# Pixel-wise Dense Detector for Image Inpainting

Ruisong Zhang[1,2] iD, Weize Quan[1,2] iD, Baoyuan Wu[3,4] iD, Zhifeng Li[5] and Dong-Ming Yan[1,2†] iD

[1]National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China
[2]School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China
[3]School of Data Science, the Chinese University of Hong Kong, Shenzhen, China
[4]Secure Computing Lab of Big Data, Shenzhen Research Institute of Big Data, China
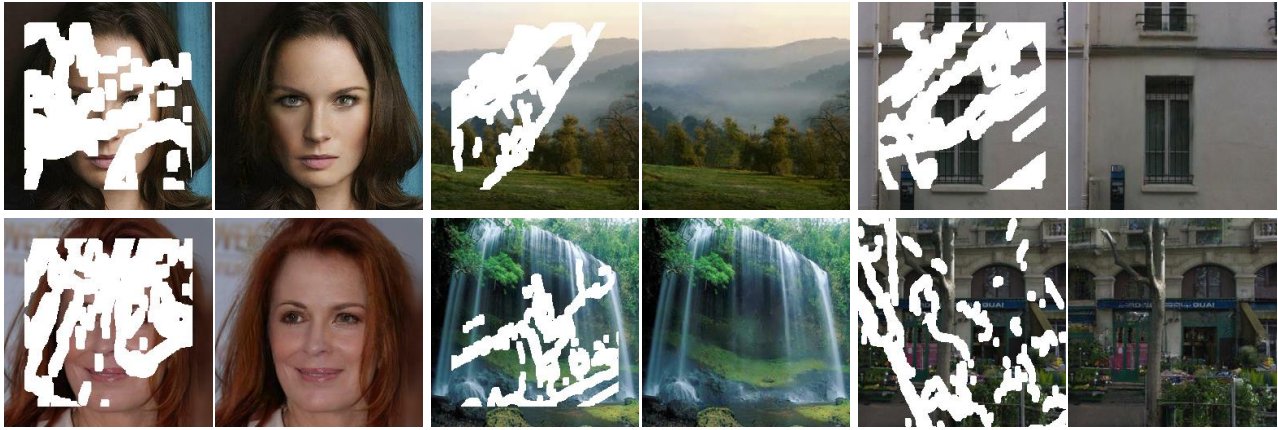[5]Tencent AI Lab, Shenzhen, China

**Figure 1:** *Image inpainting results by our proposed method.*

## Abstract

*Recent GAN-based image inpainting approaches adopt an average strategy to discriminate the generated image and output a scalar, which inevitably lose the position information of visual artifacts. Moreover, the adversarial loss and reconstruction loss (e.g., $\ell_1$ loss) are combined with tradeoff weights, which are also difficult to tune. In this paper, we propose a novel detection-based generative framework for image inpainting, which adopts the min-max strategy in an adversarial process. The generator follows an encoder-decoder architecture to fill the missing regions, and the detector using weakly supervised learning localizes the position of artifacts in a pixel-wise manner. Such position information makes the generator pay attention to artifacts and further enhance them. More importantly, we explicitly insert the output of the detector into the reconstruction loss with a weighting criterion, which balances the weight of the adversarial loss and reconstruction loss automatically rather than manual operation. Experiments on multiple public datasets show the superior performance of the proposed framework. The source code is available at https://github.com/Evergrow/GDN_Inpainting.*

**CCS Concepts**
• *Computing methodologies* → *Image processing;*

## 1. Introduction

Image inpainting is a technique of filling the semantically correct and visually plausible contents in the missing regions of corrupted images, as shown in Fig. 1. It has various applications, such as repairing deteriorated photographs, removing undesired objects from images, and even editing specified contents of images. Different from general generative tasks, image inpainting deals with corrupted images with plenty of contextual background, which is not

---

† Corresponding Author

only prior information to assist in reconstruction but also is a constraint to limit generated contents.

Early works attempt to fill missing regions with some optimization algorithms, *e.g.*, propagating the isophotes from boundaries [BBC*01; BSCB00; EF01] or copying the matching information from background patches into missing regions [BSFG09; DSB*12; HKAK14]. These methods with low-level features achieve good results especially inpainting background or some repetitive patterns. However, as they cannot extract high-level semantic information, they often fail to generate reasonable structures with novel patterns in real-world scenarios. Moreover, high computational cost also limits their deployment in practical applications.

With the advance of deep learning, *i.e.*, convolutional neural networks (CNN) and generative adversarial networks (GAN) [GPM*14], recent works [PKD*16; YLL*17; ISI17; LRS*18; YLL*18; ZFCG19] model image inpainting as a conditional generation problem to learn the mapping between the corrupted input images and the ground-truth images. These deep inpainting techniques extract rich semantic information from large scale training data, based on contextual background, to fill missing regions with plausible contents and textures. Joint adversarial training also improves the visually realistic effect of image recovery. Unfortunately, these methods often suffer from distorted structures, blurry artifacts and distinct incoherence with surrounding areas, especially for complex scenes.

To address above problems, some deep inpainting methods [YLY*18; YLY*19; NNJ*19; RYZ*19] utilize two-stage networks that rough out the missing structures or contents in the first stage, and then recover refined textures using coarse information in the second stage. Contextual attention [YLY*18] inserts the attention module into the second stage to encourage spatial coherency of attention. Based on the structure image from the first stage, StructureFlow [RYZ*19] generates the textures to fill missing regions, which shows reasonable structures and vivid details. Compared with single encoder-decoder networks, however, two-stage networks are much deeper and thus cause extra computational cost and inference time.

In addition, most deep inpainting methods [PKD*16; ISI17; YLY*18] follow an adversarial framework, which not only creates realistic results but also provides an agonizing multi-objective optimization. The first objective is the coherence and similarity with ground-truth images formulated as the reconstruction loss, and another objective is the perceptually realistic results programmed as the adversarial loss. Several tradeoff parameters are empirically set up to optimize these two objectives simultaneously. To our best knowledge, there are seldom works that improve this strategy with explicable theories.

To address the limitation discussed above, we propose a novel inpainting framework, comprising a generative model and a detective model, with a unique objective function to optimize. Generative network is an encoder-decoder architecture with residual blocks [HZRS16] for repairing the corrupted images. Inspired by the image segmentation task [LSD15], we design and implement a fully convolutional network as detective network to evaluate the inpainting results in a pixel-wise manner. The binary mask is the target of the detective network using an approach of weakly supervised learning to capture visual artifacts of entire generated image. Compared with a single scalar from the standard discriminator, the location information of artifacts with up-sampling builds a dense mapping function between the output image and ground-truth image for more accurate valuation. Under the guidance of accurate valuation, inpainting techniques pay more attention to artifacts such as distortion, blurriness and incoherence, and reduce them, thus making the valuation from the detector trend inaccurate. The adversarial learning of the generator and detector moves from "whether" to "where". Moreover, we propose a reconstruction loss weighted the valuation as inpainting objective function to solve the multi-objective optimization. This loss function without the hyper-parameter is better for describing the image inpainting task. The optimal balance between vraisemblance and similarity is taken by networks rather than empiricism. Our contributions are summarized as follows:

- We propose a novel framework that merges the generator and the detector, and both experimental results and explicable theory show that this framework contributes to reduce artifacts.
- The detector using weakly supervised learning makes a dense estimate of the entire inpainted image, which is similar to human perceptual evaluation.
- We design a weighted reconstruction loss to optimize the combined inpainting objectives including vraisemblance and similarity automatically, which achieves superior results.

## 2. Related Work

Existing image inpainting approaches can be classified into two categories: low-level features based approaches and deep semantic features based approaches. The former usually involves some geometric techniques for texture synthesis or structure propagation in low-level features. The latter often solves the inpainting problem by deep neural networks to extract global semantic features.

### 2.1. Methods Based on Low-level Features

Approaches based on low-level features are roughly divided into diffusion-based methods and patch-based methods. Traditional diffusion-based methods [BBC*01; BSCB00; EF01] typically utilize variational algorithms to propagate neighboring appearance information (*e.g.*, the isophotes) into the missing regions. Due to the limited extended prediction of partial differential equation, these methods could not produce good results when the missed regions are broad. The restoration regions generated by this kind of methods also lack meaningful structure information.

Unlike diffusion-based methods just focusing on the surrounding pixels of missing regions, patch-based methods [DSB*12; HKAK14; DAFC19] measure the similarity between missing regions and each patch from the whole context, and recover target region by copying the matched patch. Bidirectional similarity measure [SCSI08] is proposed to model visual data for summarization and reduce visual artifacts. However, dense computation of patch similarity often comes at an expensive computation cost. To accelerate these searching algorithms, PatchMatch [BSFG09] captures patch matches via random sampling and natural coherence

in the imagery, and it is widely used in the interactive editing tools. Low-level features based approaches lacking deep understanding of whole image just generate repetitive content from background without unique filling information.

## 2.2. Methods Based on Deep Semantic Features

Deep semantic features based approaches attempt to perceive the semantic structure of the corrupted image by deep neural networks for better restoration results. Context Encoders [PKD*16] first introduce CNNs for inpainting missing regions. The proposed encoder-decoder architecture is trained via incorporated reconstruction loss and adversarial loss [GPM*14]. However, this network excessively concerns about entire consistency and it often results in visual artifacts in detailed regions. To generate high-frequency details, Yang *et al.* [YLL*17] propose a multi-scale neural patch synthesis based on joint optimization of image content and texture constraints, and Iizuka *et al.* [ISI17] unite global and local discriminator to assess completed image from generative network with dilated convolutions [YK15]. However, local discriminator fails to deal with irregular missing regions.

Recently, deep inpainting techniques show the multiplex development. Attention is an important mechanism for image inpainting to build long-term correlations between missing regions and distant contextual information [YLY*18; LJXY19; ZFCG19; YLY*19]. Yu *et al.* [YLY*18] design a coarse-to-fine network and first introduce contextual attention into refined network. However, effect of attention mechanism mainly depends on results of coarse network, and poor coarse reconstruction often causes wrong match. To avoid the interference of corrupted regions, some works modify conventional convolutional operation, such as partial convolution [LRS*18] and gated convolution [YLY*19], calculating convolution only on valid pixels. These variants succeed on reduction of blurry artifacts. EdgeConnect [NNJ*19] and Structure-Flow [RYZ*19] generate reasonable structures with prior information and then synthesize fine texture. This easy-to-difficult process can obtain satisfactory visual effects, but it also needs sophisticated preprocessing to extract edge maps [NNJ*19] or edge-preserved smooth images [RYZ*19]. In addition, several works [LRS*18; NNJ*19; XLL*19] project images into high-dimensional features space built by pretrained VGG-16 [SZ14] on ImageNet [DDS*09] and then measure the similarity by perceptual loss [JAF16] and style loss [GEB16] to improve inpainting results. A possible limitation is that it reduces filling generalization of the scene outside ImageNet [DDS*09].

Deep inpainting approaches reviewed above mostly follow the adversarial framework used in Context Encoders [PKD*16]. In this framework, the discriminator takes the inpainted image as the input to evaluate in level of whole image or its patches (*e.g.*, Patch-GAN [IZZE17]) , and discards the meaningful location information of blurry artifacts in the adversarial loss when training the generator. Specially, in PatchGAN, the patch-level valuations constitute a tensor, which is unseen for the generator, while the generator can only access the adversarial loss after average. On the contrary, we introduce a detective network to detect artifacts pixel by pixel, and the output of the detector assists the generator in eliminating color discrepancy and blurriness via an adaptive weighting strategy.

## 3. Proposed Method

Our proposed method utilizes a generative network to reconstruct corrupted images and a detective network to evaluate outputs of the generator to perform image inpainting, as shown in Fig. 2. In this section, we first review GAN-based image inpainting. Then, we present the proposed novel detection-based generative framework, *i.e.*, the coupling of the generator and the detector in training stage. At last, the details of network architecture and the loss function of our method are explained.

## 3.1. GAN-based Image Inpainting

Generative Adversarial Network [GPM*14] is an advanced generative framework including two nets: generator $G$ and discriminator $D$. The contest training strategy drives both networks to improve their performance until reaching to a global optimum. Mathematically, this competitive process between $G$ and $D$ can be described as a min-max optimization problem with value function $V(G,D)$:

$$\min_{G} \max_{D} V(G,D) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] \\ + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))], \quad (1)$$

where $\mathbf{x}$ is real data, and $\mathbf{z}$ is input noise. The distribution is omitted for simplicity in following formulas. In Eq. (1), the discriminator $D$ tells the difference between real data distribution $\mathbf{x}$ and generated data distribution $G(\mathbf{z})$ as much as possible. At the same time, the generator $G$ tries to mimic real data distribution and fool the discriminator $D$.

Image inpainting is a special generative problem with plenty of prior information (*i.e.*, corrupted image) as input rather than random noise. To improve the reconstruction quality of corrupted regions, and make whole image more realistic and vivid, most of deep image inpainting methods follow the GAN-based framework. The corresponding optimization for image inpainting is written as:

$$\min_{G} \max_{D} V(G,D) = \mathbb{E}[\log D(\mathrm{I}_{gt})] + \mathbb{E}[\log(1 - D(G(\mathrm{I}_{in}, \mathrm{M})))], \quad (2)$$

where $\mathrm{I}_{gt}$ is the ground-truth image, $\mathrm{I}_{in}$ is the input corrupted image, and $\mathrm{M}$ is the binary mask (1: the missing region and 0: valid regions). Usually, $\mathrm{I}_{in} = \mathrm{I}_{gt} \odot (1 - \mathrm{M}) + \mathrm{M}$, where $\odot$ denotes pixel-wise product.

For the generator $G$, the final objective function is the combination of the adversarial loss $\mathcal{L}_{adv}$ and the reconstruction loss $\mathcal{L}_{\ell_1}$, which measures the coherence and similarity between predicted image and ground-truth image. Corresponding formulation can be expressed as

$$\mathcal{L}_G = \lambda_{adv} \cdot \mathcal{L}_{adv} + \lambda_{\ell_1} \cdot \mathcal{L}_{\ell_1} \quad (3) \\ = \lambda_{adv} \cdot \mathbb{E}[\log(1 - D(G(\mathrm{I}_{in}, \mathrm{M})))] + \lambda_{\ell_1} \cdot ||\mathrm{I}_{out} - \mathrm{I}_{gt}||_1,$$

where $\mathrm{I}_{out}$ is the prediction of generator $G$, $\lambda_{adv}$ and $\lambda_{\ell_1}$ are the tradeoff parameters setting empirically.

## 3.2. Detection-based Generative Framework

As discussed above, the discriminator $D$ in GAN actually is a classifier, which just outputs a single scalar (label or probability). This scalar may not properly evaluate the quality of the generated image
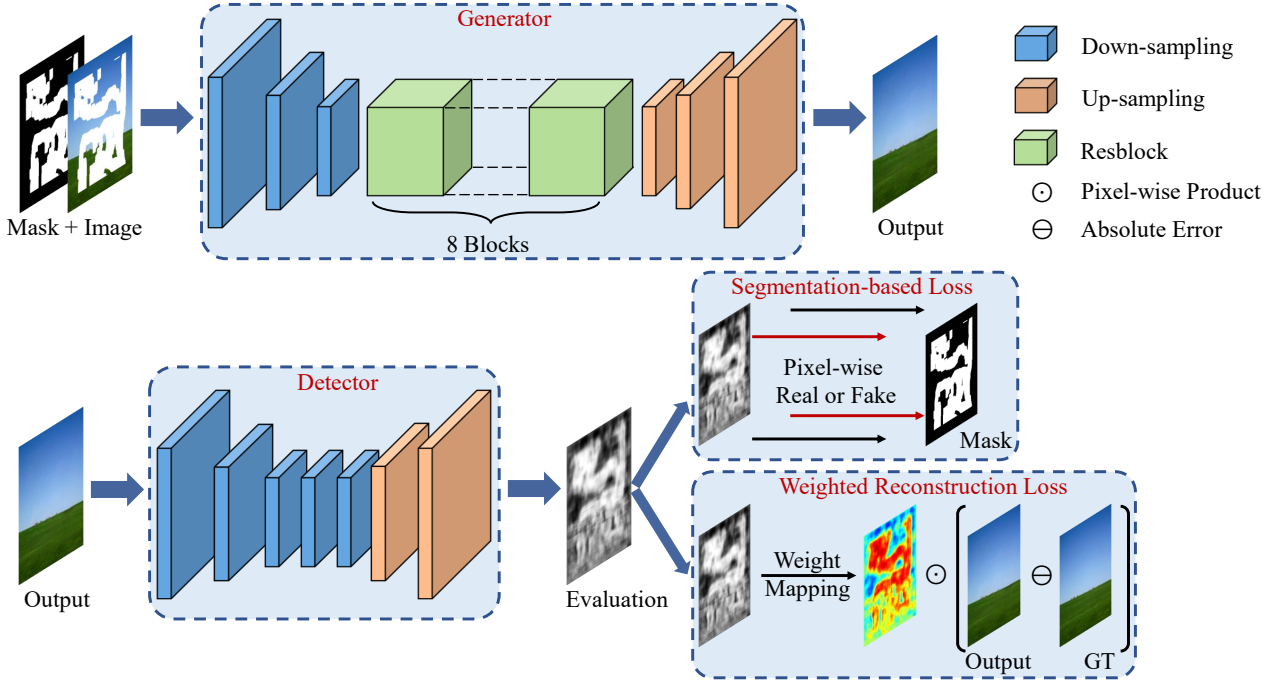
**Figure 2:** *Illustration of our proposed detection-based generative framework.*

for image inpainting, where non-masked regions with much information easily has better reconstructed quality than the missing regions. The common average criterion in $D$ to some extent weakens this difference, more precisely, loses the exact position information of "authentic" artifacts. In extreme cases, such scalar might misguide the generator $G$ that pays less attention to the missing region.

To solve the above problems, we proposed a novel detection-based generative framework for image inpainting. This framework consists of a generator $G$ and a detector $D_{et}$. The generator $G$ repairs the missing region to be harmony with contextual background. The detector $D_{et}$ evaluates the output $I_{out}$ of $G$ in pixel-wise manner and also localizes the unreasonable completion region, *e.g.*, various artifacts, blurry patch, etc. Compared with a single scalar of $D$ in GAN, the precise localization significantly assists the generator $G$ in reconstructing of corrupted images.

Unfortunately, it is difficult to provide the ground-truth mask for $D_{et}$. In this work, we adopt a weakly supervised strategy: we consider the binary mask M as the "proxy" of the ground-truth mask. This is reasonable because the missing regions are usually difficult to predict than the non-masked regions (just reproducing the non-masked regions of the input image) and also have more artifacts. In addition, the similar pattern of corrupted regions in the non-masked regions are also captured by $D_{et}$ as the training progresses. The processing of this weakly supervised learning can be written as V $= D_{et}(I_{out})$, where V is the output of detector $D_{et}$ to evaluate the prediction $I_{out}$. The valuation V is the same size as the mask M, and its value, from 0 to 1, reflects the realistic degree of each pixel, *i.e.*, the lower value indicates the more realistic completion result. Unlike the standard discriminator, the detector $D_{et}$ gives

a human-like evaluation of the inpainted image, rather than a rough score. More analysis about the output of the detector $D_{et}$ will be shown in Section 4.5.

The coupled training should be executed between the generator $G$ and the detector $D_{et}$. As the coupled intermediate, $I_{out}$ is the input of the detector $D_{et}$, while valuation V, inserted in the reconstruction loss (see Eq. (3)), constitutes a weighted reconstruction loss, which optimizes the generator $G$. A segmentation loss considering the imbalance between masked and non-masked regions is used to optimize the detector $D_{et}$. The above description is a typically adversarial process where the detector accurately locates artifacts, *i.e.*, minimization of the segmentation-based loss, while the generator tries to deceive the detector that the probability of the artifacts in any location is the same, *i.e.*, maximization of the segmentation-based loss. We will validate the proposed framework in Section 4.3, and discuss the details of the loss functions and their implementation in Section 3.4.

### 3.3. Network Architecture

The proposed image inpainting framework consists of two networks: the generator $G$ and the detector $D_{et}$.

#### 3.3.1. Generator

The generator following an architecture similar to EdgeConnect [NNJ*19] comprises three components: two down-sampling layers, eight residual blocks [HZRS16], and two up-sampling layers, as shown in Fig. 2. A convolutional layer with kernel size = 4 and stride = 2 executes down-sample, and up-sample is conducted by a deconvolutional layer with kernel size = 4 and stride

= 2. Each of residual block holds two dilated convolutions [YK15] with kernel size = 3 , stride = 1 and dilation factor = 2. After all convolutions/deconvolutions except the last convolution in whole network, instance normalization [UVL17] and ReLU is followed, successively.

### 3.3.2. Detector

The detector is a seven-layer fully convolutional network, which is augmented by in-network up-sampling and pixel-wise loss for dense evaluation of the inpainted image. The first five convolutional layers down-sample images twice, followed by two deconvolutional layers to up-sample images back to the original size. The final softmax layer transforms the output to the probability map V. Leaky ReLU with $\alpha = 0.2$ is used in down-sampling stage, and all convolutional kernel size is 4.

### 3.4. Loss Functions

During the adversarial training, the generator targets both pixel-wise reconstruction precision and plausible visual result, while the detector aims at fine-grained evaluation for the prediction. We incorporate segmentation-based loss and weighted reconstruction loss to train our detector and generator, separately.

### 3.4.1. Segmentation-based Loss

As described in Section 3.2, we use the weakly supervised learning to train our detector with the binary mask as ground-truth, therefore, a natural option is the standard pixel-wise binary cross entropy loss,

$$\mathcal{L}_{CE} = -\frac{1}{N}\sum_{i=1}^{N} M_i \log V_i + (1-M_i)\log(1-V_i), \qquad (4)$$

where $M_i$ and $V_i$ respectively are the one element of mask M and valuation V, and $N$ is the number of elements in M. In Eq. (4), a smaller value of the loss function approximates more accurate valuation from the detector due to the weak supervision. Generally, the missing region in input image is smaller than the valid region, and this causes the imbalance between positive and negative samples. To balance the two classes, we introduce a weight $\alpha$, which is the mask ratio for input image. The balanced version of $\mathcal{L}_{CE}$ is formulated as

$$\mathcal{L}_{BCE} = -\frac{1}{N}\sum_{i=1}^{N} (1-\alpha)M_i \log V_i + \alpha(1-M_i)\log(1-V_i). \quad (5)$$

In addition, we also consider a recent segmentation loss, *i.e.*, focal loss [LGG*17], which is an enhanced version of Eq. (5) with tunable focusing parameter $\gamma \geq 0$,

$$\begin{aligned} \mathcal{L}_{Focal} = &-\frac{1}{N}\sum_{i=1}^{N} (1-\alpha)(1-V_i)^{\gamma}M_i \log V_i \\ &+ \alpha V_i^{\gamma}(1-M_i)\log(1-V_i). \end{aligned} \qquad (6)$$

In our experiments, we compare these three optional loss functions for training the detector, and find that the focal loss can yield the best performance, as shown in Section 4.4.

### 3.4.2. Weighted Reconstruction Loss

The generator learns with two objectives: the realistic visual quality and the consistency with the ground-truth image. The traditional adversarial framework combines these objectives with two hyper-parameters shown in Eq. (3), but it is still multi-objective optimization essentially. Compared with the single-objective optimization, multi-objective optimization is often difficult, *e.g.*, the balance between the maximum margin and the minimum error is a tough problem in Soft-Margin SVM algorithm [CWYZ04]. More importantly, tradeoff parameters without the physical significance to explain reduces generalization in different inpainting cases, *i.e.*, the dataset of face, scene and street view needs corresponding parameters. To address this problem, we refer to Boosting algorithm [DSS93], the idea is to increase the weight of weak samples and decrease the weight of strong samples. For a single image inpainting, the valuation V distinguishes the weak or strong pixels appropriately. Therefore, our proposed weighted reconstruction loss merging two objectives can be written as

$$\mathcal{L}_w = \frac{1}{N}\sum_{i=1}^{N} W_i \cdot ||I_{out}^i - I_{gt}^i||_1, \qquad (7)$$

where $W_i$, $I_{out}^i$, and $I_{gt}^i$ are the pixel-wise weight W, prediction $I_{out}$ and ground-truth $I_{gt}$, respectively. Note that the weight W maps to the valuation V, called weight mapping (please see following Section 3.4.3 for details). Minimization of the weighted reconstruction loss for the generator is to multiply the artifact by a smaller weight, which is just the opposite target of the detector, described in Eq. (4). Competition between the generator and detector drives to improve their capabilities until the artifacts are hard to be perceived by naked eyes or the detector, which is the fundamental purpose of image inpainting task.

### 3.4.3. Weight Mapping

For high accuracy of the missing region reconstruction, we adopt to enhance the weight of weak pixels instead of subtracting the weight of strong pixels, so the range of the weight W is $[1,+\infty)$. The weight mapping is a transition function from the valuation $[0,1]$ to the weight $[1,+\infty)$, and the candidate is linear or exponential functions.

- Linear transition can be written as $W = 1 + V$, and the range of the weight W is $[1,2]$. Although the simple form is easy to implement, the low upper bound value causes less enhancement for awful reconstruction.
- Exponential transition can be written as $W = x^V$, where $x$ $(x > 1)$ is a base number of exponential function, and the range of the weight W is $[1,x]$. Compared with linear transition, the exponential transition reduces the relative weight on well-inpainted pixels, putting more focus on artifacts. Unlike the tradeoff parameters $\lambda_{adv}$ and $\lambda_{\ell_1}$, $x$ controls the refined texture completion in the missing region purposefully, as shown in Fig. 3. The larger value of the base number indicates the clearer and more exquisite edge or texture information.

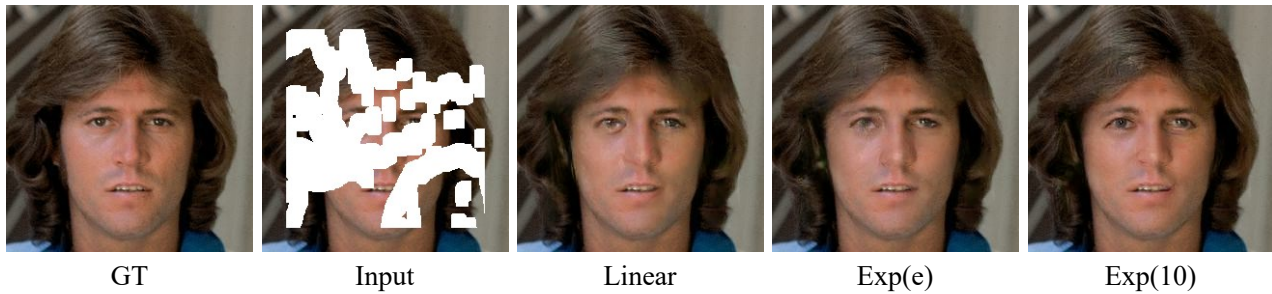To sum up, we combine the weighted reconstruction loss and

| GT | Input | Linear | Exp(e) | Exp(10) |

**Figure 3:** *From left to right: the ground-truth image, the corrupted image, result using linear function, result using exponential function with e and 10 as base separately. The edge ratios after canny algorithm [Can86] with* σ = 1.0 *are 10.14% (GT), 8.61% (Line), 8.72% (Exp(e)) and 9.39% (Exp(10)).*

focal loss as the final loss function to train the whole framework, which introduces the min-max adversarial process as follows:

$$\begin{cases} \min_{G} ||x^{Det(G(\mathrm{I}_{in},\mathrm{M}))} \odot (G(\mathrm{I}_{in},\mathrm{M}) - \mathrm{I}_{gt})||_1, \\ \max_{Det} -\mathcal{L}_{Focal}(Det(G(\mathrm{I}_{in},\mathrm{M})),\mathrm{M}). \end{cases} \quad (8)$$

From Eq. (8), the generator and detector in our framework just dispute about the weight value in the corrupted regions (not involved in non-masked regions), which is an improvement of the global competition in the GAN-based framework to solve image inpainting problem.

## 4. Experimental Results

To validate our proposed method, we quantitatively and qualitatively compare our method with several recent state-of-the-art methods on three public datasets including CelebA-HQ [LLWT15; KALL17], Places2 [ZLK*17], and Paris StreetView [DSG*12]. Moreover, comparisons with relative inpainting frameworks under the same generator verify the effectiveness of our detection-based framework. Ablation study is conducted to choose the appropriate loss for the training of our detector, and the visualization of the valuation is also analyzed to validate the ability of our detector.

### 4.1. Implementation Details

We first describe the details of three public datasets used in our experiments. CelebA-HQ [LLWT15; KALL17] is a high quality face dataset with 30K images, and we randomly select 27K images for training and the remaining 3K images for testing. For Places2 [ZLK*17], we select 30 categories for our training and testing. We randomly sample 2K images from per category in training split of Places2 to construct our training set (60K images). The corresponding 3K images in testing split of Places2 are directly as our testing set. For Paris StreetView [DSG*12], we keep original splits, *i.e.*, 14,900 images for training and 100 images for testing.

For the irregular training masks, we create 180K masks with/without border constraints from the source of QD-IMD [Isk18] that is a collection of 50 million human drawings.

Several data augmentation operations: rotation, dilation, and cropping are adopted in sequence during the mask generation. The irregular mask set from Liu *et al.* [LRS*18] including 12K masks is used as testing masks. Note that both training and testing masks are classified by ranges of the mask ratio from [0.01,0.1] to [0.5,0.6] with the step of 0.1.

In view of computational cost, we resize images and masks to $256 \times 256$ as the input of networks. The Adam optimizer [KB14] is used to optimize the parameters of inpainting models with the learning rate of $10^{-4}$ and $\beta_1 = 0, \beta_2 = 0.9$. We train the model for 100 epochs with the batch size of 8. During the training period, the focusing parameter γ in focal loss (Eq. (6)) is set to 2.0, and the weight mapping mentioned in Section 3.4.3 is set to the exponential transition with 10.0 as the base number.

### 4.2. Comparison with State-of-the-arts

We compare the proposed method with three representative state-of-the-art works that are different kinds of deep inpainting techniques:

- PConv [LRS*18]: The method uses a novel partial convolution instead of standard convolution to solve inpainting problem.
- PEN [ZFCG19]: A pyramid-context encoder network with a series of attention modules for inpainting.
- GConv [YLY*19]: A two-stage inpainting network with gated convolution, which is an upgrade of Contextual Attention [YLY*18].

For a fair comparison, these three existing methods and our method are all trained on three public datasets mentioned above. Officially released source codes of PEN and GConv are obtained from their respective project page [Zen19; Yu19]. As the source code of PConv is not available at the time of experiments, we use an unofficial implementation [Gru19] to train the model under our careful matching with their paper [LRS*18]. For each inpainting technique, the mean inference time of a $256 \times 256$ image on three datasets are recorded: PConv (27.97ms), PEN (61.08ms), GConv (18.91ms) and our method (18.80ms). Because of the one-stage network without special modules like attention layer leveraged in the generator, our method has an advantage on the time cost. Moreover, due to the short inference time of our detector (13.07ms), it is high-efficiency to train our whole framework.
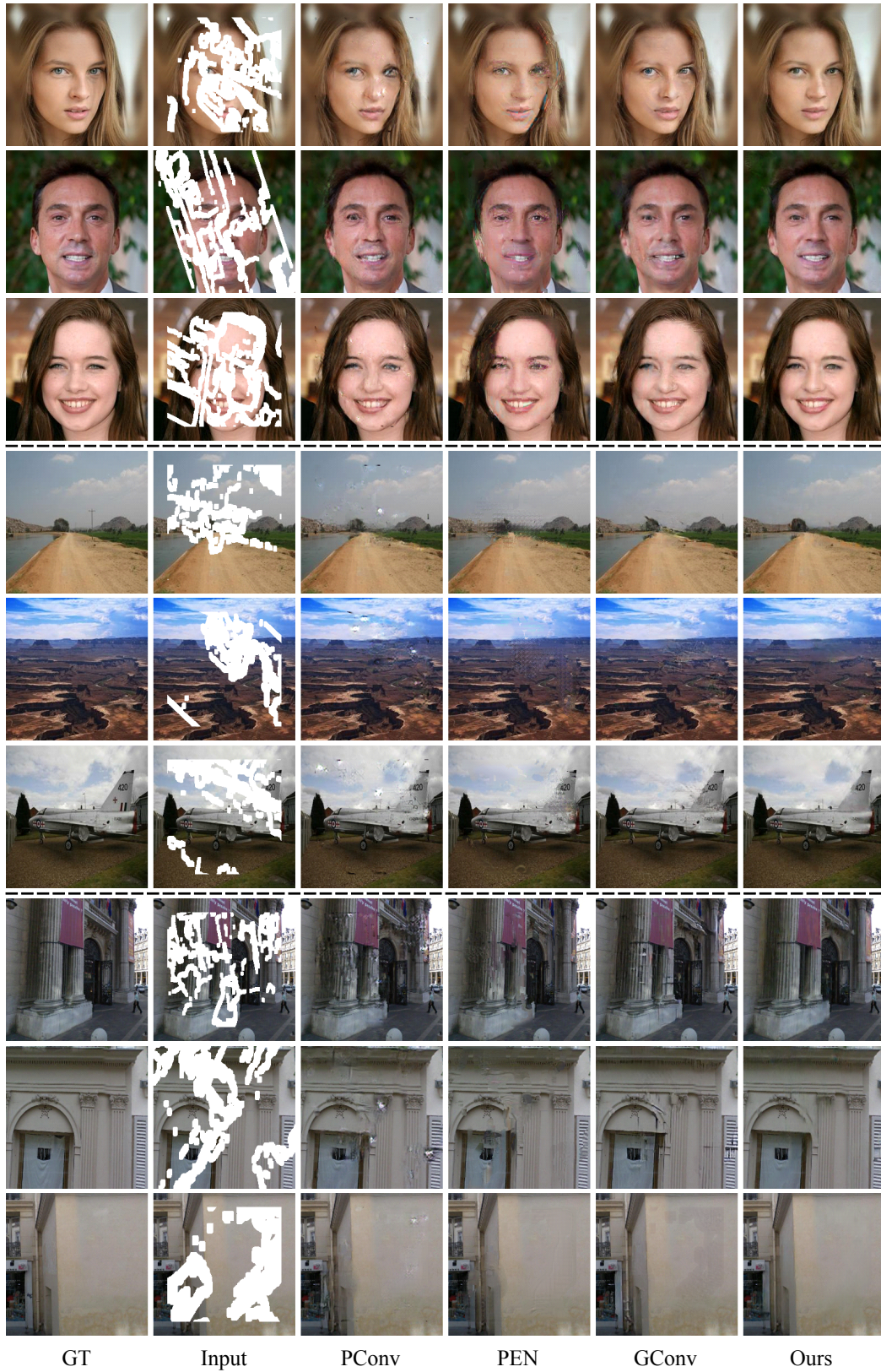
| GT | Input | PConv | PEN | GConv | Ours |

**Figure 4:** *Example results of qualitative comparison. From top to bottom splited three groups from CelebA-HQ, Places2 and Paris StreetView testing set, respectively.*
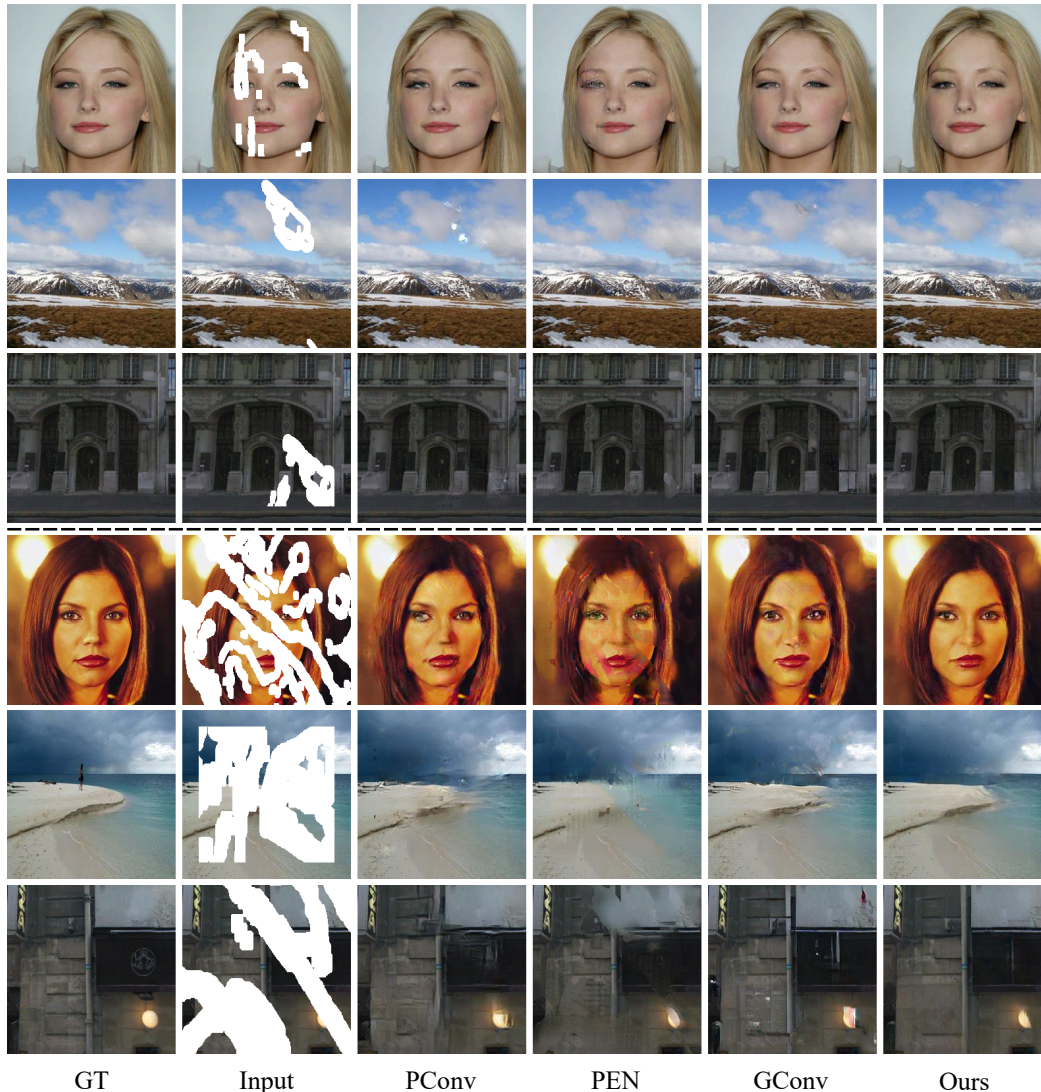
| GT | Input | PConv | PEN | GConv | Ours |

**Figure 5:** *Example results for [0.01, 0.1] (top) and [0.5, 0.6] (bottom) mask from CelebA-HQ, Places2 and Paris StreetView testing set.*

### 4.2.1. Qualitative Comparisons

Fig. 4 shows some example outputs of four different models, *i.e.*, PConv [LRS*18], PEN [ZFCG19], GConv [YLY*19] and our method. There is no post-processing operation to ensure fairness. We observe that PConv [LRS*18] sometimes suffers from obvious visual artifacts and produces some meaningless textures. PEN [ZFCG19] generates checkerboard artifacts in corrupted regions. It also shows poor coherence with background because of over-smoothing results and color inconsistency. GConv [YLY*19] produces better results but still exhibits imperfect details. Our method achieves more plausible results, especially in face image cases. More comparisons are provided in supplemental material. In addition, we specifically illustrate some inpainting results on extreme mask range to verify our framework. As shown in Fig. 5, for the trivial details, such as eyes and brows in face images, our results

are slightly superior to that of other three methods on small mask range. Our method also achieves competitive results on large mask range in Fig. 5 and the following quantitative analysis.

### 4.2.2. Quantitative Comparisons

Multiple reasonable contents combined with contextual background constitute a realistic image, which may be different from the ground-truth image. Because the nature of the non-unique solution of image inpainting problem, numerical metrics are difficult to evaluate the quality of a single inpainting case. However, the mean value of metrics on whole dataset could measure the performance of inpainting techniques. In our work, we follow the previous inpainting works and measure the results with four metrics: $\ell_1$ error, peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) [WBSS04] and Fréchet Inception Distance (FID) [HRU*17]. The first three metrics calculate pixel-wise de-
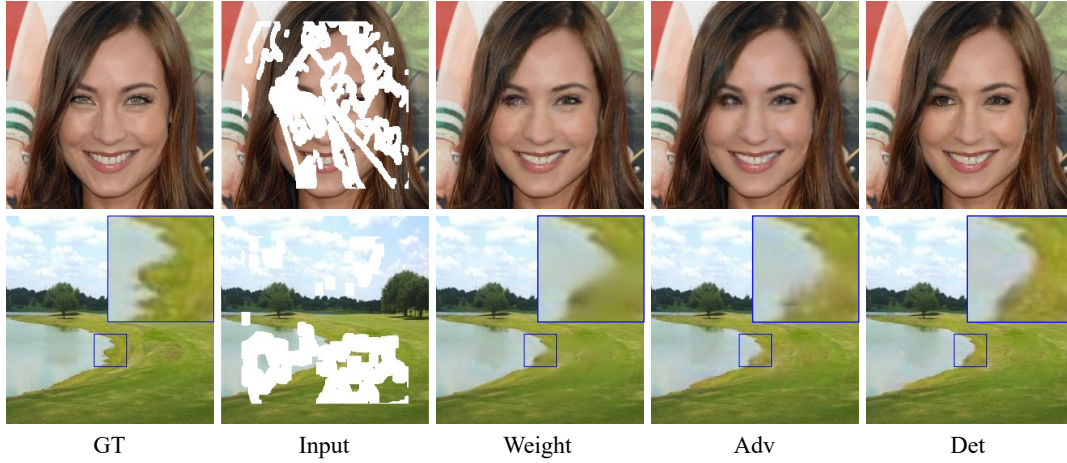
| GT | Input | Weight | Adv | Det |

**Figure 6:** *Comparison between results of hard-weighted $\ell_1$ loss based framework (Weight), typical adversarial framework (Adv) and detection-based generative framework (Det).*

**Table 1:** *Quantitative comparison of four different inpainting methods on CelebA-HQ and Places2. "C" stands for CelebA-HQ and "P" stands for Places2. Note that each statistic is calculated over all testing set in a fixed mask order. † Lower is better. ¶ Higher is better.*

|  | Mask | (0.01-0.1] | | (0.1-0.2] | | (0.2-0.3] | | (0.3-0.4] | | (0.4-0.5] | | (0.5-0.6] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Data | C | P | C | P | C | P | C | P | C | P | C | P |
| $\ell_1$ (%)† | PConv | 0.84 | 1.27 | 2.05 | 3.07 | 3.57 | 5.16 | 5.43 | 7.37 | 7.59 | 9.94 | 12.00 | 14.14 |
|  | PEN | 0.80 | 1.11 | 2.15 | 2.96 | 3.88 | 5.25 | 5.83 | 7.59 | 8.02 | 10.19 | 11.77 | 13.95 |
|  | GConv | **0.65** | **1.00** | 1.81 | 2.71 | 3.41 | 4.98 | 5.33 | 7.44 | 7.53 | 10.22 | 12.05 | 14.97 |
|  | Ours | 0.69 | **1.00** | **1.62** | **2.36** | **2.82** | **4.09** | **4.28** | **5.98** | **6.06** | **8.18** | **10.05** | **12.10** |
| PSNR¶ | PConv | 34.88 | 31.99 | 30.00 | 27.27 | 27.20 | 24.79 | 24.95 | 23.07 | 23.12 | 21.58 | 20.33 | 19.56 |
|  | PEN | 35.34 | 33.24 | 29.76 | 27.72 | 26.79 | 24.85 | 24.70 | 23.04 | 23.06 | 21.58 | 20.85 | 19.84 |
|  | GConv | **37.14** | **34.05** | 31.02 | 28.26 | 27.57 | 25.00 | 25.03 | 22.84 | 23.10 | 21.16 | 20.22 | 18.88 |
|  | Ours | 36.37 | 33.51 | **32.02** | **29.43** | **29.13** | **26.64** | **26.81** | **24.68** | **24.89** | **23.04** | **21.79** | **20.74** |
| SSIM¶ | PConv | 0.983 | 0.961 | 0.963 | 0.917 | 0.936 | 0.866 | 0.898 | 0.809 | 0.851 | 0.738 | 0.744 | 0.615 |
|  | PEN | 0.988 | 0.975 | 0.965 | 0.929 | 0.933 | 0.867 | 0.894 | 0.801 | 0.849 | 0.723 | 0.764 | 0.605 |
|  | GConv | **0.991** | **0.978** | 0.971 | 0.936 | 0.941 | 0.876 | 0.902 | 0.808 | 0.856 | 0.73 | 0.750 | 0.594 |
|  | Ours | 0.990 | 0.977 | **0.977** | **0.949** | **0.958** | **0.907** | **0.930** | **0.856** | **0.894** | **0.792** | **0.793** | **0.665** |
| FID† | PConv | 2.79 | 4.44 | 4.42 | 8.38 | 5.60 | 12.7 | 7.73 | 17.46 | 11.02 | 24.3 | 15.16 | **32.75** |
|  | PEN | 1.41 | 3.1 | 4.19 | 8.76 | 8.38 | 17.31 | 12.68 | 29.84 | 18.73 | 47.05 | 23.38 | 66.7 |
|  | GConv | **0.78** | **2.27** | 2.05 | 6.02 | 3.93 | 11.02 | 5.86 | 16.39 | 8.64 | **22.75** | **12.75** | 33.17 |
|  | Ours | 1.08 | 2.82 | **1.86** | **5.78** | **3.34** | **10.38** | **5.16** | **16.18** | **7.84** | 23.89 | 15.34 | 37.18 |

viation under the assumption that recovery regions target to the ground-truth images. FID based on semantic measurement scales the Wasserstein-2 distance between distributions of real and reconstructed images with a pre-trained Inception-V3 model [SVI*16]. Table 1 reports the quantitative comparison results of all four methods on CelebA-HQ and Places2 dataset. Our method performs best for all measure in the range [0.1, 0.4], and we could still achieve competitive results with other works in too low/too high intervals. As the detector is not sensitive to the mask with much small mask ratio, our model has less poor performance at the mask range of [0.01, 0.1] , compared with GConv model.

### 4.3. Comparisons with Relative Inpainting Frameworks

We also compare the proposed detection-based generative framework with two relative image inpainting frameworks:

- Weight: The framework just involves the generator without the discriminator. Its objective function usually adopts hard-weighted $\ell_1$ loss, *i.e.*, assign heavy weight ($\lambda_1$) to the corrupted regions, and light weight ($\lambda_2$) to non-masked regions. We set $\lambda_1 = 6$ and $\lambda_2 = 1$, which is the same as PConv [LRS*18].
- Adv: The framework, first used by Context Encoders [PKD*16], combines reconstruction loss with adversarial loss to train the model. Most recent deep inpainting methods including PEN [ZFCG19] and GConv [YLY*19] adopts it as well.
- Det: The detection-based generative framework we proposed applies in image inpainting task.

For a fair comparison, the generator in above three frameworks are the same, and detailed structure is described in Section 3.3. Sampled results of different frameworks are shown in Fig. 6. In the first row, Weight and Adv fail to reconstruct reasonable structures of eyes. As we see zoom-in regions, Weight produces the over-
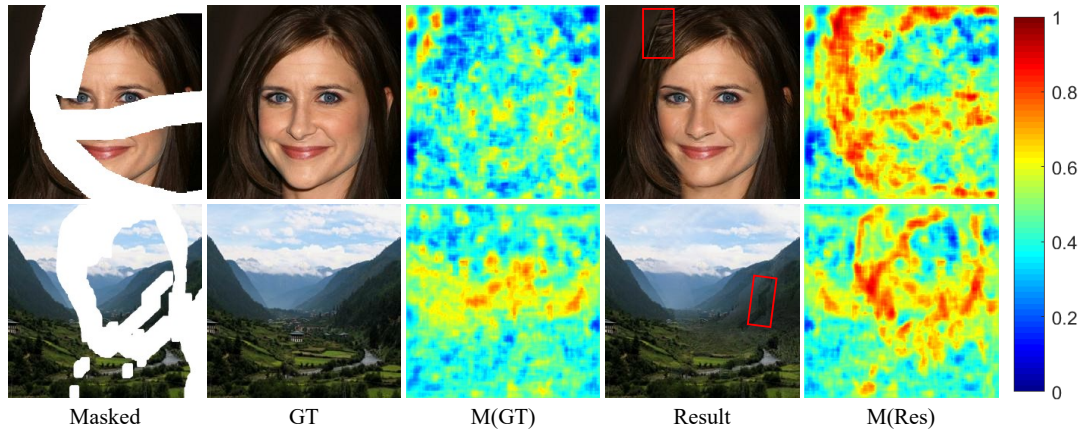
**Figure 7:** *Visualization of the output of the detector.*

**Table 2:** *The evaluation results of ablation study about loss function of detector: the cross entropy loss (CE), the balanced cross entropy loss (BCE), and the focal loss (Focal). Higher is better for both two metrics.*

| Mask | (0.01-0.1] | | (0.1-0.2] | | (0.2-0.3] | | (0.3-0.4] | | (0.4-0.5] | | (0.5-0.6] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| CE | 35.88 | 0.989 | 31.77 | 0.976 | 28.9 | 0.956 | 26.56 | 0.927 | 24.64 | 0.890 | 21.49 | 0.784 |
| BCE | 36.29 | 0.990 | 31.88 | 0.976 | 28.95 | 0.956 | 26.6 | 0.927 | 24.67 | 0.889 | 21.49 | 0.784 |
| Focal | **36.37** | **0.990** | **32.02** | **0.977** | **29.13** | **0.958** | **26.81** | **0.93** | **24.89** | **0.894** | **21.79** | **0.793** |

smoothing result at the border of bank and lake, and the result of Adv is also blurry. Det outperforms all the other frameworks in detail, it is largely because of the advanced competition about "Where are the artifacts?" instead of "Is the image real or fake?".

### 4.4. Ablation Study

In this experiment, we conduct an ablation study to evaluate the performance of several optional detection loss on CelebA-HQ, *i.e.*, the standard cross entropy loss, the balanced loss, and the focal loss. We use the six ranges of the mask ratio mentioned in Section 4.1. The corresponding results are reported in Table 2. We find that the focal loss finally achieves the highest performance, especially for PSNR. Therefore, we choose the focal loss as our loss function to train the detector. However, the results of focal loss are slightly better than that of standard cross entropy loss. The reason is that although the focal loss can solve imbalance problem to improve the detection accuracy [LGG*17], but such improvement implicitly affect the final inpainting results.

### 4.5. Visualization of Detection

The detector is an important component in our proposed framework, and its results are close related to the inpainting quality. Through visualizing the evaluation of the inpainted images and the corresponding ground-truth images, we analyze how the detector perceives the artifacts by weakly supervised learning. In detail, we respectively feed the inpainted image and corresponding ground-truth image into the trained detector, and then visualize the colormaps of the output of the detector. The visualization results are shown in Fig. 7, where red color means the region has artifacts with

high probability, and blue color means low probability. We find that the colormaps of the ground-truth images ("M(GT)" in Fig. 7) are irregular showing a bit relation with image contents, whereas the colormaps of inpainted images ("M(Res)" in Fig. 7) show strong correlation with input masks ("Masked" in Fig. 7), *i.e.*, most of red pixels are in/beside the mask regions. In the meanwhile, the obvious visual artifacts (the regions highlighting with the red rectangle in inpainted results) match with the hottest regions in the colormaps of inpainted results. Specially, in the second row of Fig. 7, the red rectangle in "Result" corresponds to the non-masked region around the missing region in "Masked". This means the detector learns the location of the artifacts rather than the binary mask. Consequently, the output of detector has the ability to indicate the position of defects, which is to some extent consistent with the perception of human eyes, and this useful position information is inserted into reconstruction loss to enhance visual artifacts.

### 4.6. Real-world Applications

We demonstrate some daily applications of our whole framework on image translation. Fig. 8 shows our results of object removal (the left column), text removal (the middle column) and old photo restoration (the right column). We leverage the model trained on Places2 dataset without fine-tuning to conduct object removal task, which appears the harmonious filling contents, especially at the mask boundaries. Due to the large domain gaps between textual images and inpainting benchmarks, we retrain our model on the dataset collected in real world to solve text removal problem. The experimental results show that our method can handle the complex illumination and noise in real scenarios. Following [WZC*20], we synthesize amount of training data to train the model and then re-

**Figure 8:** *Example results of daily applications. From left to right, three split groups are results of object removal, text removal and old photo restoration, respectively. For each pair of images, the top image is the input and the bottom image is the image translation result.*

pair the old photos, and the results indicate that our method can recover unstructured degradation and structured scratches.

## 5. Discussion and Conclusion

In this paper, we proposed a detection-based generative framework to address image inpainting problem. To perceptually localize visual artifacts in inpainted images, we introduced a dense detector, which is trained by weakly supervised learning, to evaluate the quality of inpainted images in a pixel-wise manner. Furthermore, the reconstruction loss is combined with such evaluation using a weighting criterion to train the generator, which avoid tuning the tradeoff parameters manually. Extensive experiments demonstrate the superiority of proposed detection-based image inpainting framework. However, semantic information may disturb our detector to omit feeble artifacts existing in inpainted images during training period. From Table 1, this obstruction is apparent in case of small mask ratio. As Fig. 9 shown, the detector easily captures artifacts from heavy scratches (the region highlighting with the red rectangle) in the second row, while the detector pays more attention to semantic information instead of tiny scratches in the first row. For large mask ratio, mask with too rough information of artifacts location may not be suitable to weakly supervised learning of the detector. Accurate artifacts localization is an open problem for the detection-based framework to solve. In future work, this framework may extend to other conditional generative tasks, *e.g.*, image synthesis and image denoising. We also plan to implement our approach with the open source platform, Jittor [HLY*20], for shorter training time.
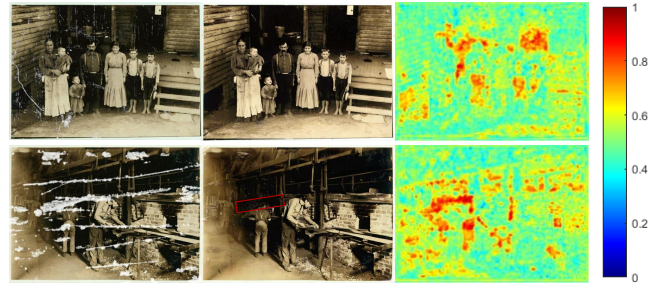
## Acknowledgements

**Figure 9:** *Visualization of the old photo restoration. From right to left are input images, repaired results and visualizations of repaired results, respectively.*

## References

[BBC*01] BALLESTER, COLOMA, BERTALMIO, MARCELO, CASELLES, VICENT, et al. "Filling-in by joint interpolation of vector fields and gray levels". *IEEE Trans. Image Process.* 10.8 (2001), 1200–1211 2.

[BSCB00] BERTALMIO, MARCELO, SAPIRO, GUILLERMO, CASELLES, VINCENT, and BALLESTER, COLOMA. "Image inpainting". *Proc. ACM SIGGRAPH*. 2000, 417–424 2.

[BSFG09] BARNES, CONNELLY, SHECHTMAN, ELI, FINKELSTEIN, ADAM, and GOLDMAN, DAN B. "PatchMatch: A randomized correspondence algorithm for structural image editing". *ACM Trans. Graph.* 28.3 (2009), 24 2.

[Can86] CANNY, JOHN. "A computational approach to edge detection". *IEEE Trans. Pattern Anal. Mach. Intell.* 6 (1986), 679–698 6.

[CWYZ04] CHEN, DI-RONG, WU, QIANG, YING, YIMING, and ZHOU, DING-XUAN. "Support vector machine soft margin classifiers: error analysis". *Journal of Machine Learning Research* 5.Sep (2004), 1143–1175 5.

[DAFC19] DANON, DOV, AVERBUCH-ELOR, HADAR, FRIED, OHAD, and COHEN-OR, DANIEL. "Unsupervised natural image patch learning". *Computational Visual Media* 5.3 (2019), 229–237 2.

[DDS*09] DENG, JIA, DONG, WEI, SOCHER, RICHARD, et al. "Imagenet: A large-scale hierarchical image database". *IEEE CVPR*. 2009, 248–255 3.

[DSB*12] DARABI, SOHEIL, SHECHTMAN, ELI, BARNES, CONNELLY, et al. "Image Melding: combining inconsistent images using patch-based synthesis". *ACM Trans. Graph. (Proc. SIGGRAPH)* 31.4 (2012), 1–10 2.

[DSG*12] DOERSCH, CARL, SINGH, SAURABH, GUPTA, ABHINAV, et al. "What Makes Paris Look like Paris?": *ACM Trans. Graph. (Proc. SIGGRAPH)* 31.4 (2012), 101:1–101:9 6.

[DSS93] DRUCKER, HARRIS, SCHAPIRE, ROBERT, and SIMARD, PATRICE. "Boosting performance in neural networks". *Advances in Pattern Recognition Systems using Neural Network Technologies*. World Scientific, 1993, 61–75 5.

[EF01] EFROS, ALEXEI A and FREEMAN, WILLIAM T. "Image quilting for texture synthesis and transfer". *Proc. ACM SIGGRAPH*. 2001, 341–346 2.

[GEB16] GATYS, LEON A., ECKER, ALEXANDER S., and BETHGE, MATTHIAS. "Image Style Transfer Using Convolutional Neural Networks". *IEEE CVPR*. 2016, 2414–2423 3.

[GPM*14] GOODFELLOW, IAN, POUGET-ABADIE, JEAN, MIRZA, MEHDI, et al. "Generative adversarial nets". *Advances in Neural Information Processing Systems*. 2014, 2672–2680 2, 3.

[Gru19] GRUBER, MATHIAS. "Unofficial implementation of "Image Inpainting for Irregular Holes Using Partial Convolutions"". https://github.com/MathiasGruber/PConv-Keras. Accessed 5 Match 2020. 2019 6.

[HKAK14] HUANG, JIA-BIN, KANG, SING BING, AHUJA, NARENDRA, and KOPF, JOHANNES. "Image completion using planar structure guidance". *ACM Trans. Graph. (Proc. SIGGRAPH)* 33.4 (2014), 1–10 2.

[HLY*20] HU, SHI-MIN, LIANG, DUN, YANG, GUO-YE, et al. "Jittor: A Noval Deep Learning Framework with Unified Graph Execution and Meta Operators". *Science China-Information Sciences* (2020). to appear. https://github.com/Jittor/Jittor. 11.

[HRU*17] HEUSEL, MARTIN, RAMSAUER, HUBERT, UNTERTHINER, THOMAS, et al. "Gans trained by a two time-scale update rule converge to a local nash equilibrium". *Advances in Neural Information Processing Systems*. 2017, 6626–6637 8.

[HZRS16] HE, KAIMING, ZHANG, XIANGYU, REN, SHAOQING, and SUN, JIAN. "Deep residual learning for image recognition". *IEEE CVPR*. 2016, 770–778 2, 4.

[ISI17] IIZUKA, SATOSHI, SIMO-SERRA, EDGAR, and ISHIKAWA, HIROSHI. "Globally and locally consistent image completion". *ACM Trans. Graph. (Proc. SIGGRAPH)* 36.4 (2017), 1–14 2, 3.

[Isk18] ISKAKOV, KARIM. "Quick Draw Irregular Mask Dataset". https://github.com/karfly/qd-imd. Accessed 5 Match 2020. 2018 6.

[IZZE17] ISOLA, PHILLIP, ZHU, JUN-YAN, ZHOU, TINGHUI, and EFROS, ALEXEI A. "Image-to-image translation with conditional adversarial networks". *IEEE CVPR*. 2017, 1125–1134 3.

[JAF16] JOHNSON, JUSTIN, ALAHI, ALEXANDRE, and FEI-FEI, LI. "Perceptual losses for real-time style transfer and super-resolution". *ECCV*. 2016, 694–711 3.

[KALL17] KARRAS, TERO, AILA, TIMO, LAINE, SAMULI, and LEHTINEN, JAAKKO. "Progressive growing of gans for improved quality, stability, and variation". *arXiv preprint arXiv:1710.10196* (2017) 6.

[KB14] KINGMA, DIEDERIK P and BA, JIMMY. "Adam: A method for stochastic optimization". *arXiv preprint arXiv:1412.6980* (2014) 6.

[LGG*17] LIN, TSUNG-YI, GOYAL, PRIYA, GIRSHICK, ROSS, et al. "Focal loss for dense object detection". *IEEE ICCV*. 2017, 2980–2988 5, 10.

[LJXY19] LIU, HONGYU, JIANG, BIN, XIAO, YI, and YANG, CHAO. "Coherent semantic attention for image inpainting". *IEEE ICCV*. 2019, 4170–4179 3.

[LLWT15] LIU, ZIWEI, LUO, PING, WANG, XIAOGANG, and TANG, XIAOOU. "Deep learning face attributes in the wild". *IEEE ICCV*. 2015, 3730–3738 6.

[LRS*18] LIU, GUILIN, REDA, FITSUM A, SHIH, KEVIN J, et al. "Image inpainting for irregular holes using partial convolutions". *ECCV*. 2018, 85–100 2, 3, 6, 8, 9.

[LSD15] LONG, JONATHAN, SHELHAMER, EVAN, and DARRELL, TREVOR. "Fully convolutional networks for semantic segmentation". *IEEE CVPR*. 2015, 3431–3440 2.

[NNJ*19] NAZERI, KAMYAR, NG, ERIC, JOSEPH, TONY, et al. "Edge-Connect: Generative image inpainting with adversarial edge learning". *arXiv preprint arXiv:1901.00212* (2019) 2–4.

[PKD*16] PATHAK, DEEPAK, KRAHENBUHL, PHILIPP, DONAHUE, JEFF, et al. "Context encoders: Feature learning by inpainting". *IEEE CVPR*. 2016, 2536–2544 2, 3, 9.

[RYZ*19] REN, YURUI, YU, XIAOMING, ZHANG, RUONAN, et al. "StructureFlow: Image Inpainting via Structure-aware Appearance Flow". *IEEE ICCV*. 2019, 181–190 2, 3.

[SCSI08] SIMAKOV, DENIS, CASPI, YARON, SHECHTMAN, ELI, and IRANI, MICHAL. "Summarizing visual data using bidirectional similarity". *IEEE CVPR*. 2008, 1–8 2.

[SVI*16] SZEGEDY, CHRISTIAN, VANHOUCKE, VINCENT, IOFFE, SERGEY, et al. "Rethinking the inception architecture for computer vision". *IEEE CVPR*. 2016, 2818–2826 9.

[SZ14] SIMONYAN, KAREN and ZISSERMAN, ANDREW. "Very deep convolutional networks for large-scale image recognition". *arXiv preprint arXiv:1409.1556* (2014) 3.

[UVL17] ULYANOV, DMITRY, VEDALDI, ANDREA, and LEMPITSKY, VICTOR. "Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis". *IEEE CVPR*. 2017, 6924–6932 5.

[WBSS04] WANG, ZHOU, BOVIK, ALAN C, SHEIKH, HAMID R, and SIMONCELLI, EERO P. "Image quality assessment: from error visibility to structural similarity". *IEEE Trans. Image Process.* 13.4 (2004), 600–612 8.

[WZC*20] WAN, ZIYU, ZHANG, BO, CHEN, DONGDONG, et al. "Bringing Old Photos Back to Life". *IEEE CVPR*. 2020, 2747–2757 10.

[XLL*19] XIE, CHAOHAO, LIU, SHAOHUI, LI, CHAO, et al. "Image Inpainting with Learnable Bidirectional Attention Maps". *IEEE ICCV*. 2019, 8858–8867 3.

[YK15] YU, FISHER and KOLTUN, VLADLEN. "Multi-scale context aggregation by dilated convolutions". *arXiv preprint arXiv:1511.07122* (2015) 3, 5.

[YLL*17] YANG, CHAO, LU, XIN, LIN, ZHE, et al. "High-resolution image inpainting using multi-scale neural patch synthesis". *IEEE CVPR*. 2017, 6721–6729 2, 3.

[YLL*18] YAN, ZHAOYI, LI, XIAOMING, LI, MU, et al. "Shift-Net: Image Inpainting via Deep Feature Rearrangement". *ECCV*. 2018, 3–19 2.

[YLY*18] YU, JIAHUI, LIN, ZHE, YANG, JIMEI, et al. "Generative image inpainting with contextual attention". *IEEE CVPR*. 2018, 5505–5514 2, 3, 6.

[YLY*19] YU, JIAHUI, LIN, ZHE, YANG, JIMEI, et al. "Free-form image inpainting with gated convolution". *IEEE ICCV*. 2019, 4471–4480 2, 3, 6, 8, 9.

[Yu19] YU, JIAHUI. "DeepFill v1/v2 with Contextual Attention and Gated Convolution". https://github.com/JiahuiYu/generative_inpainting. Accessed 5 Match 2020. 2019 6.

[Zen19] ZENG, YANHONG. "Official implement for Learning Pyramid-Context Encoder Network for High-Quality Image Inpainting". https://github.com/researchmm/PEN-Net-for-Inpainting. Accessed 5 Match 2020. 2019 6.

[ZFCG19] ZENG, YANHONG, FU, JIANLONG, CHAO, HONGYANG, and GUO, BAINING. "Learning pyramid-context encoder network for high-quality image inpainting". *IEEE CVPR*. 2019, 1486–1494 2, 3, 6, 8, 9.

[ZLK*17] ZHOU, BOLEI, LAPEDRIZA, AGATA, KHOSLA, ADITYA, et al. "Places: A 10 million image database for scene recognition". *IEEE Trans. Pattern Anal. Mach. Intell.* 40.6 (2017), 1452–1464 6.