entropy error and a semantic error to resolve these two issues respectively. The semantic error (SME) is defined as the hinge loss for the confidence that a composite image with a synthesized region should be assigned the groundtruth label of its real counterpart, where the confidence value is estimated by a pre-trained image classifier. In our experiments, images synthesized by our inpainting model can successfully reduce the semantic error estimated by a powerful image classifier. This indicates that our model is capable of inferring semantically valid content from the image context.

In summary, this paper has the following contributions:

- We present a fully convolutional generative adversarial network without a fully-connected layer for maintaining the original spatial information in the input image. This network can process images with a variable size.
- We introduce a novel context-aware loss function including a perceptual loss term, which measures the similarity between a composite image and its corresponding groundtruth real image.
- We propose two novel measures, a local entropy error based on middle-level statistics and a semantic error based on high-level features, for evaluating inpainting results.

## II. RELATED WORK

Recently, deep neural networks including generative adversarial networks have exhibited great performance in image generation, image transformation and image completion. This section discusses previous work relevant to image inpainting and our proposed method.

### A. Image Inpainting

Many algorithms on recovering holes in images or videos have been proposed [8], [9], [10], [11], [12], [13], [14], [15]. Some existing methods for image completion are related to texture synthesis [16], [17] or patch-based synthesis [18], [19], [20]. Efros and Leung [16] proposed a method for predicting pixels from the context boundary while [17] searches for matching patches and quilts them properly. Drori et.al. [21] computed a confidence map to guide filling while Komodakis et.al. [22] proposed a priority belief propagation method. However, these exemplar based approaches struggle to generate globally consistent structures despite producing seamless high-frequency textures. Hays and Efros [23] filled large missing regions using millions of photographs and presented seamless

results. However, in this method, missing regions need to be prepared carefully by completely removing partially occluded objects. Synthesizing content for arbitrary missing regions remains a challenging task (e.g., recovering body parts for a partially occluded object).
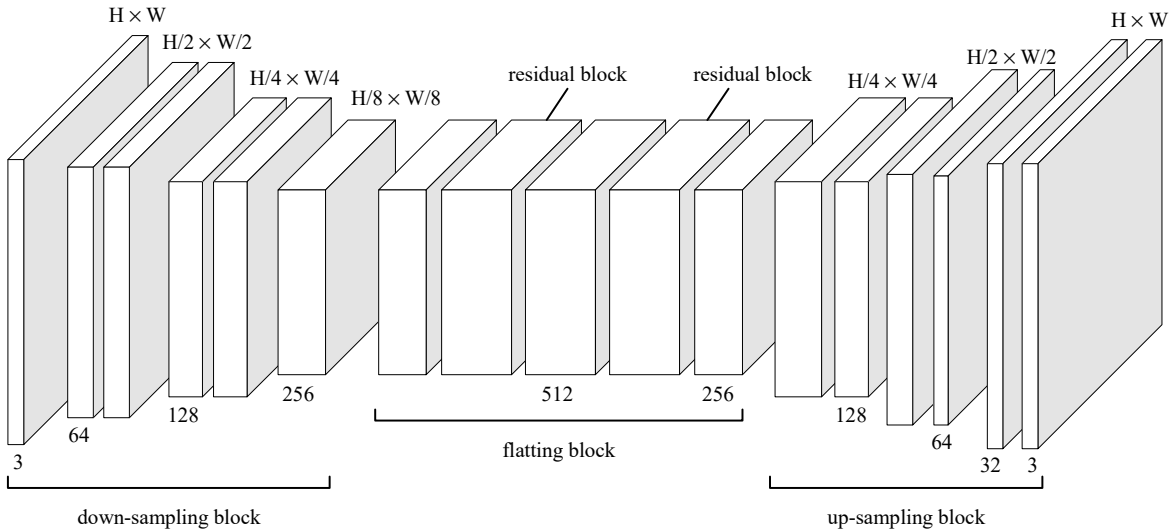
### B. Generative Adversarial Networks

Generative adversarial networks (GAN), which estimate generative models by simultaneously training two adversarial models were first introduced by Goodfellow et.al. [5] for image generation. Radford et.al. [24] further developed a more stable set of architectures for training generative adversarial networks, called deep convolutional generative adversarial networks (DCGAN). Recently GAN has widely applied to image generation [25], image transformation[26], image completion [2] and texture synthesis [27]. Context Encoder [2] uses a novel channel-wise fully connected layer for feature learning but keeps the traditional fully connected layer for semantic inpainting. Yeh et.al. [6] employed GAN with both a perceptual loss and a contextual loss to solve inpainting. Notice that the perceptual loss in [6] is essentially an adversarial loss and the contextual loss considers the context only (excluding the synthesized region). Yang et.al. [7] conducted online optimization upon a pre-trained inpainting model primarily inherited from Context Encoder. The optimization is too expensive for real-time or interactive applications. Common disadvantages exist in these GAN based approaches. First, the fully connected layer in the encoder-decoder framework cannot preserve accurate spatial information. Second, the discriminator in current GANs only evaluates the synthesized region but not the semantic and appearance consistency between the predicted region and the image context.

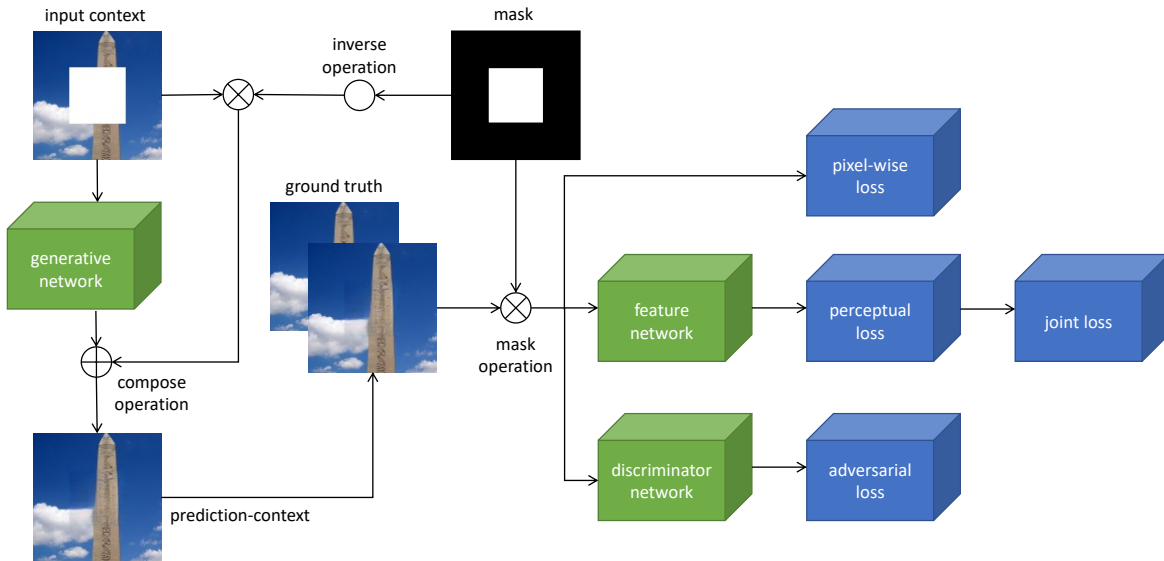### C. Fully Convolutional Networks

Fully convolutional networks (FCNs), which was first used in [28] for semantic image segmentation, provides an end-to-end learnable neural network solution for pixel-level image comprehension. Without fully connected layers, FCNs occupy less memory and can learn and predict more efficiently. Besides, FCNs preserve spatial information and extract location sensitive features. Recently FCNs have achieved excellent results on semantic segmentation [28], edge detection [29], saliency detection [30] and other pixel-wise labeling tasks. In this paper, we exploit the idea of FCN in GAN-based inpainting to better capture object contours, preserve spatial information in features, and infer coherent visual content from context.

### D. Context-Aware Perceptual Loss

Perceptual loss is a feature reconstruction loss defined by deep neural networks [31]. It guides neural models to generate images visually similar to their corresponding targets (e.g., ground truth) and has been widely utilized in style transfer [32]. Dosovitskiy and Brox et.al. [25] presented a similar concept, called DeePSiM, which successfully generates images with sharp details. So far perceptual loss has been

(a) Fully Convolutional Generative Network



(b) Context-Aware Semantic Inpainting Pipeline

Fig. 2: Network Architecture

applied to style transfer [32], [31], super resolution [31] and texture synthesis[33]. However, these topics primarily use the "texture network", a part of the VGG network [34] to extract middle-level features while high-level features from the fully connected layers have not been investigated for image completion. In this paper we exploit high-level deep features in the definition of perceptual loss to synthesize regions semantically consistent with their contexts.

## III. METHOD

As shown in Figure 2b, our proposed Context-Aware Semantic Inpainting method (CASI) is composed of an inpainting generation pipeline (on the left) and a joint loss function (on the right). The fully convolutional generative network takes an image context as the input, where the missing region is filled with the mean pixel value. The missing region is generated by point-wise multiplication (denoted as 'mask operation') with

a mask. The inverse operation turns one into zero, and zero into one. The output of the generative network is a synthesized image with the same size as the input. Then this output image is cropped using the boundary of the missing region and placed within the image context to form a composite image (denoted as 'prediction-context'), via a point-wise addition (denoted as 'compose operation'). The discriminator network receives the synthesized content within the missing region and the ground truth within the same region respectively, and attempts to classify the received content as either 'real' or 'fake'. The classification error is formulated as the adversarial loss, one of the components in the proposed loss. Our joint loss function is a linear combination of a pixel-wise $L2$ loss, the adversarial loss and a perceptual loss.

### A. Fully Convolutional Generative Network

The fully convolutional generative network consists of three blocks: down-sampling, flatting and up-sampling. First, the down-sampling block plays the role of an encoder, which reduces each spatial dimension to 1/8 of the input size. The flatting block discovers and maintains essential edges without further changing the spatial size. Finally, the up-sampling block plays the role of a decoder, which transforms the feature map to an RGB image with the same resolution as the input.

The down-sampling block has three convolutional layers using 4×4 kernels and two convolutional layers using 3×3 kernels. The first layer of this block performs 4×4 convolution. Then these two types of convolutional layers alternate and the block ends with a 4×4 convolutional layer. The 4×4 convolutions use a stride of 2 and 1 pixel padding to reduce the spatial size by half while doubling the number of channels in the feature map. Reduced spatial dimensions allow convolution kernels to have larger receptive fields in the input image. The 3×3 convolutions use a stride of 1 and 1 pixel padding to keep the same spatial size and channel number. Such layers enhance the recognition capacity of the network. The flatting block has three convolutional layers using 3×3 kernels and two residual blocks. These residual blocks enhance prediction accuracy for semantic inpainting. The middle layer doubles the number of channels while the last layer reduces it by half. Thus the flatting block keeps the number of channels the same in the input and output feature maps. The up-sampling block has three de-convolutional layers using 4×4 kernels and three convolutional layers using 3×3 kernels. Similar to the down-sampling block, the two types of layers alternate, and the first layer performs 4×4 deconvolution. In the up-sampling block, 4×4 deconvolution acts as parameterized interpolation which doubles the spatial size while each 3×3 convolutional layer reduces the number of channels by half. The last layer of the up-sampling block generates an RGB image with the same size as the input.

Our proposed generative network does not have a bottleneck fully connected layer, and enjoys the benefits of fully convolutional architecture. It is capable of locating essential boundaries, maintaining fine details and yield consistent structures in missing regions.

### B. Discriminative Network

Our discriminator shares a similar but shallower structure with the down-sampling block in the generator network. Compared with the down-sampling block, the discriminator removes all 3×3 convolutional layers to avoid overfitting. Otherwise, the capacity of the discriminator would be so large that the generator does not have a chance to confuse the discriminator and improve itself. A fully connected layer is employed to perform binary classification at the end of the discriminator.

Normalization and non-linear activations are used in CASI. Except for the last layer, every convolutional layer in the generator and the discriminator is followed with a batch normalization (batchnorm) layer. Rectified linear units (ReLU)

follow each batchnorm layer in the generator while Leaky-rectified Linear Units (LeakyReLU) are used in the discriminator according to the architecture guidelines in DCGAN. A Sigmoid layer is adopted in the last layer of the generator and the discriminator to map pixel and confidence values respectively.

### C. Loss Function

Given the analysis in Section I, existing GAN based semantic inpainting methods fail to grasp high-level semantics and synthesize semantically consistent content for the missing region. In this paper, we propose to composite the synthesized region and its image context together as a whole, and measures the visual similarity between this composite image and the ground truth using a perceptual loss. Our overall loss function consists of a pixel-wise $L2$ loss, an adversarial loss and a perceptual loss. It can be formulated as follows,

$$L_{inp} = \lambda_{pix} l_{pix} + \lambda_{adv} l_{adv} + \lambda_{per} l_{per}, \quad (1)$$

where $L_{inp}$ denotes the overall inpainting loss. $l_{per}$, $l_{adv}$, $l_{pix}$ denote our perceptual loss, adversarial loss and pixel-wise $L2$ loss respectively while $\lambda_{per}$, $\lambda_{adv}$ and $\lambda_{pix}$ are the weights of the respective loss terms.

Pixel-wise $L2$ loss, $l_{pix}$, is a straightforward and widely used loss in image generation. It measures the pixel-wise differences between the synthesized region and its corresponding ground truth. $l_{pix}$ is defined in Eq. (2),

$$l_{pix}(x, z) = ||M \odot (x - z)||_2{}^2, \quad (2)$$

where $M$ is a binary mask where a value of 1 indicates the missing region and a value of 0 indicates the known context region, $\odot$ is the element-wise product, $x$ is the ground-truth image and $z$ is the corresponding inpainting result computed as in Eq. (3),

$$z = ((1 - M) \odot x) \oplus (M \odot G((1 - M) \odot x)), \quad (3)$$

where $\oplus$ is the element-wise addition, G is the CASI generator, $(1 - M) \odot x$ is the context region of $x$, and $M \odot G(\cdot)$ is the missing region in the generator's output. $\oplus$ in Eq. (3) merges the known context region and the synthesized missing region to obtain the final inpainting result.

However, calculating loss within the image space cannot guarantee to generate an image perceptually similar to the ground truth as neural networks tend to predict pixel values close to the mean of the training data. In practice, the pixel-wise $L2$ loss only produces blurred images without clear edges or detailed textures. Thus we exploit an adversarial loss and a novel perceptual loss to overcome this problem.

The adversarial loss $l_{adv}$ is defined on the objective function of the discriminator. As the discriminator aims at distinguishing synthesized content from its corresponding ground truth, its objective is to minimize a binary categorical entropy $e$ in Eq. (4).

$$\begin{aligned} &e(D(M \odot x), D(M \odot z)) \\ &= -[log(D(M \odot x)) + log(1 - D(M \odot z))], \end{aligned} \quad (4)$$

where $e$ denotes binary categorical entropy and $D$ is the CASI discriminator. The discriminator $D$ predicts the probability that the input image is a real image rather than a synthesized one. If the binary categorical entropy is smaller, the accuracy of the discriminator is better. Note that $D$ is not a pre-trained or constant model during the training stage. Instead, $G$ and $D$ are trained alternatively. As minimizing the binary categorical entropy $e$ is equivalent to maximizing the negative of the binary categorical entropy, the final objective value of the discriminator is described in the right side of Eq. (5). As the generator acts as an adversarial model of the discriminator, it tends to minimize the negative of the binary categorical entropy. Thus the adversarial loss of the generator $l_{adv}$ can be formally described as

$$l_{adv} = \max_D[log(D(M \odot x)) + log(1 - D(M \odot z))]. \quad (5)$$

$l_{adv}$ makes the synthesized region deviate from the overly smooth result obtained using the pixel-wise $L2$ loss as real images are not very smooth and typically have fine details. Although the adversarial loss promotes fine details in the synthesized result, it also has disadvantages. First, existing discriminators are unaware of the image context and do not explicitly consider the composite image consisting of both the synthesized region and the image context. Second, binary classification is not challenging enough for the discriminator to learn the appearance of different objects and parts. Note that semantic inpainting needs to not only synthesize textures consistent with the context but also recover missing object parts, which requires high-level features extracted from the image context. Thus we propose a perceptual loss based on high-level semantic features.

Our perceptual loss, $l_{per}$, is defined in Eq. (6),

$$l_{per}(x, z) = e(F(x), F(z))$$
$$= \frac{1}{C_j H_j W_j}||F_j(x) - F_j(z)||_2^2, \quad (6)$$

where $F$ is a pre-trained feature network that extracts a generic global feature from the input, $F_j$ denotes the activations of the $j$-th layer of $F$, $F_j(x)$ and $F_j(z)$ are a $C_j \times H_j \times W_j$ tensor respectively. In our experiments, we use ResNet-18 pre-trained over the ImageNet dataset [35] as the feature network $F$, and the 512-dimensional feature from the second last layer of ResNet-18 as $F_j$. Similar high-level features extracted by $F$ give rise to similar generated images, as suggested in [25]. In addition, a perceptual loss based on high-level features makes up for the missing global information typically represented in a fully connected layer in the generator. Different from DeepSiM, our feature is extracted from the composite image consisting of the synthesized region and the image context rather than from the synthesized region alone.

## IV. IMPLEMENTATION

Let us discuss the details of our inpainting pipeline. Training images for CASI require no labels. As shown in Algorithm 1, the training stage consists of a limited number of iterations. During each training iteration, the discriminator is updated $Diters$ times and the generator is trained once. In each

iteration that updates the discriminator, each training image is separated into an image center and an image context. The image center has the same size of the central region, and the image context is the image filled with the mean pixel value in the central region. The image center and image context of a training image form a training pair. The generator takes the image context as the input and synthesizes the image center. The discriminator attempts to distinguish the synthesized content from the ground-truth image center. The adversarial loss is calculated and then the parameters of the discriminator are updated. In the rest of each training iteration, the pixel-wise $L2$ loss is computed, the feature network extracts a feature from the composite image, and three loss functions are combined to obtain the joint inpainting loss. The generator is finally updated according to the joint loss. This process is repeated until the joint loss converges. In the testing stage, each testing image is first filled with the mean pixel value in the center and then passed to the CASI generator. The central region of the generator's output is cropped and pasted back into the testing image to yield the final inpainting result.

Our CASI is implemented on top of DCGAN [24] and Context Encoder [2] in Torch and Caffe [36]. ADAM [37] is adopted to perform stochastic gradient descent. As in [2], CASI predicts a larger region which overlaps with the context region (by 4px). $10\times$ weight is used for the pixel-wise $L2$ loss in the overlapping area. Using a TITAN XP GPU, training on a dataset of 20000 images costs 3 to 4 days. Inpainting a single image takes less than 0.2 seconds. Recovering a batch of 20 images costs less than 1 second.

---

**Algorithm 1**

---

1: $F \leftarrow$ LOADMODEL()
2: $G \leftarrow$ INITWEIGHT()$, D \leftarrow$ INITWEIGHT()
3: **for** $i \leftarrow 1, maxIterations$ **do**
4:     $x, z, M$
5:     **for** $j \leftarrow 1, Diters$ **do**
6:         $x \leftarrow$ SAMPLEBATCH()
7:         Compute $z$ using Eq. (3)
8:         Compute $l_{adv}$ using Eq. (4)
9:         Update $D$
10:     **end for**
11:     $l_{pix} \leftarrow$ MSE$(x, z)$
12:     $f_x \leftarrow F(x), f_z \leftarrow F(z)$
13:     Compute $l_{per}$ using Eq. (6)
14:     Compute $L_{inp}$ using Eq. (1)
15:     Update $G$
16: **end for**

---

## V. EVALUATION

This section evaluates our proposed deep neural network architecture and joint loss function on a subset of ImageNet [35] and the Paris StreetView dataset [2], [38]. This subset contains 20 randomly sampled categories, denoted as "ImageNet-20". ImageNet-20 consists of $25,000$ training images and $1,000$ testing images. Paris StreetView contains $14,900$ training samples and 100 testing samples.

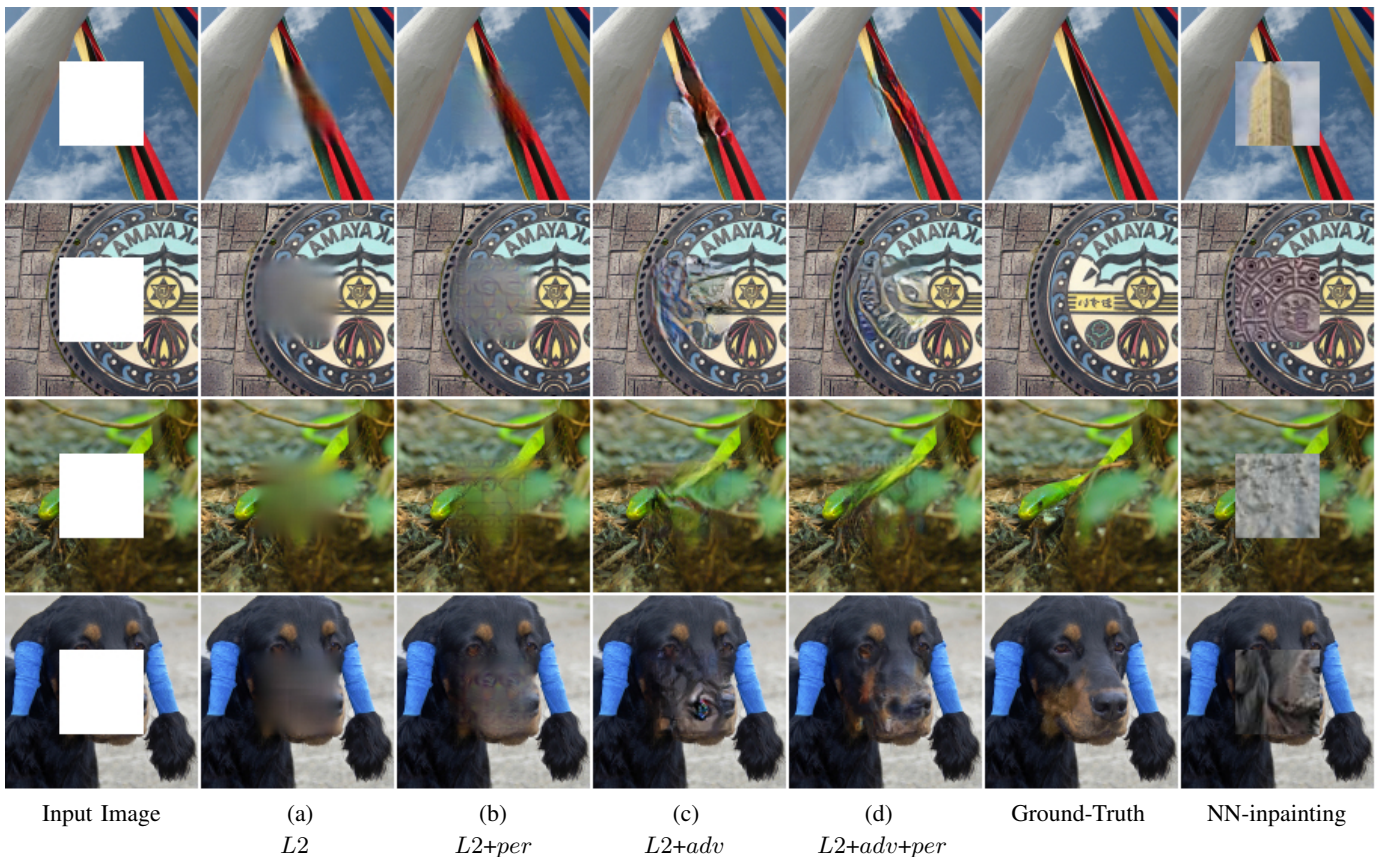| Input Image | (a) | (b) | (c) | (d) | Ground-Truth | NN-inpainting |
| :---: | :---: | :---: | :---: | :---: | :---: | :---: |
| | $L2$ | $L2+per$ | $L2+adv$ | $L2+adv+per$ | | |

Fig. 3: Comparison among different combinations of loss functions and Nearest-Neighbor(NN)-inpainting. The adversarial loss promotes low-level sharp details while the perceptual loss improves high-level semantic consistency.

### A. Effectiveness of Perceptual Loss

We first verify whether adding a perceptual loss improves the results. CASI is trained using 4 different loss functions respectively to compare their performance. For these loss functions, the hyper-parameters of CASI are set in the same way, and the perceptual loss is defined using the same feature extracted using the same feature network. The four loss functions are: (a) pixel-wise $L2$ loss, (b) $L2$ loss + perceptual loss, (c) $L2$ loss + adversarial loss, (d) $L2$ loss + adversarial loss + perceptual loss. In the following we use (a)-(d) to refer to these loss functions.

Figure 3 show qualitative results of the above loss functions. The resolution of each images is $128 \times 128$. This result includes 4 samples representing different cases. All the missing regions are at the center of the image. From left to right, each column corresponds to a loss function from (a) to (d), respectively. As shown in this figure, (a) and (b) generate over-smooth results while (c) and (d) present sharper details. This conforms that the adversarial loss indeed alleviate the blurriness caused by the $L2$ loss. Between (a) and (b), (a) is more blurry while subtle textures or wrinkles can be observed in (b). Between (c) and (d), although they both preserve sharp edges, (d) is more semantically consistent with the context region. These results reveal that the adversarial loss works in the middle level to yield patches with consistent sharp details while the perceptual loss synthesizes consistent high-level contents.

TABLE I: Quantitative results on ImageNet-20. CASIs without the adversarial loss achieve lower mean $L2$ error and higher PSNR but generate blurry results, which indicates that mean $L2$ error and PSNR inaccurately assess over-smooth cases.

| Method | mean $L1$ error | mean $L2$ error | PSNR |
| :--- | :---: | :---: | :---: |
| Context Encoder | 12.15% | 3.31% | 15.59dB |
| CASI,$L2$ | 11.07% | **2.57%** | **17.08dB** |
| CASI,$L2 + per$ | 11.21% | 2.64% | 16.95dB |
| CASI,$L2 + adv$ | 11.15% | 2.93% | 16.68dB |
| CASI,$L2+adv+per$ | **10.89%** | 2.83% | 16.81dB |

Table I shows quantitative results from this experiment. It presents numerical errors between synthesized contents and their ground truth using three commonly employed measures, mean $L1$ error, mean $L2$ error and PSNR. Notations (a)-(d) are used to denote four trained CASI models. As shown in Table I, (a) achieves the smallest mean $L2$ error and PSNR while (d) achieves the smallest mean $L1$ error. Mean $L2$ error is smaller for solutions close to the mean value but such solutions are overly smooth and undesirable (see (a) and (b) in Figure 3). Models trained without the adversarial loss have advantage in mean $L2$ error due to their blurry results. Similar results have been reported in [25]. Between (c) and (d), (d) has smaller

TABLE II: Investigation of Perceptual Loss

| Method | mean $L1$ error | mean $L2$ error | PSNR |
|---|---|---|---|
| $\lambda_{per} = 0$ | 11.15% | 2.93% | 16.68dB |
| $\lambda_{per} = 0.2$ | **10.89%** | **2.83%** | **16.81dB** |
| $\lambda_{per} = 0.4$ | 11.12% | 2.93% | 16.60dB |
| $\lambda_{per} = 0.7$ | 11.43% | 3.06% | 16.44dB |

TABLE III: Effectiveness of Fully Convolutional Architecture

| Method | mean $L1$ error | mean $L2$ error | PSNR |
|---|---|---|---|
| CASI+$fc$ | 9.70% | 1.71% | 18.83dB |
| CASI | **7.49%** | **1.37%** | **20.37dB** |

TABLE IV: Effectiveness of Residual Block

| Method | mean $L1$ error | mean $L2$ error | PSNR |
|---|---|---|---|
| CASI- | 11.09% | 2.93% | 16.31dB |
| CASI | **10.89%** | **2.83%** | **16.81dB** |
| CASI- | 7.79% | 1.43% | 20.14dB |
| CASI | **7.49%** | **1.37%** | **20.37dB** |

TABLE V: High-resolution Case on Paris StreetView

| Method | mean $L1$ error | mean $L2$ error | PSNR |
|---|---|---|---|
| ContextEncoder | 9.04% | 1.82% | 18.90dB |
| CASI | **8.04%** | **1.53%** | **19.79dB** |
| NeuralPatch | 9.59% | 2.07% | 18.42dB |
| CASI+ | **8.62%** | **1.73%** | **19.18dB** |

mean $L2$ error than (c). And (d) also has smaller mean $L1$ error than (c). Thus the perceptual loss is effective in improving our CASI model.

### B. Investigation of Perceptual Loss

This section investigates how the parameter of perceptual loss effect the performance of our method. We set the hyper-parameters in our algorithm as follows. The summation of the weights of all loss terms is 1.0. The weight of the adversarial loss is 0.001, as suggested by [2]. We determine the weight of the perceptual loss $\lambda_{per}$ by cross validation on the ImageNet-20 dataset. As shown in Table II, setting the weight of the perceptual loss to 0.2 achieves the lowest mean L1 error, mean L2 error and the highest PSNR value among four different parameter settings.

### C. Effectiveness of Fully Convolutional Architecture

This section investigates whether applying fully convolutional architecture benefits semantic inpainting. We design a CASI+$fc$ model by inserting two fully connected layers after the third layer of the CASI flatting block (described in Figure 2a). The first fully connected layer takes a convolutional feature map as input and outputs a 2048-d feature vector which is followed by a Tanh layer. The second fully connected layer takes the output of the activation layer as input and output a feature map with spatial dimensions. Then the fourth layer of the CASI flatting block takes the feature map as input. We compared CASI+$fc$ model and CASI model on Paris Street View dataset. As Table III shows, CASI outperforms CASI+$fc$ by 2.21% in mean $L1$ error, 0.34% in mean $L2$ error and 1.54dB with regards to PSNR although CASI+$fc$ contains more parameters than CASI. The result suggests applying fully convolutional architecture is more conducive for generative network as the fully connected layers could collapse the spatial structure of the image features.

### D. Effectiveness of Residual Block

This section verifies whether adding residual blocks enhance the performance. We design a CASI- model by removing the two residual blocks in CASI model and demonstrate comparison results between them. As shown in the upper part in Table IV, CASI outperforms CASI- by 0.2% in mean $L1$

error, 0.1% in mean $L2$ error and 0.5dB in PSNR, on the ImageNet-20 dataset. As the lower part in Table IV shows, CASI presents better performance than CASI- in mean $L1$ error, mean $L2$ error and PSNR value, on the Paris Street View dataset. The above results suggest that adding residual blocks improves prediction accuracy for the CASI model.

### E. High-resolution Case

This section investigates how our method performs on high-resolution cases. The motivation of investigation on high-resolution cases is that most existing neural network based inpainting methods can only deal with input images not larger than $128 \times 128$. This section demonstrates how the proposed method perform with input images of $512 \times 512$. Two groups of experiments are presented. The first group compare our method to [2] by scaling image to match with the input size of [2]. As shown in upper part of Table V, our CASI model presents lower mean $L1$ error, lower mean $L2$ error and higher PSNR value than ContextEncoder [2] in high-resolution Paris Street View dataset. The second group investigates whether adding a post-optimization based on our model deals with high-resolution cases. One concurrent work, NeuralPatch [7], trains its network to synthesize content at the image center and presents high-resolution object removal results during testing. We have integrated our method with post-optimization in [7] (denoted as CASI+) and demonstrate better performance than NeuralPatch [7]. As the lower part in Table V shows, the CASI+ method achieves lower mean $L1$ error, lower mean $L2$ error and higher PSNR value in comparison to NeuralPatch [7], which suggests that the proposed CASI can provide more accurate reference content for post-optimization based image completion methods. Figure 4 is a qualitative comparison between [7] and CASI+. As Figure 4 shows, CASI+ extends more reasonable edges and preserves more details than [7]. More comparison results can be found in the supplementary document.

### F. General and In-the-wild Case

This section investigates how the proposed method perform on general and in-the-wild cases. The first experiment in this
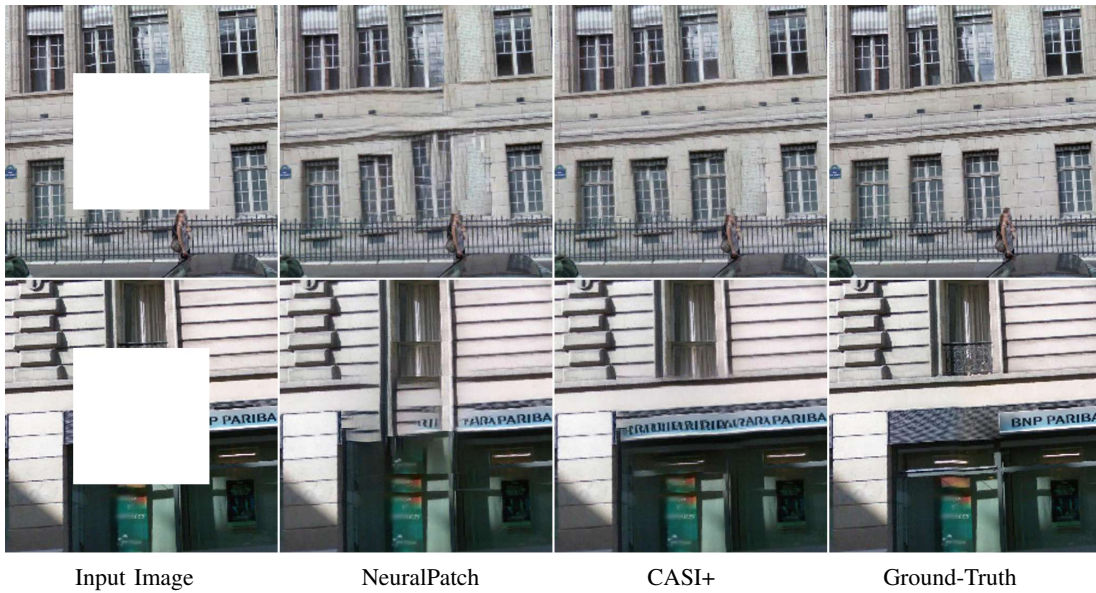
| Input Image | NeuralPatch | CASI+ | Ground-Truth |

Fig. 4: High-resolution Cases on Paris StreetView



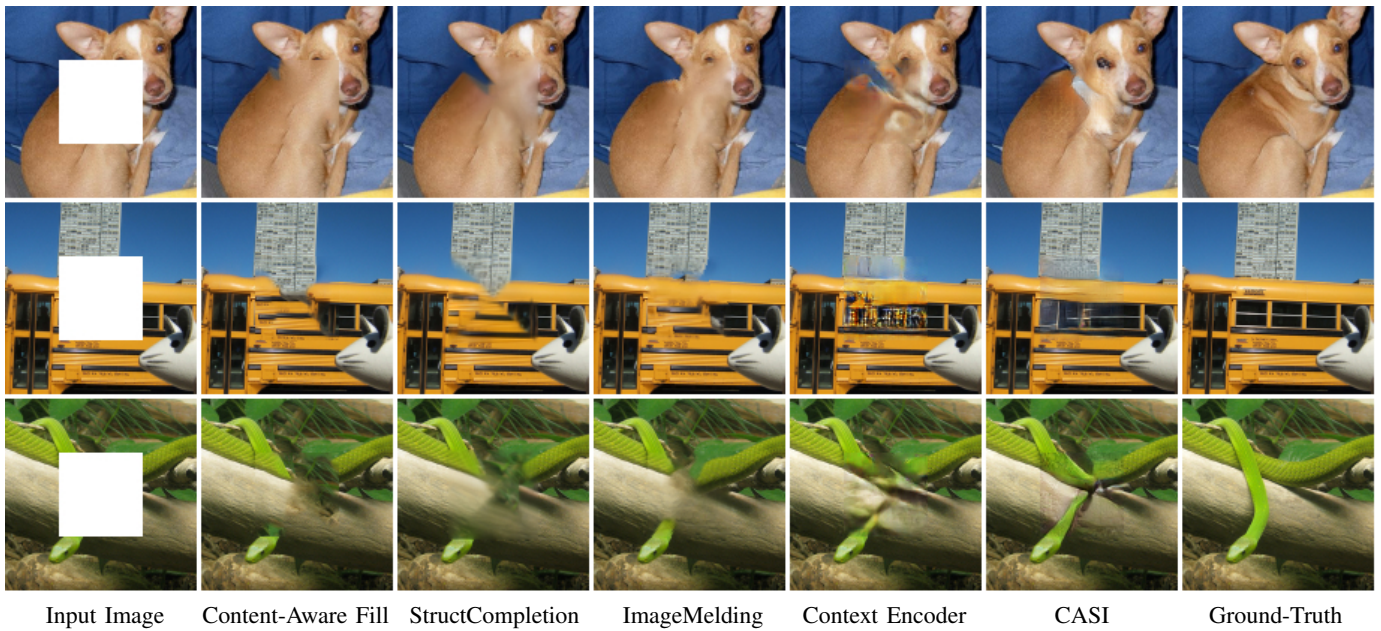| Input Image | Content-Aware Fill | StructCompletion | ImageMelding | Context Encoder | CASI | Ground-Truth |

Fig. 5: Comparison on ImageNet-20 dataset

section is to test the proposed method on high-resolution real images that are collected out of ImageNet and Paris StreetView dataset. The qualitative results of the first experiment are shown in Figure 6. The resolution of the input images in Figure 6 are $430 \times 645$, $708 \times 1062$ and $426 \times 570$. The results verify that our proposed method could perform well on in-the-wild cases.

The second experiment in this section is to test the proposed method on real images with irregular corrupted region. The qualitative results of the second experiment are displayed in Figure 7. These input images are also collected in-the-wild out of ImageNet and Paris StreetView datasets and their resolutions are $357 \times 500$, $332 \times 450$ and $332 \times 450$ respectively. The results suggest that the proposed algorithm is capable of

repairing images with irregular corrupted region.

### G. Investigation of Generalization Ability

This section investigates the generalization ability of the CASI model. If the CASI model has weak generalization ability and overfits the training data, it may predict what it memorize from the training data. Thus we conduct a nearest neighbor inpainting (NN-inpainting) experiment. For each testing input image, we search for the most matching patch from the training dataset to complete the image, using the algorithm proposed in [23]. The qualitative results of NN-inpainting are displayed in Figure 3. The CASI results (in Figure 3d) are quite different from the NN-inpainting results

Input Image                  CASI+                  Ground-Truth

Fig. 6: Inpainting Results of CASI with In-the-wild Cases

and demonstrate the superiority in preserving both appearance and structure coherence, which indicates that the CASI model does not simply copy or memorize patch from the training dataset while repairing the input images.

### H. Comparison with the State of the Art

We compare our proposed CASI model trained using the joint loss with other 4 state-of-the-art image inpainting methods, including Content-Aware Fill [39], StructCompletion [8], ImageMelding [20] and Context Encoder [2]. As shown in Figure 5, methods [39], [8], [20] without using neural network fail to recover the dog face in the first sample, extend the bus window in the second sample and connect the snake body in the third sample. These methods fail to recover the high-level semantics. Context Encoder struggles to display clear structure while the proposed CASI shows visually acceptable results in Figure 5.

The second experiment in this section compares our method with other state-of-the-art inpainting methods [1], [39], [7], [8], [20], [2] on the Paris StreetView dataset. Table VI shows the quantitative results. Results from PatchMatch [39], Neural Patch Synthesis (NeuralPatch) and Context Encoder are collected from [7], [7] and [2], respectively. As shown in Table VI, our results exceed others by a considerable margin under all three measures. Our method outperforms the second

TABLE VI: Quantitative results on Paris StreetView

| Method | mean $L1$ error | mean $L2$ error | PSNR |
|---|---|---|---|
| PatchMatch | 12.59% | 3.14% | 16.82dB |
| NeuralPatch | 10.01% | 2.21% | 18.00dB |
| StructCompletion | 9.67% | 2.07% | 18.03dB |
| ImageMelding | 9.55% | 2.19% | 18.05dB |
| Context Encoder | 9.37% | 1.96% | 18.58dB |
| CASI | **7.49%** | **1.37%** | **20.37dB** |

best by 1.58% in mean $L1$ error, 0.53% in mean $L2$ error and 1.56dB in PSNR.

### I. Investigation of Criteria for Inpainting

In this section, we use more criteria to evaluate CASI and Context Encoder, and propose two new criteria for semantic inpainting. There are three major experiments. In the first experiment, we evaluate inpainting methods using structural similarity index (SSIM) [40] and feature similarity index (FSIM) [41]. These indices are originally applied to image quality assessment (IQA) that attempts to quantify the visibility of differences between two images. Here we investigate the visual differences between inpainting results and their corresponding ground truth. Thus we test inpainting methods
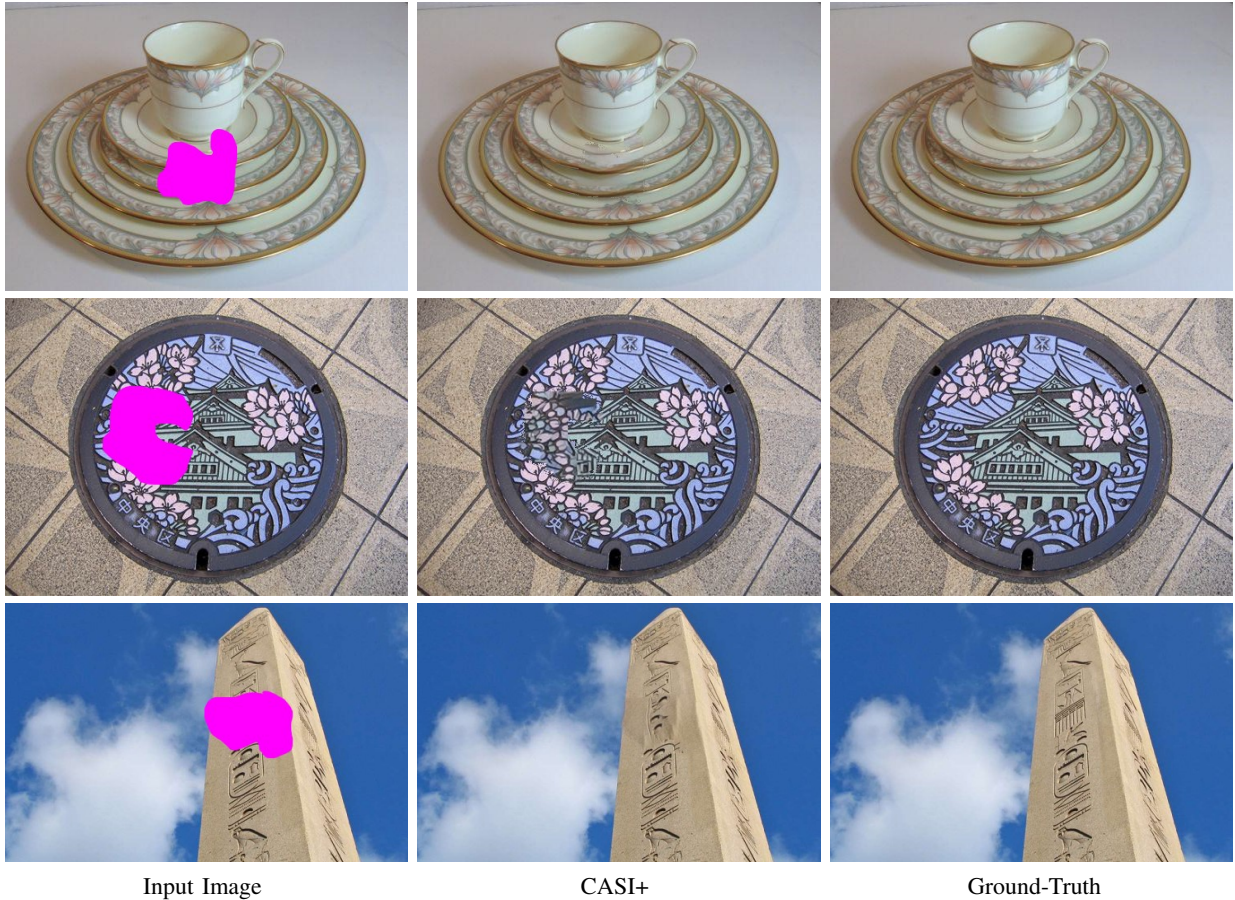
| Input Image | CASI+ | Ground-Truth |

Fig. 7: Examples of Inpainting Results on Irregular Corrupted Regions

TABLE VII: Similarity Indices on ImageNet-20.

| Method | SSIM | FSIM | FSIMc |
|---|---|---|---|
| Context Encoder | 0.2579 | 0.6977 | 0.6899 |
| CASI,$L2$ | 0.5196 | 0.6255 | 0.6202 |
| CASI,$L2 + per$ | 0.4927 | 0.6843 | 0.6779 |
| CASI,$L2 + adv$ | 0.5141 | 0.7202 | 0.7148 |
| CASI,$L2+adv+per$ | **0.5198** | **0.7239** | **0.7187** |
| $\lambda_{per} = 0$ | 0.5141 | 0.7202 | 0.7148 |
| $\lambda_{per} = 0.2$ | 0.5198 | 0.7239 | 0.7187 |
| $\lambda_{per} = 0.4$ | 0.5093 | 0.7203 | 0.7149 |
| $\lambda_{per} = 0.7$ | 0.4951 | 0.7163 | 0.7108 |

TABLE VIII: Local Entropy Errors on ImageNet-20.

| Method | LEMSE | LEMAE |
|---|---|---|
| Context Encoder | 0.5872 | 0.5391 |
| CASI,$L2$ | 1.8926 | 1.0795 |
| CASI,$L2 + per$ | 0.8454 | 0.7219 |
| CASI,$L2 + adv$ | 0.4869 | 0.4945 |
| CASI,$L2+adv+per$ | **0.4611** | **0.4847** |
| $\lambda_{per} = 0$ | 0.4869 | 0.4945 |
| $\lambda_{per} = 0.2$ | 0.4611 | 0.4847 |
| $\lambda_{per} = 0.4$ | 0.4470 | 0.4759 |
| $\lambda_{per} = 0.7$ | 0.4492 | 0.4771 |

using the two IQA indices. SSIM is a classical index defined by structural similarity while FSIM is the state of the art based on two low-level features, phase congruency (PC) and gradient magnitude. FSIM is defined in Eq. (7),

$$FSIM = \frac{\sum S_{PC}(x) \cdot S_G(x) \cdot PC_m(x)}{\sum PC_m(x)}, \qquad (7)$$

where $S_{PC}(x)$ and $S_G(x)$ are PC similarity and gradient similarity respectively at position $x$, and $PC_m(x)$ is the PC value of $x$ as a weight. As shown in Table VII, all CASI models achieve higher similarity with the ground truth than Context Encoder under SSIM, FSIM and FSIMc (FSIM for color image). It indicates that our method not only recovers more consistent structures but also synthesizes content with higher visual quality. However, SSIM and FSIM are still biased towards blurry results of CASI, $L2$ $(+l_{per})$.

In the second experiment, we introduce a novel local entropy error to rate blurry predictions more accurately. Entropy in texture analysis is a statistic characterizing the texture within an image region, as defined in [42]. The local entropy at a pixel is defined as the entropy within a $9 \times 9$ neighborhood of the pixel. We define local entropy error as the mean squared error (denoted as LEMSE) or the mean absolute error (LEMAE) of local entropy within the synthesized region. As shown in Table VIII, our proposed CASI delivers the lowest LEMSE and LEMAE among all methods. In addition, CASI

TABLE IX: Semantic Errors on ImageNet-20

| Method | SME-r50 | SME-r101 | SME-r152 | SME-r200 | SME-v16 | SME-v19 |
|---|---|---|---|---|---|---|
| baseline | 0.2063 | 0.1735 | 0.1852 | 0.2063 | 0.1794 | 0.2086 |
| Context Encoder | 0.1467 | 0.1462 | 0.1442 | 0.1467 | 0.1001 | 0.1123 |
| CASI,$L2$ | 0.1862 | 0.1908 | 0.1886 | 0.1877 | 0.1444 | 0.1652 |
| CASI,$L2 + per$ | 0.1542 | 0.1631 | 0.1671 | 0.1626 | 0.1213 | 0.1384 |
| CASI,$L2 + adv$ | 0.1276 | 0.1359 | 0.1349 | 0.1362 | 0.0846 | 0.0952 |
| CASI,$L2+adv+per$ | **0.1070** | **0.1180** | **0.1201** | **0.1200** | **0.0721** | **0.0775** |
| $\lambda_{per} = 0$ | 0.1276 | 0.1360 | 0.1350 | 0.1363 | 0.0846 | 0.0952 |
| $\lambda_{per} = 0.2$ | 0.1070 | 0.1180 | 0.1201 | 0.1200 | 0.0721 | 0.0775 |
| $\lambda_{per} = 0.4$ | 0.1074 | 0.1125 | 0.1218 | 0.1215 | 0.0704 | 0.0767 |
| $\lambda_{per} = 0.7$ | 0.0994 | 0.1126 | 0.1117 | 0.1131 | 0.0632 | 0.0702 |

with $L2$ loss and CASI with $L2 + per$ loss achieve the largest and second largest errors under both LEMSE and LEMAE, which is consistent with most of the visual results (a subset is given in Figure 3) and confirms that our proposed local entropy error is capable of rating over-smooth results accurately.

In the third experiment, we propose a high-level criterion, semantic error, which aims at measuring how successful an inpainting method recovers the semantics. Semantic error (SME) is defined with respect to a pre-trained image classifier that outputs a probability of the image being part of each possible category. SME is based on two probabilities that the groundtruth image and the synthesized image belong to the groundtruth category respectively. It is formulated as in the following equation,

$$SME = \frac{1}{n} \sum_{i=1}^{n} max(0, P_{x_i}^{y_i} - P_{z_i}^{y_i}), \qquad (8)$$

where $n$ is the number of testing samples, $x_i$, $z_i$ and $y_i$ are the groundtruth image, synthesized image (with real context) and the groundtruth category of the $i$-th sample. $P_{x_i}^{y_i}$ is the probability that image $x_i$ belongs to category $y_i$, estimated by a pre-trained classifier (e.g., residual network [43] or VGG network[34]). Here we associate the probability of assigning the correct label with our semantic error because we focus on to what extent a corruption "makes a dog unlike a dog" and to what extent the restored content "makes a dog look like a dog again". A baseline model simply fills the missing region with the mean pixel value. The SME of this baseline measures how much a corrupted region harms the semantic information of an image. In Table IX, SME-r$L$ represents the SME achieved by applying an $L$-layer residual network as the classifier while SME-v$L$ represents the SME achieved by adopting an $L$-layer VGG network as the classifier. Notice that our feature network is simpler than the ResNets used for estimating SME, which implies that harvesting knowledge using a low-capacity model can reduce the SME estimated by a high-capacity classifier. As shown in Table IX shows, our proposed network outperforms other inpainting methods by achieving the smallest semantic error.

Perceptual loss weight is also investigated on the above new criteria for semantic inpainting, as shown in the lower part of Table VII, Table VIII and Table IX. $\lambda_{per} = 0.7$ performs better on similarity indices and semantic errors while $\lambda_{per} = 0.4$ demonstrates better results on local entropy error. To compromise different criteria, $\lambda$ is chosen from 0.2 to 0.4.

## VI. CONCLUSION

In this paper, we have presented a fully convolutional generative adversarial network with a context-aware loss function for semantic inpainting. This network employs a fully convolutional architecture in the generator, which does not have a fully connected layer as the bottleneck layer. The joint loss includes a perceptual loss to capture semantic information around the synthesized region. In addition, we have developed two new measures for evaluating sharpness and semantic validity respectively. In summary, our method delivers state-of-the-art results in qualitative comparisons and under a wide range of quantitative criteria.

## REFERENCES

[1] https://research.adobe.com/project/content-aware-fill.

[2] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544.

[3] J. Xie, L. Xu, and E. Chen, "Image denoising and inpainting with deep neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 341–349. [Online]. Available: http://papers.nips.cc/paper/4686-image-denoising-and-inpainting-with-deep-neural-networks.pdf

[4] A. Fawzi, H. Samulowitz, D. Turaga, and P. Frossard, "Image inpainting through neural networks hallucinations," in *2016 IEEE 12th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, July 2016, pp. 1–5.

[5] 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*, 2014, pp. 2672–2680.

[6] R. Yeh, C. Chen, T. Y. Lim, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with perceptual and contextual losses," *arXiv preprint arXiv:1607.07539*, 2016.

[7] C. Yang, X. Lu, Z. Lin, E. Shechtman, O. Wang, and H. Li, "High-resolution image inpainting using multi-scale neural patch synthesis," *arXiv preprint arXiv:1611.09969*, 2016.

[8] J.-B. Huang, S. B. Kang, N. Ahuja, and J. Kopf, "Image completion using planar structure guidance," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 4, p. 129, 2014.

[9] J. Jia and C.-K. Tang, "Image repairing: Robust image synthesis by adaptive nd tensor voting," in *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 1. IEEE, 2003, pp. I–I.

[10] A. Rares, M. J. Reinders, and J. Biemond, "Edge-based image restoration," *IEEE Transactions on Image Processing*, vol. 14, no. 10, pp. 1454–1468, 2005.

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

[12] Y. Pritch, E. Kav-Venaki, and S. Peleg, "Shift-map image editing," in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 151–158.

[13] Y. Wexler, E. Shechtman, and M. Irani, "Space-time video completion," in *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, vol. 1. IEEE, 2004, pp. I–I.

[14] E. I. M. Wexler, Yonatan Shechtman, "Space-time completion of video," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 3, no. 29, pp. 463–476, 2007.

[15] J. Sun, L. Yuan, J. Jia, and H.-Y. Shum, "Image completion with structure propagation," *ACM Transactions on Graphics (ToG)*, vol. 24, no. 3, pp. 861–868, 2005.

[16] A. A. Efros and T. K. Leung, "Texture synthesis by non-parametric sampling," in *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, vol. 2. IEEE, 1999, pp. 1033–1038.

[17] 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*. ACM, 2001, pp. 341–346.

[18] A. Criminisi, P. Perez, and K. Toyama, "Object removal by exemplar-based inpainting," in *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 2. IEEE, 2003, pp. II–II.

[19] Z. Xu and J. Sun, "Image inpainting by patch propagation using patch sparsity," *IEEE transactions on image processing*, vol. 19, no. 5, pp. 1153–1165, 2010.

[20] S. Darabi, E. Shechtman, C. Barnes, D. B. Goldman, and P. Sen, "Image melding: Combining inconsistent images using patch-based synthesis." *ACM Trans. Graph.*, vol. 31, no. 4, pp. 82–1, 2012.

[21] I. Drori, D. Cohen-Or, and H. Yeshurun, "Fragment-based image completion," in *ACM Transactions on graphics (TOG)*, vol. 22, no. 3. ACM, 2003, pp. 303–312.

[22] N. Komodakis and G. Tziritas, "Image completion using efficient belief propagation via priority scheduling and dynamic pruning," *IEEE Transactions on Image Processing*, vol. 16, no. 11, pp. 2649–2661, 2007.

[23] J. Hays and A. A. Efros, "Scene completion using millions of photographs," *ACM Transactions on Graphics (SIGGRAPH 2007)*, vol. 26, no. 3, 2007.

[24] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[25] A. Dosovitskiy and T. Brox, "Generating images with perceptual similarity metrics based on deep networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 658–666.

[26] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *arxiv*, 2016.

[27] C. Li and M. Wand, "Precomputed real-time texture synthesis with markovian generative adversarial networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 702–716.

[28] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

[29] S. Xie and Z. Tu, "Holistically-nested edge detection," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.

[30] G. Li and Y. Yu, "Deep contrast learning for salient object detection," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 478–487.

[31] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711.

[32] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," *arXiv preprint arXiv:1508.06576*, 2015.

[33] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky, "Texture networks: Feed-forward synthesis of textures and stylized images," in *Int. Conf. on Machine Learning (ICML)*, 2016.

[34] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[36] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.

[37] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[38] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. A. Efros, "What makes paris look like paris?" *ACM Transactions on Graphics (SIGGRAPH)*, vol. 31, no. 4, pp. 101:1–101:9.

[39] C. Barnes, D. B. Goldman, E. Shechtman, and A. Finkelstein, "The patchmatch randomized matching algorithm for image manipulation," *Communications of the ACM*, vol. 54, no. 11, pp. 103–110, 2011.

[40] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[41] L. Zhang, L. Zhang, X. Mou, and D. Zhang, "Fsim: A feature similarity index for image quality assessment," *IEEE transactions on Image Processing*, vol. 20, no. 8, pp. 2378–2386, 2011.

[42] R. C. Gonzalez, R. E. Woods, and S. L. Eddins, *Digital Image Processing Using MATLAB®*. New Jersey, Prentice Hall, 2003, ch. 11.

[43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.