# MAT: Mask-Aware Transformer for Large Hole Image Inpainting

Wenbo Li[1]    Zhe Lin[2]    Kun Zhou[3]    Lu Qi[1]    Yi Wang[4*]    Jiaya Jia[1]

[1]The Chinese University of Hong Kong    [2]Adobe Inc.

[3]The Chinese University of Hong Kong (Shenzhen)    [4]Shanghai AI Laboratory

{wenboli,luqi,leojia}@cse.cuhk.edu.hk

zlin@adobe.com    kunzhou@link.cuhk.edu.cn    wangyi@pjlab.org.cn
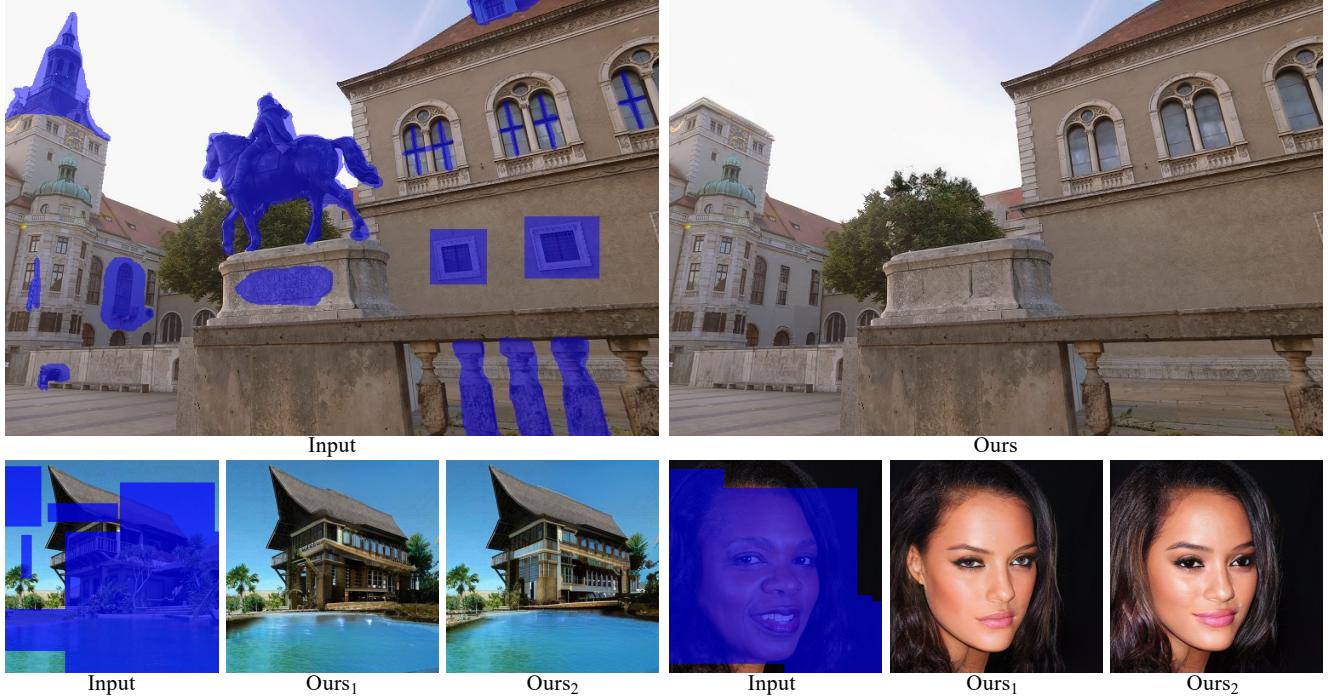
Figure 1. The proposed MAT supports photo-realistic and pluralistic large hole image inpainting. The first example is a real-world high-resolution image and the other two examples ($512 \times 512$) are from Places [78] and FFHQ [26] datasets.

## Abstract

*Recent studies have shown the importance of modeling long-range interactions in the inpainting problem. To achieve this goal, existing approaches exploit either standalone attention techniques or transformers, but usually under a low resolution in consideration of computational cost. In this paper, we present a novel transformer-based model for large hole inpainting, which unifies the merits of transformers and convolutions to efficiently process high-resolution images. We carefully design each component of our framework to guarantee the high fidelity and diversity of recovered images. Specifically, we customize an inpainting-oriented transformer block, where the attention module aggregates non-local information only from partial valid tokens, indicated by a dynamic mask. Extensive experiments demonstrate the state-of-the-art performance of the new model on multiple benchmark datasets. Code is released at https://github.com/fenglinglwb/MAT.*

## 1. Introduction

Image completion (a.k.a. inpainting) is a fundamental problem in computer vision, which aims to fill missing regions with plausible contents. It has many applications including image editing [23], image re-targeting [9], photo

---
*Corresponding author

restoration [53, 54] and object removal [3].

In inpainting, modeling the contextual information is crucial, especially for large masks. Creating reasonable structures and textures for the missing areas demands contextual understanding, using distant information according to non-local priors [4, 7, 38, 56] in images. Previous works employ stacked convolutions to reach large receptive fields and model long-range relationships, which works well on aligned (*e.g.*, faces, bodies) and texture-heavy (*e.g.*, forests, water) data. When processing images with complicated structures (*i.e.*, the first example in the $2_{nd}$ row in Figure 1), it is difficult for fully convolutional neural networks (CNNs) to characterize the semantic correspondences between distant areas. This is mainly due to the inherent properties of CNNs, the slow growth of the effective receptive field and the inevitable dominance of nearby pixels. To explicitly model long-range dependencies in inpainting, [61, 65, 66] propose to employ attention modules in the CNN-based generator. However, limited by the quadratic computational complexity, the attention module is merely applied to relatively small-scale feature maps with a few times, where long-range modeling is not fully exploited.

In contrast to applying attention modules to CNNs, transformer [52] is a natural architecture to handle non-local modeling, where attention is a basic component in every block. Recent advances [55, 68, 77] adopt transformer structures to address the inpainting problem. Nonetheless, affected by the complexity issue, existing works only employ transformers to infer low-resolution predictions (*e.g.* $32 \times 32$) for subsequent processing, hence the produced image structure is coarse, compromising the final image quality, especially on large-scale masks.

In this paper, we develop a new inpainting transformer, capable of generating high-resolution completed results for large mask inpainting. Due to the lack of useful information in some regions (this is common when the given mask rules out most pixels), we find the commonly utilized transformer block (LN→MSA→LN→FFN) exhibits inferior performance in adversarial training. In this regard, we customize the vanilla Transformer block to increase optimization stability and also improve performance, by removing the conventional layer normalization [1] and replacing the residual learning with fusion learning using feature concatenation. We analyze why these modifications are crucial for learning and empirically demonstrate they are non-trivial. Also, to handle possible heavy interactions between all tokens extracted from the high-resolution input, we propose a new variant of multi-head self-attention (MSA), named multi-head contextual attention (MCA). It computes non-local relations only using partial valid tokens. The selection of adopted tokens is indicated by a dynamic mask, which is initialized by the input mask and updated with spatial constraints and long-range interactions,

improving the efficiency at no cost of effectiveness. Additionally, we incorporate a novel style manipulation module into the proposed framework, inherently supporting pluralistic generation. As shown in Fig. 1, our method successfully fills large holes with visually realistic and exceptionally diverse contents. Our contributions are summarized as:

- We develop a novel inpainting framework MAT. It is the first transformer-based inpainting system capable of directly processing high-resolution images.

- We meticulously design components of MAT. The proposed multi-head contextual attention conducts long-range dependency modeling efficiently by exploiting valid tokens, indicated by a dynamic mask. We also propose a modified transformer block to make training large masks more stable. Moreover, we design a novel style manipulation module to improve diversity.

- MAT sets new state of the arts on multiple benchmark datasets including Places [78] and CelebA-HQ [25]. It also enables pluralistic completion.

## 2. Related Work

Image completion has been a longstanding problem in computer vision. Early diffusion-based methods [2, 6] propagate neighboring undamaged information to the holes. Within an internal or external searching space, patch-based or exemplar-based approaches [10–12, 19, 28, 30, 50] borrow patches with similar appearance based on human-defined distance metrics to complete missing regions. Patch-Match [3] proposes a multi-scale patch searching strategy to accelerate the inpainting process. Moreover, partial differential equation [5, 17] and global or local image statistics [14, 15, 31] are vastly studied in the literature. Though traditional methods can often obtain visually realistic results, the lack of high-level understanding hinders them from generating semantically reasonable contents.

In the few years, deep learning has achieved great success on the image completion. Pathak *et al*. [42] bring the adversarial training [16] to inpainting and utilize an encoder-decoder-based architecture to fill holes. Afterwards, numerous variants [34, 57, 64, 69] of the U-Net structure [45] have been developed for image completion. Besides, more sophisticated networks or learning strategies are proposed to generate high-quality images, including global and local discrimination [22], contextual attention [35, 61, 65, 66], partial [33] and gated [67] convolution, *etc*. Multi-stage generation has also received a great amount of attention, where intermediate clues like object edges [40], foreground contours [63], structures [44] and semantic segmentation maps [49] are extensively exploited. To allow for high-resolution image inpainting, a few attempts have been made to develop progressively generation systems, such as [18, 32, 41, 71, 72].
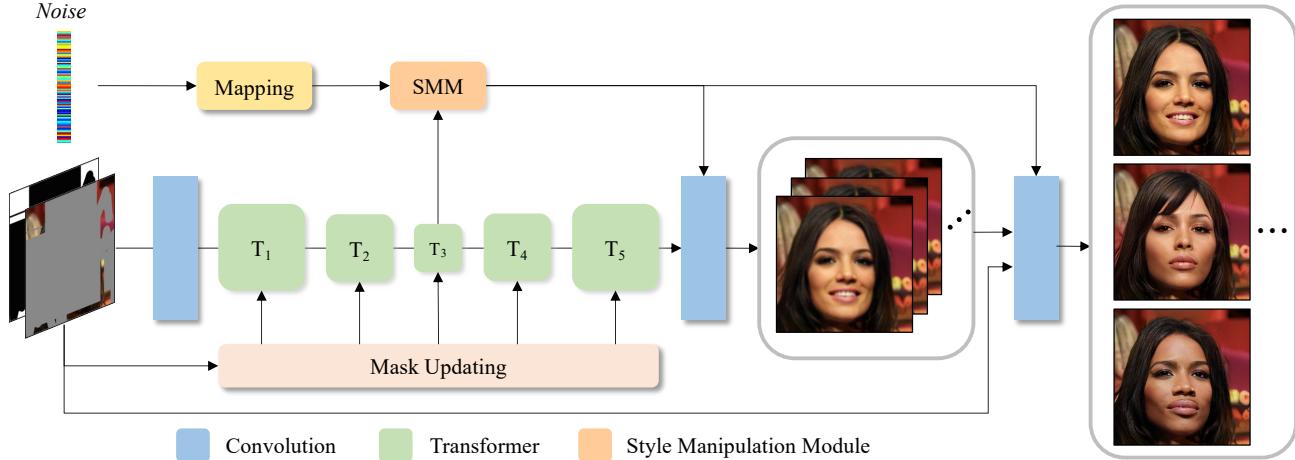
Figure 2. The proposed mask-aware transformer (MAT) for pluralistic inpainting, which consists of a convolutional head, a transformer body and a convolutional tail for reconstruction together with a Conv-U-Net for refinement. The mask updating strategy is described in Sec. 3.3.2.

Recently, researchers switch their focus to more challenging settings, among which the most representative problems are pluralistic generation and large hole filling. For the former, Zheng *et al.* [76] propose a probabilistically principled framework with two parallel paths, capable of producing multiple plausible solutions. UCTGAN [74] projects the instance image space and masked image space into a common low-dimensional manifold space via optimizing the KL-divergence to allow diverse generations of missing contents. Later on, [55] and [68] take advantage of bidirectional attention or auto-regressive transformers to accomplish this goal. Although these methods improve the diversity, their completion and inference performances are limited due to the variational training and raster-scan-order-based generation. On the other hand, some works [37, 51, 75, 77] are proposed to solve the large hole inpainting problem. For example, CoModGAN [75] leverages the modulation techniques [8, 26, 27] to handle large-scale missing regions. In this work, we develop a novel framework to simultaneously achieve high-quality pluralistic generation and large hole filling, bringing the best of long-range context interaction and unconditional generation to the image completion task.

## 3. Method

Given a masked image, formulated as $\mathbf{I}_M = \mathbf{I} \odot \mathbf{M}$, image completion aims to hallucinate visually appealing and semantically appropriate contents for missing areas. In this work, we present a mask-aware transformer (MAT) for large mask inpainting, supporting conditional long-range interactions. Besides, in light of the ill-posed nature of image completion problem, *i.e.*, there could be numerous plausible solutions to fill the large holes, our approach is designed to support pluralistic generation.

## 3.1. Overall Architecture

As shown in Fig. 2, our proposed MAT architecture consists of a convolutional head, a transformer body, a convolutional tail and a style manipulation module, bringing the merits of transformers and convolutions into full play. Specifically, a convolutional head is used to extract tokens, then the main body with five stages of transformer blocks at varying resolutions (with different numbers of tokens) models long-range interactions via the proposed multi-head contextual attention (MCA). For the output tokens from the body, a convolution-based reconstruction module is adopted to upsample the spatial resolution to the input size. Moreover, we adopt another Conv-U-Net to refine high-frequency details, leaning upon the local texture refinement capability and efficiency of CNNs. At last, we design a style manipulation module, enabling the model to deliver diverse predictions by modulating the weights of convolutions. All components in our method are detailed below.

## 3.2. Convolutional Head

The convolutional head takes in the incompleted image $\mathbf{I}_M$ and the given mask $\mathbf{M}$, and produces $^1/_8$ sized feature maps used for tokens. It contains four convolutional layers, one for changing the input dimension and others for downsampling the resolution.

We utilize a convolutional head mainly for two reasons. First, the incorporation of local inductive priors in early visual processing remains vital for better representation [43] and optimizability [60]. On the other hand, it is designed for fast downsampling to reduce computational complexity and memory cost. Also, we empirically find this design is better than the linear projection head used in ViT [13], as validated in the supplementary material.
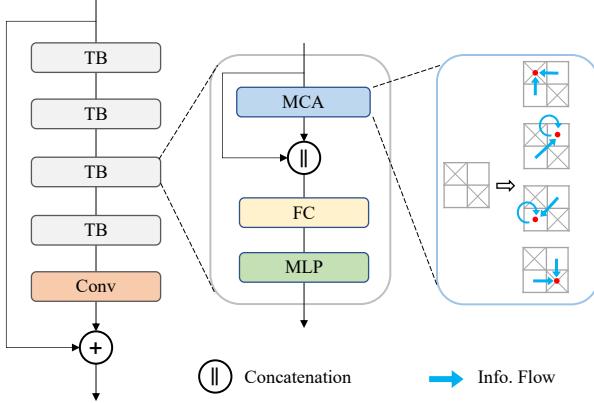
Figure 3. Structure of a single transformer stage. "TB" refers to an adjusted transformer block and "MCA" represents the proposed multi-head contextual attention. The valid tokens are denoted as □ and invalid tokens are ⊠. The blue arrow indicates the output of attention is computed as the weighted sum of valid tokens (indicated by blue arrows) while ignoring invalid tokens.

## 3.3. Transformer Body

The transformer body processes tokens by building long-range correspondences. It contains five stages of the proposed adjusted transformer blocks, with an efficient attention mechanism guided by an additional mask.

### 3.3.1 Adjusted Transformer Block

We propose a new transformer block variant to handle the optimization of masks with large holes. In detail, we remove the layer normalization (LN) [1] and employ fusion learning (using feature concatenation) instead of residual learning. As shown in Fig. 3, we concatenate the input and output of attention and use a fully connected (FC) layer:

$$\mathbf{X}'_{k,\ell} = \text{FC}(\,[\text{MCA}(\mathbf{X}_{k,\ell-1}),\ \mathbf{X}_{k,\ell-1}]\,)\,, \qquad (1)$$

$$\mathbf{X}_{k,\ell} = \text{MLP}(\mathbf{X}'_{k,\ell})\,, \qquad (2)$$

where $\mathbf{X}_{k,\ell}$ is the output of the MLP module of the $\ell$-th block in the $k$-th stage. After several transformer blocks, as illustrated in Fig. 3, we adopt a convolution layer with a global residual connection. Note that we abandon the positional embedding in the transformer block since [59, 62] have shown that $3 \times 3$ convolutions are sufficient to provide positional information for transformers. Thus, the flowing only depends on the feature similarity, which promotes long-range interactions.

**Analysis.** The general architecture of transformer [52] contains two sub-modules, a multi-head self-attention (MSA) module and an MLP module. Layer normalization is applied before every module and a residual connection [20]

after every module. Whereas, we observe unstable optimization using the general block when handling large-scale masks, sometimes incurring gradient exploding. We attribute this training issue to the large ratio of invalid tokens (their values are nearly zero). In this circumstance, layer normalization may magnify useless tokens overwhelmingly, leading to unstable training. Besides, residual learning generally encourages the model to learn high-frequency contents. However, considering most tokens are invalid at the beginning, it is difficult to directly learn high-frequency details without proper low-frequency basis in GAN training, which makes the optimization harder. Replacing such residual learning with concatenation leads to obviously superior results, as verified in Sec. 4.3.

### 3.3.2 Multi-Head Contextual Attention

To handle a large number of tokens (up to 4096 tokens for $512 \times 512$ images) and low fidelity in the given tokens (at most 90% tokens are useless), our attention module exploits shifted windows [36] and a dynamical mask, capable of conducting non-local interactions using a few feasible tokens. The output is computed as the weighted sum of valid tokens, as shown in Fig. 3, which is formulated as

$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^T + \mathbf{M}'}{\sqrt{d_k}})\mathbf{V}\,, \qquad (3)$$

where $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are the query, key, value matrices and $\frac{1}{\sqrt{d_k}}$ is the scaling factor. The mask $\mathbf{M}'$ is expressed as:

$$\mathbf{M}'_{ij} = \begin{cases} 0, & \text{if token } j \text{ is valid}\,, \\ \text{-}\tau, & \text{if token } j \text{ is invalid}\,, \end{cases} \qquad (4)$$

where $\tau$ is a large positive integer (100 in our experiments). In this case, the aggregation weights of invalid tokens are nearly 0. After each attention, we shift the positions of $w \times w$ sized windows by $(\lfloor \frac{w}{2} \rfloor, \lfloor \frac{w}{2} \rfloor)$ pixels, enabling cross-window connections.

**Mask Updating Strategy.** The mask ($\mathbf{M}'$) points out whether a token is valid or invalid, which is initialized by the input mask and automatically updated during propagation. The updating follows a rule that all tokens in a window are updated to be valid after attention as long as there is at least one valid token before. If all tokens in a window are invalid, they remain invalid after attention. As shown in Fig. 4, going through an attention from (a) to (b), all tokens in the top left window become valid, while tokens in other windows are still invalid. After several times of window shift and attention, the mask is updated to be fully valid.

**Analysis.** For images dominated by missing regions, the default attention strategy not only fails to borrow visible information to inpaint the holes, but also undermines the
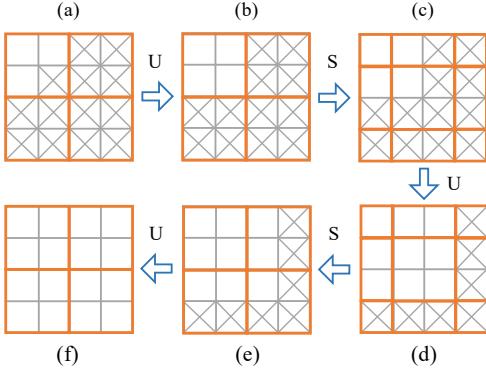
Figure 4. Toy example of mask updating. The feature map is initially partitioned into $2 \times 2$ windows (in orange). "U" means a mask updating after attention and "S" indicates the window shift.

effective valid pixels. To reduce color discrepancy or blurriness, we propose to only involve valid tokens (selected by a dynamic mask) for computing relations. The effectiveness of our design is manifested in Sec 4.3.

### 3.4. Style Manipulation Module

Inspired by [8, 26, 27], we design a style manipulation module to endow our framework with pluralistic generation. It manipulates the output by changing the weight normalization of convolution layers in the reconstruction procedure with an additional noise input. To enhance the representation ability of noise inputs, we enforce the image-conditional style $\mathbf{s}_c$ to learn from both the image feature $\mathbf{X}$ and the noise-unconditional style $\mathbf{s}_u$, formulated as

$$\mathbf{s}_u = \mathcal{E}(\mathbf{n}), \tag{5}$$

$$\mathbf{X}^{'} = \mathbf{B} \odot \mathbf{X} + (\mathbf{1} - \mathbf{B}) \odot \text{Resize}(\mathbf{s}_u), \tag{6}$$

$$\mathbf{s}_c = \mathcal{F}(\mathbf{X}^{'}), \tag{7}$$

where $\mathbf{B}$ is a random binary mask, on which values are set to 1 with a probability of $p$ and to 0 with $1 - p$, $\mathcal{E}$ and $\mathcal{F}$ are mapping functions. As shown in Fig. 2, the style representation is obtained by fusing both style representations:

$$\mathbf{s} = \mathcal{A}(\mathbf{s}_u, \mathbf{s}_c), \tag{8}$$

where $\mathcal{A}$ is a mapping function. Then the weights $\mathbf{W}$ of convolutions are baked as

$$\mathbf{W}^{'}_{ijk} = \mathbf{W}_{ijk} \cdot \mathbf{s}_i, \tag{9}$$

$$\mathbf{W}^{''}_{ijk} = \mathbf{W}^{'}_{ijk} \Big/ \sqrt{\sum_{i,k} \mathbf{W}^{'}_{ijk}{}^2 + \epsilon}, \tag{10}$$

where $i, j, k$ denotes the input channels, output channels and spatial footprint of the convolution, respectively, and $\epsilon$ is a small constant. The modulation of different style representations leads to pluralistic outputs. Also, we incorporate the noise injection [26] into our framework to further enhance the diversity of generation.

### 3.5. Loss Functions

To improve the quality and diversity of the generation, we adopt the non-saturating adversarial loss [16] for both two stages to optimize our framework, regardless of the pixel-wise MAE or MSE loss that usually leads to averaged blurry results. We also use the $R_1$ regularization [39, 46], written as $R_1 = E_x \|\nabla D(x)\|$. Besides, we adopt the perceptual loss [24] with an empirically low coefficient since we notice it enables easier optimization.

**Adversarial Loss.** We calculate the adversarial loss as

$$\mathcal{L}_{\mathrm{G}} = -\mathbb{E}_{\hat{x}} \left[ \log \left( D \left( \hat{x} \right) \right) \right], \tag{11}$$

$$\mathcal{L}_{\mathrm{D}} = -\mathbb{E}_{x} \left[ \log \left( D \left( x \right) \right) \right] - \mathbb{E}_{\hat{x}} \left[ \log \left( 1 - D \left( \hat{x} \right) \right) \right], \tag{12}$$

where $x$ and $\hat{x}$ are the real and generated images. We apply adversarial loss to both two-stage generations in Fig. 2.

**Perceptual Loss.** The perceptual loss is formulated as

$$\mathcal{L}_{\mathrm{P}} = \sum_{i} \eta_i \left\| \phi_i \left( \hat{x} \right) - \phi_i \left( x \right) \right\|_1, \tag{13}$$

where $\phi_i(\cdot)$ denotes the layer activation of a pre-trained VGG-19 [48] network. We only consider the high-level features of $conv_{4\_4}$ and $conv_{5\_4}$, allowing for variations of inpainted results, with scaling coefficients $\eta_i$ as $\frac{1}{4}$ and $\frac{1}{2}$.

**Overall Loss.** The overall loss of the generator is

$$\mathcal{L} = \mathcal{L}_{\mathrm{G}} + \gamma R_1 + \lambda \mathcal{L}_{\mathrm{P}}. \tag{14}$$

where $\gamma = 10$ and $\lambda = 0.1$.

## 4. Experiments

### 4.1. Datasets and Metrics

We conduct experiments on the Places365-Standard [78] and the CelebA-HQ [25] datasets at $512 \times 512$ resolution. Specifically, on the Places dataset, we use the 1.8 million and 36.5 thousand images from train and validation sets to train and evaluate our models, respectively. Images are randomly cropped or padded to $512 \times 512$ size during training while centrally cropped or padded for evaluation. For CelebA-HQ, train and validation splits are organized with 24,183 and 2,993 images. Though trained on $512 \times 512$ images, we show our model generalizes well to a larger resolution in the supplementary material.

In terms of the large hole setting, following [75], we opt for perceptual metrics including FID [21], P-IDS [75] and U-IDS [73] for evaluation. We suggest that it is inappropriate to use the pixel-wise L1 distance, PSNR and SSIM [58], since preliminary works [29, 47] have shown that these metrics correlate weakly with human perception regarding image quality, especially for the ill-posed large-scale image completion problem. Though LPIPS [73] is calculated in
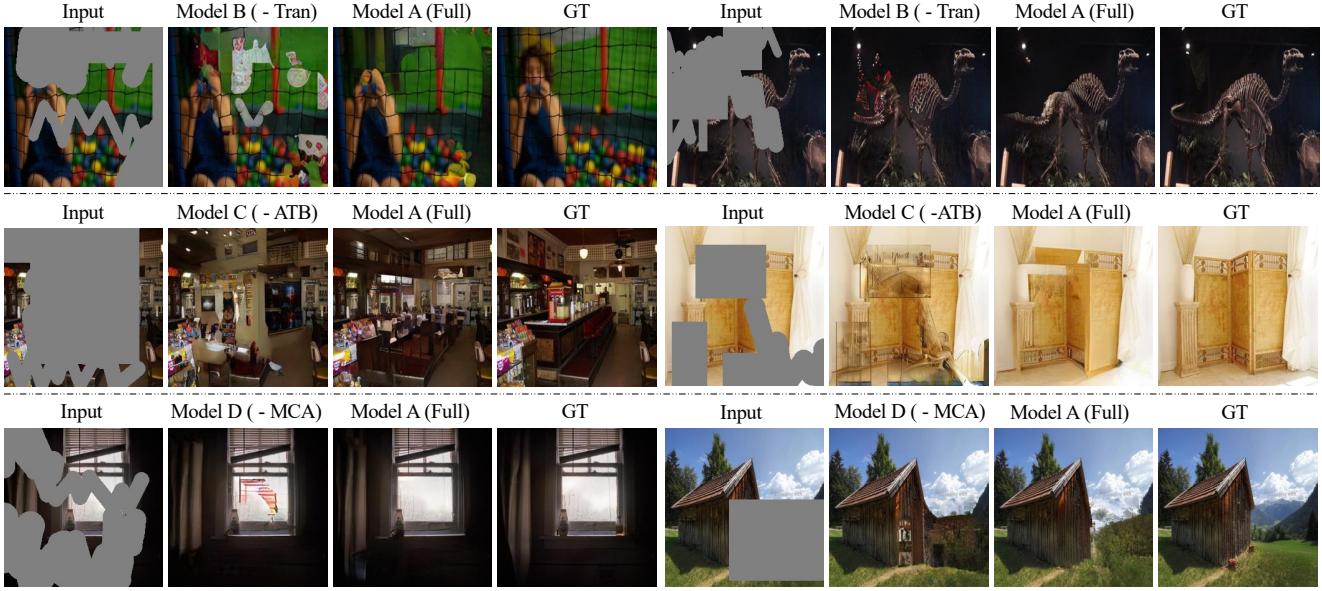
Figure 5. Visual examples for ablation study. Model A is our full model, while model B, C, D refer to models replacing transformers with convolutions, using the conventional transformer block and multi-head attention, respectively.

the deep feature space, the pixel-wise evaluation still greatly punishes diverse inpainting systems for large holes. Thus we only use it for reference in the supplementary material.

## 4.2. Implementation Details

In our framework, we set the numbers of convolution channels and FC dimensions to 180 for the head, body, and reconstruction modules. The block numbers and window sizes of 5-level transformer groups are $\{2, 3, 4, 3, 2\}$ and $\{8, 16, 16, 16, 8\}$, respectively. The last Conv-U-Net firstly downsamples the resolution to $\frac{1}{32}$ and then upsamples to the original size, where the numbers of convolution layers and channels at different scales are borrowed from [27]. The mapping network consists of 8 FC layers and the style manipulation module is implemented with convolutions followed with an AvgPool layer. Different from [55, 68, 77], our transformer architecture is *without* pre-training.

All experiments are carried out on 8 NVidia V100 GPUs. Following [75], we train our models for 50M images on Places365-Standard and 25M images on CelebA-HQ. The batch size is 32. We adopt an Adam optimizer with $\beta_1 = 0$ and $\beta_2 = 0.99$ and set the learning rate to $1 \times 10^{-3}$. The free-form mask is described in the supplementary file.

## 4.3. Ablation Study

In this section, we tease apart which components of our framework contribute most to the final performance. To enable a quick exploration, we only use 100K training images in Places [78] ($\approx 5.6\%$) at $256 \times 256$ resolution and train the models for 5M samples ($10\%$ of the full setting). We adopt the first 10K validation images for evaluation. The

| Type | Model | FID↓ | P-IDS (%)↑ | U-IDS(%)↑ |
|------|-------|------|------------|-----------|
| A | Full Model | **5.97** | **13.17** | **29.23** |
| B | - Tran. | 6.21 | 11.30 | 27.39 |
| C | - Adjusted Tran. Block | 6.36 | 12.30 | 28.05 |
| D | - MCA | 6.08 | 13.13 | 29.19 |
| E | - Style Mani. Module | 6.10 | 11.88 | 27.94 |
| F | - High-Res. Gen. | 6.32 | 12.57 | 28.21 |

Table 1. Ablation study of the framework components. "A" represents our full model. "B" replaces transformers with convolutions. "C" replaces our adjusted transformer block with the original design [52]. "D" means using the conventional attention strategy. "E" removes the noise style manipulation. "F" limits the output size of first-stage generation to $64 \times 64$.

quantitative comparison is shown in Table 1.

**Conv-Transformer Architecture.** We explore whether the long-range context relations modeled by transformers are useful for filling large holes. Replacing the transformer blocks with convolution blocks (model "B" in Table 1), we find an obvious performance drop on all metrics, especially on P-IDS and U-IDS, indicating that the inpainted images lose some fidelity. Moreover, we show some visual examples in Fig. 5. Compared to the fully convolutional network, our MAT takes advantage of distant context to reconstruct the structure of net and texture of dinosaur skeleton well, showing the effectiveness of long-range interactions.

**Adjusted Transformer Block.** In our framework, we develop a novel transformer block since the conventional design easily leads to unstable optimization, in which case we need to lower the learning rate of transformer body. As illustrated in Table 1, our design (model "A") obtains superior
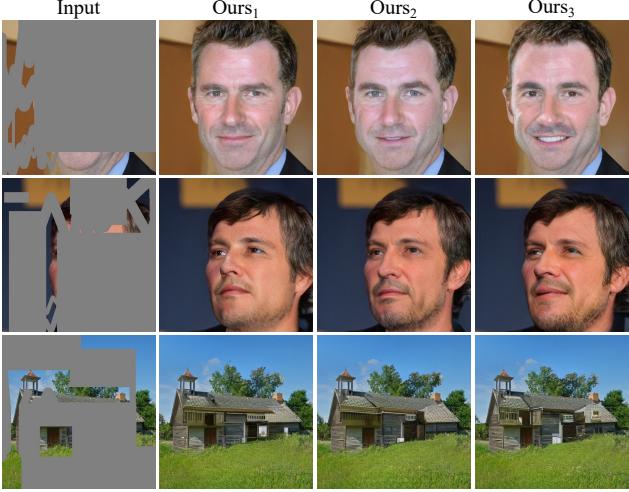
Figure 6. Visual examples with different style representations.



Figure 7. Failure cases of our method (MAT).

performance, 0.39 improvement on FID, than model "C" with the original transformer block. As illustrated in Fig. 5, we notice our design produces more visually appealing results, supporting high-quality image completion. Especially for the first example, even though the missing area is extremely large, our method can still recover a semantically consistent and visually realistic indoor scene.

**Multi-Head Contextual Attention.** To quickly fill the missing regions with realistic contents, we propose a multi-head contextual attention (MCA). To make a deeper understanding, we build a model without partial aggregation from valid tokens. The quantitative results are shown as model "D" in Table 1. It is noted that FID drops by 0.1 yet other metrics do not change too much. We suggest the proposed contextual attention is helpful for maintaining color consistency and reducing blurriness. As illustrated in Fig. 5, the model without MCA generates contents with incorrect colors for the first example, while producing blurry artifacts for the second example. Both the quantitative and qualitative results validate the power of our MCA.

**Style Manipulation Module.** To deal with large masks, apart from the conditional long-range interaction, we also introduce unconditional generation. To quantify the unconditional generative capability of our framework, we strip the noise style manipulation. From the results of model "E" in Table 1, we find a large gap on P-IDS and U-IDS, showing the modulation of stochastic noise styles further improves the naturalness of completed images.

**High Resolution in Reconstruction.** Due to quadratically increased computational complexity, existing works [55,68, 77] adopt transformers to synthesize low-resolution results, typically $32 \times 32$, for subsequent processing. By contrast, our MAT architecture takes advantage of its computational efficiency to enable high-resolution outputs in the recon-
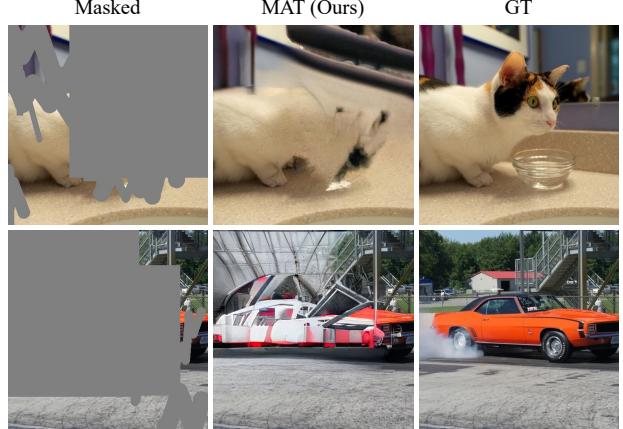
struction stage. As illustrated in Table 1, our full model "A" achieves significant improvement over model "F", demonstrating the importance of high-resolution prediction.

## 4.4. Comparison with State of the Arts

We compare the proposed MAT with a number of state-of-the-art approaches. For a fair comparison, we use publicly available models to test on the same masks. As illustrated in Table 2, MAT achieves state-of-the-art performance on both CelebA-HQ and Places. Especially, even if we only use a subset Places365-Standard (1.8M images) to train our model, much fewer than CoModGAN [75] (8M images) and Big LaMa [51] (4.5M images), MAT still yields promising results. Besides, our method is much more parameter-efficient than the second-best CoModGAN and transformer-based ICT [55]. As illustrated in Fig 8, compared to other methods, the proposed MAT restores more photo-realistic images with fewer artifacts. For example, our method successfully recovers visually pleasing flowers and regular building structures.

## 4.5. Pluralistic Generation

The inherent diversity of our framework mainly sources from the style manipulation. As shown in Fig. 6, style variants lead to different completions. From the first example in Fig. 6, we observe a change from a pursed smile to a toothy laugh. And the second example shows different face contours and appearances. As for the final one, we find different window and roof structures.

## 4.6. Limitations and Failure Cases

Trained without semantic annotations, MAT usually struggles when processing objects with a variety of shapes, e.g., running animals. As shown in Fig. 7, our method fails to recover the cat and car due to the lack of semantic context understanding. Also, limited by the downsampling and pre-

| Method | #Param. ×10⁶ | Places (512 × 512) | | | | | | CelebA-HQ (512 × 512) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Small Mask | | | Large Mask | | | Small Mask | | | Large Mask | | |
| | | FID↓ | P-IDS(%)↑ | U-IDS(%)↑ | FID↓ | P-IDS(%)↑ | U-IDS(%)↑ | FID↓ | P-IDS(%)↑ | U-IDS(%)↑ | FID↓ | P-IDS(%)↑ | U-IDS(%)↑ |
| **MAT (Ours)**[†] | 62 | 0.78 | 31.72 | 43.71 | 1.96 | 23.42 | 38.34 | 2.86 | 21.15 | 32.56 | 4.86 | 13.83 | 25.33 |
| **MAT (Ours)** | | 1.07 | 27.42 | 41.93 | 2.90 | 19.03 | 35.36 | | | | | | |
| CoModGAN [75][†] | 109 | 1.10 | 26.95 | 41.88 | 2.92 | 19.64 | 35.78 | 3.26 | 19.65 | 31.41 | 5.65 | 11.23 | 22.54 |
| LaMa [51][†] | 51/27 | 0.99 | 22.79 | 40.58 | 2.97 | 13.09 | 32.29 | 4.05 | 9.72 | 21.57 | 8.15 | 2.07 | 7.58 |
| ICT [55] | 150 | - | - | - | - | - | - | 6.28 | 2.24 | 9.99 | 12.84 | 0.13 | 0.58 |
| MADF [79] | 85 | 2.24 | 14.85 | 35.03 | 7.53 | 6.00 | 23.78 | 3.39 | 12.06 | 24.61 | 6.83 | 3.41 | 11.26 |
| AOT GAN [70] | 15 | 3.19 | 8.07 | 30.94 | 10.64 | 3.07 | 19.92 | 4.65 | 7.92 | 20.45 | 10.82 | 1.94 | 6.97 |
| HFill [65] | 3 | 7.94 | 3.98 | 23.60 | 28.92 | 1.24 | 11.24 | - | - | - | - | - | - |
| DeepFill v2 [67] | 4 | 3.02 | 9.17 | 32.56 | 9.27 | 4.01 | 21.32 | 10.11 | 3.11 | 9.52 | 24.42 | 0.17 | 0.42 |
| EdgeConnect [40] | 22 | 4.03 | 5.88 | 27.56 | 12.66 | 1.93 | 15.87 | 10.58 | 4.14 | 12.45 | 39.99 | 0.10 | 0.22 |

Table 2. Quantitative comparison on Places [78] and CelebA-HQ [25]. "†": Our Mat, CoModGAN [75] and LaMa [51] use 8M, 8M and 4.5M training images on Places, respectively, while our other model (without "†") is only trained on a subset (1.8M images). The LaMa models on Places and CelebA are different in size. The results of LPIPS and 256 × 256 CelebA are provided in the supplementary. The best and second best results are in red and blue.
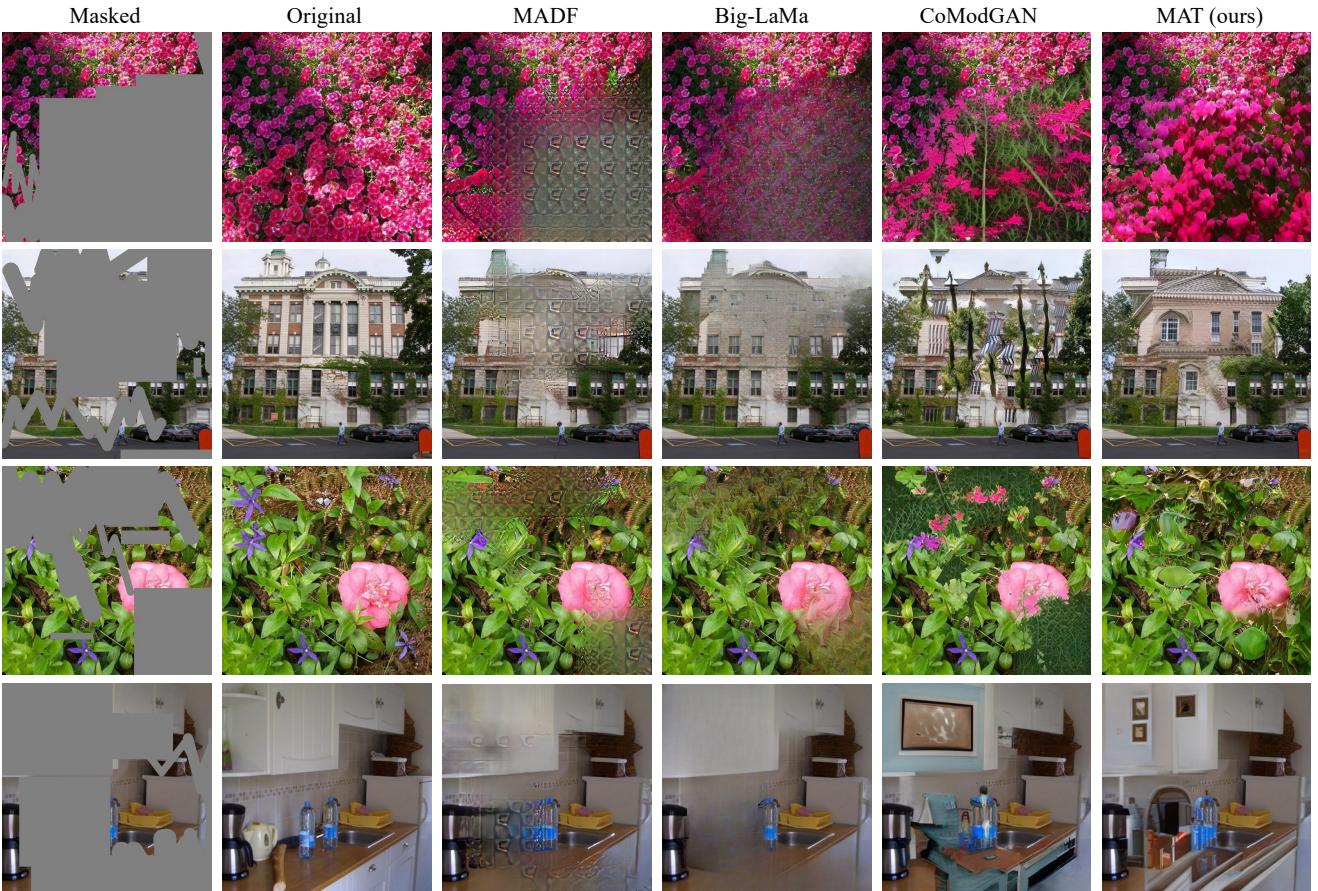


| Masked | Original | MADF | Big-LaMa | CoModGAN | MAT (ours) |

Figure 8. Qualitative comparison (512×512) with state-of-the-art methods. Our results are more visually realistic, containing more details.

defined window sizes in attention, we need to pad or resize an image to make its size a multiple of 512.

## 5. Conclusion

We have presented a mask-aware transformer (MAT) for pluralistic large hole image inpainting. Taking advantage of

the proposed adjusted transformer architecture and partial attention mechanism, the proposed MAT achieves state-of-the-art performance on multiple benchmarks. Also, we design a style modulation module to improve the diversity of generation. Extensive qualitative comparisons have demonstrated the superiority of our framework in terms of image quality and diversity.

# References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 2, 4

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

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

[4] Dana Berman, Shai Avidan, et al. Non-local image dehazing. In *CVPR*, pages 1674–1682, 2016. 2

[5] Marcelo Bertalmio. Strong-continuation, contrast-invariant inpainting with a third-order optimal pde. *TIP*, 15(7):1934–1938, 2006. 2

[6] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424, 2000. 2

[7] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *CVPR*, volume 2, pages 60–65. IEEE, 2005. 2

[8] Ting Chen, Mario Lucic, Neil Houlsby, and Sylvain Gelly. On self modulation for generative adversarial networks. In *ICLR*, 2018. 3, 5

[9] Donghyeon Cho, Jinsun Park, Tae-Hyun Oh, Yu-Wing Tai, and In So Kweon. Weakly-and self-supervised learning for content-aware deep image retargeting. In *ICCV*, pages 4558–4567, 2017. 1

[10] Antonio Criminisi, Patrick Perez, and Kentaro Toyama. Object removal by exemplar-based inpainting. In *CVPR*, volume 2, pages II–II. IEEE, 2003. 2

[11] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama. Region filling and object removal by exemplar-based image inpainting. *TIP*, 13(9):1200–1212, 2004. 2

[12] Ding Ding, Sundaresh Ram, and Jeffrey J Rodríguez. Image inpainting using nonlocal texture matching and nonlinear filtering. *TIP*, 28(4):1705–1719, 2018. 2

[13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 3, 12

[14] Mohamed-Jalal Fadili, J-L Starck, and Fionn Murtagh. Inpainting and zooming using sparse representations. *The Computer Journal*, 52(1):64–79, 2009. 2

[15] Mrinmoy Ghorai, Soumitra Samanta, Sekhar Mandal, and Bhabatosh Chanda. Multiple pyramids based image inpainting using local patch statistics and steering kernel feature. *TIP*, 28(11):5495–5509, 2019. 2

[16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NIPS*, 27, 2014. 2, 5

[17] Harald Grossauer. A combined pde and texture synthesis approach to inpainting. In *ECCV*, pages 214–224. Springer, 2004. 2

[18] Zongyu Guo, Zhibo Chen, Tao Yu, Jiale Chen, and Sen Liu. Progressive image inpainting with full-resolution residual network. In *ACMMM*, pages 2496–2504, 2019. 2

[19] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ToG*, 26(3):4–es, 2007. 2

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 4

[21] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NIPS*, 30, 2017. 5, 13

[22] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ToG*, 36(4):1–14, 2017. 2

[23] Youngjoo Jo and Jongyoul Park. Sc-fegan: Face editing generative adversarial network with user's sketch and color. In *ICCV*, pages 1745–1753, 2019. 1

[24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, pages 694–711. Springer, 2016. 5

[25] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 2, 5, 8, 13, 14

[26] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 1, 3, 5, 14

[27] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, pages 8110–8119, 2020. 3, 5, 6

[28] Olivier Le Meur, Josselin Gautier, and Christine Guillemot. Examplar-based inpainting based on local geometry. In *ICIP*, pages 3401–3404. IEEE, 2011. 2

[29] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, pages 4681–4690, 2017. 5

[30] Joo Ho Lee, Inchang Choi, and Min H Kim. Laplacian patch-based image synthesis. In *CVPR*, pages 2727–2735, 2016. 2

[31] Anat Levin, Assaf Zomet, and Yair Weiss. Learning how to inpaint from global image statistics. In *ICCV*, volume 1, pages 305–312, 2003. 2

[32] Jingyuan Li, Ning Wang, Lefei Zhang, Bo Du, and Dacheng Tao. Recurrent feature reasoning for image inpainting. In *CVPR*, pages 7760–7768, 2020. 2

[33] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *ECCV*, pages 85–100, 2018. 2

[34] Hongyu Liu, Bin Jiang, Yibing Song, Wei Huang, and Chao Yang. Rethinking image inpainting via a mutual encoder-decoder with feature equalizations. In *ECCV*, pages 725–741. Springer, 2020. 2

[35] Hongyu Liu, Bin Jiang, Yi Xiao, and Chao Yang. Coherent semantic attention for image inpainting. In *ICCV*, pages 4170–4179, 2019. 2

[36] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *ICCV*, 2021. 4

[37] Yuqing Ma, Xianglong Liu, Shihao Bai, Lei Wang, Aishan Liu, Dacheng Tao, and Edwin Hancock. Region-wise generative adversarial imageinpainting for large missing areas. *arXiv preprint arXiv:1909.12507*, 2019. 3

[38] Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Non-local sparse models for image restoration. In *ICCV*, pages 2272–2279. IEEE, 2009. 2

[39] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *ICML*, pages 3481–3490. PMLR, 2018. 5

[40] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z Qureshi, and Mehran Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. *arXiv preprint arXiv:1901.00212*, 2019. 2, 8, 13

[41] Seoung Wug Oh, Sungho Lee, Joon-Young Lee, and Seon Joo Kim. Onion-peel networks for deep video completion. In *ICCV*, pages 4403–4412, 2019. 2

[42] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, pages 2536–2544, 2016. 2

[43] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *arXiv preprint arXiv:2108.08810*, 2021. 3

[44] Yurui Ren, Xiaoming Yu, Ruonan Zhang, Thomas H Li, Shan Liu, and Ge Li. Structureflow: Image inpainting via structure-aware appearance flow. In *ICCV*, pages 181–190, 2019. 2

[45] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2

[46] Andrew Slavin Ross and Finale Doshi-Velez. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *AAAI*, 2018. 5

[47] Mehdi SM Sajjadi, Bernhard Scholkopf, and Michael Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *ICCV*, pages 4491–4500, 2017. 5

[48] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 5

[49] Yuhang Song, Chao Yang, Yeji Shen, Peng Wang, Qin Huang, and C-C Jay Kuo. Spg-net: Segmentation prediction and guidance network for image inpainting. *arXiv preprint arXiv:1805.03356*, 2018. 2

[50] Jian Sun, Lu Yuan, Jiaya Jia, and Heung-Yeung Shum. Image completion with structure propagation. In *ACM SIGGRAPH 2005 Papers*, pages 861–868. 2005. 2

[51] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. *arXiv preprint arXiv:2109.07161*, 2021. 3, 7, 8, 13

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008, 2017. 2, 4, 6

[53] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Bringing old photos back to life. In *CVPR*, pages 2747–2757, 2020. 2

[54] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Old photo restoration via deep latent space translation. *arXiv preprint arXiv:2009.07047*, 2020. 2

[55] Ziyu Wan, Jingbo Zhang, Dongdong Chen, and Jing Liao. High-fidelity pluralistic image completion with transformers. *arXiv preprint arXiv:2103.14031*, 2021. 2, 3, 6, 7, 8, 13

[56] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, pages 7794–7803, 2018. 2

[57] Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative multi-column convolutional neural networks. *NIPS*, 2018. 2

[58] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *TIP*, 13(4):600–612, 2004. 5

[59] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, pages 22–31, 2021. 4

[60] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *arXiv preprint arXiv:2106.14881*, 2021. 3

[61] Chaohao Xie, Shaohui Liu, Chao Li, Ming-Ming Cheng, Wangmeng Zuo, Xiao Liu, Shilei Wen, and Errui Ding. Image inpainting with learnable bidirectional attention maps. In *ICCV*, pages 8858–8867, 2019. 2

[62] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *NeurIPS*, 34, 2021. 4

[63] Wei Xiong, Jiahui Yu, Zhe Lin, Jimei Yang, Xin Lu, Connelly Barnes, and Jiebo Luo. Foreground-aware image inpainting. In *CVPR*, pages 5840–5848, 2019. 2

[64] Zhaoyi Yan, Xiaoming Li, Mu Li, Wangmeng Zuo, and Shiguang Shan. Shift-net: Image inpainting via deep feature rearrangement. In *ECCV*, pages 1–17, 2018. 2

[65] Zili Yi, Qiang Tang, Shekoofeh Azizi, Daesik Jang, and Zhan Xu. Contextual residual aggregation for ultra high-resolution image inpainting. In *CVPR*, pages 7508–7517, 2020. 2, 8, 13

[66] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *CVPR*, pages 5505–5514, 2018. 2

[67] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *ICCV*, pages 4471–4480, 2019. 2, 8, 12, 13

[68] Yingchen Yu, Fangneng Zhan, Rongliang Wu, Jianxiong Pan, Kaiwen Cui, Shijian Lu, Feiying Ma, Xuansong Xie, and Chunyan Miao. Diverse image inpainting with bidirectional and autoregressive transformers. *arXiv preprint arXiv:2104.12335*, 2021. 2, 3, 6, 7

[69] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for high-quality image inpainting. In *CVPR*, pages 1486–1494, 2019. 2

[70] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Aggregated contextual transformations for high-resolution image inpainting. *arXiv preprint arXiv:2104.01431*, 2021. 8, 13

[71] Yu Zeng, Zhe Lin, Jimei Yang, Jianming Zhang, Eli Shechtman, and Huchuan Lu. High-resolution image inpainting with iterative confidence feedback and guided upsampling. In *ECCV*, pages 1–17. Springer, 2020. 2

[72] Haoran Zhang, Zhenzhen Hu, Changzhi Luo, Wangmeng Zuo, and Meng Wang. Semantic image inpainting with progressive generative networks. In *ACMMM*, pages 1939–1947, 2018. 2

[73] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. 5, 13

[74] Lei Zhao, Qihang Mo, Sihuan Lin, Zhizhong Wang, Zhiwen Zuo, Haibo Chen, Wei Xing, and Dongming Lu. Uctgan: Diverse image inpainting based on unsupervised cross-space translation. In *CVPR*, pages 5741–5750, 2020. 3

[75] Shengyu Zhao, Jonathan Cui, Yilun Sheng, Yue Dong, Xiao Liang, I Eric, Chao Chang, and Yan Xu. Large scale image completion via co-modulated generative adversarial networks. In *ICLR*, 2020. 3, 5, 6, 7, 8, 12, 13

[76] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Pluralistic image completion. In *CVPR*, pages 1438–1447, 2019. 3

[77] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Tfill: Image completion via a transformer-based architecture. *arXiv preprint arXiv:2104.00845*, 2021. 2, 3, 6, 7

[78] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *PAMI*, 40(6):1452–1464, 2017. 1, 2, 5, 6, 8, 12, 13, 14

[79] Manyu Zhu, Dongliang He, Xin Li, Chao Li, Fu Li, Xiao Liu, Errui Ding, and Zhaoxiang Zhang. Image inpainting by end-to-end cascaded refinement with mask awareness. *TIP*, 30:4855–4866, 2021. 8, 13

# MAT: Mask-Aware Transformer for Large Hole Image Inpainting
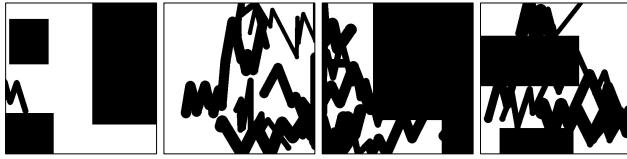## (Supplementary Material)



Figure A.1. Examples of free-form masks ($512 \times 512$). Visible and invisible pixels are in white and black colors.

## A. Network Architecture

As illustrated in Sec. 3.1, the proposed MAT is a two-stage framework, where the first stage consists of a convolutional head, a transformer body and a convolutional reconstruction tail while the second stage is a Conv-U-Net. And the discriminator follows the design of CoModGAN [75].

Given an $H \times W$ input, the head first applies a convolution to change the number of channels from 4 (image 3 + mask 1) to 180 and then adopts three strided convolutions (stride = 2) to downsample the feature size to $\frac{H}{8} \times \frac{W}{8}$. The feature is transformed to tokens as input to the transformer body. The body is composed of five stages of transformer blocks, where the block numbers are $\{2, 3, 4, 3, 2\}$ and the corresponding feature sizes are $\{\frac{H}{8} \times \frac{W}{8}, \frac{H}{16} \times \frac{W}{16}, \frac{H}{32} \times \frac{W}{32}, \frac{H}{16} \times \frac{W}{16}, \frac{H}{8} \times \frac{W}{8}\}$. The downsampling and upsampling are realized by convolutions. The detailed structure of a transformer block is shown in Sec. 3.3. Then the output tokens from the body are converted to a 2D feature, passed to the reconstruction tail. The convolutional tail upsamples the feature size from $\frac{H}{8} \times \frac{W}{8}$ to $H \times W$ and generates a completed image, during which style modulation is applied to all layers to enable pluralistic generation.

The second-stage Conv-U-Net takes in the coarse prediction and the input mask for subsequent high-fidelity detail rendering. It first downsamples the feature size to $\frac{H}{32} \times \frac{W}{32}$ and then upsamples the size back to $H \times W$. Shortcut connections are adopted at each resolution. The number of convolution channels in the encoder starts from 64 and is doubled after each downsampling, with a maximum of 512, while the decoder uses a symmetrical setting. Besides, all decoding layers are modulated by the image-conditional and noise-unconditional style representations.

## B. Free-Form Mask Sampling and Statistics

Referring to DeepFill v2 [67], we sample rectangles and brush strokes with random sizes, shapes and locations to
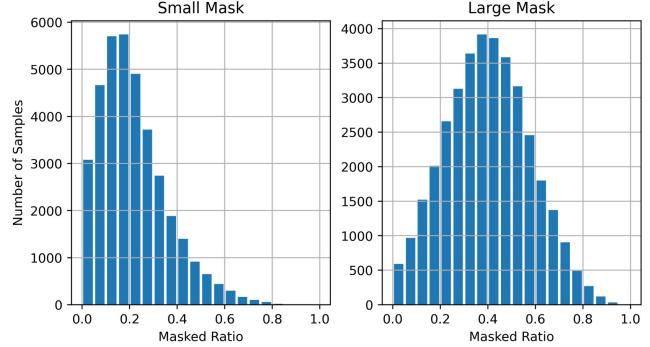


Figure A.2. Small and large mask ($512 \times 512$) statistics on the Places Val set [78]. The are totally 36500 masks.

generate free-form masks. During training, we use a large mask sampling strategy. The number of up to full-size or half-size rectangles is uniformly sampled within $[0, 3]$ or $[0, 5]$. The number of strokes is randomly sampled within $[0, 9]$, with a random brush width within $[12, 48]$ and vertex number within $[4, 18]$. During testing, apart from the large mask setup, we also introduce a small mask sampling strategy, where the number of up to full-size or half-size rectangles is within $[0, 2]$ or $[0, 3]$ and the number of strokes is within $[0, 4]$, while other settings remain unchanged. Note that our model is trained on large masks and is evaluated on both small and large mask settings. As shown in Fig. A.2, we present the mask statistics on the Places Val set [78] that is used for evaluation. It is observed that large masks are very aggressive and diverse.

## C. Tokenization

As described in Sec. A, we adopt a stack of convolutions (the convolutional head) to extract tokens for the transformer body, which is specially tailored to the inpainting problem. Compared to the linear projection of ViT [13], our design owns two merits. First, stacked convolutions can gradually fill the holes, producing more effective tokens. Second, the multi-scale downsampled features can be passed to the decoder through shortcut connections, improving the optimization. As illustrated in Table C.1 and Fig. C.3, stacked convolutions obtain obviously superior results. The model using linear projection is more likely to generate unpleasing artifacts and fail to borrow surrounding textures to fill the holes, while our MAT successfully recovers high-fidelity contents. Both the quantitative and qualitative results demonstrate the effectiveness of our MAT.

| Model | FID↓ | P-IDS (%)↑ | U-IDS(%)↑ |
|---|---|---|---|
| Stacked Conv. (Ours) | **5.97** | **13.17** | **29.23** |
| Linear Projection | 10.54 | 5.77 | 20.86 |

Table C.1. Quantitative comparison between linear projection and stacked convolutions for token extraction. We use the same training setting as the ablation study (Sec. 4.3).
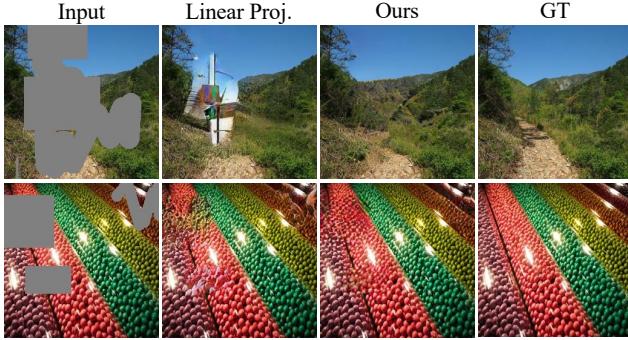


Figure C.3. Qualitative comparison between linear projection and stacked convolutions (ours) for tokenization.

| Model | Feature Dim. | Block Num. | Window Size | FID↓ |
|---|---|---|---|---|
| Ours | 180 | {2, 3, 4, 3, 2} | {8, 16, 16, 16, 8} | **5.97** |
| V1 | 90 | {2, 3, 4, 3, 2} | {8, 16, 16, 16, 8} | 6.28 |
| V2 | 180 | {1, 1, 2, 1, 1} | {8, 16, 16, 16, 8} | 6.18 |
| V3 | 180 | {2, 3, 4, 3, 2} | {8, 8, 8, 8, 8} | 6.09 |

Table D.2. Ablation study on model configuration.

## D. Model Configuration

Following the same experimental setting as ablation study, we explore several model variants in terms of feature width, block number and window size of the transformer body, leaving Conv-U-Net unchanged. The results are shown in the Table D.2. The performance is positively correlated to the model capacity and attention range.

## E. CelebA-HQ 256 × 256 Results

We provide the quantitative results on $256 \times 256$ CelebA-HQ [25]. As illustrated in Table F.3, our MAT yields significant improvements on FID [21], P-IDS [75] and U-IDS [73] metrics over other methods.

## F. LPIPS Results

As discussed in Sec. 4.1, LPIPS [73] is not an appropriate measure for large mask inpainting, especially for pluralistic generation systems, since there could be numerous plausible solutions to fill the holes. Therefore, we provide the LPIPS results only for reference. As shown in Table F.4, our method achieves superior or comparable performance on the CelebA-HQ [25] and Places [78] datasets. *Note that we only use 22.5% of full data to train our Places model.*

| Method | Small Mask | | | Large Mask | | |
|---|---|---|---|---|---|---|
| | FID↓ | P-IDS↑ | U-IDS↑ | FID↓ | P-IDS↑ | U-IDS↑ |
| MAT (Ours) | **2.94** | **20.88** | **32.01** | **5.16** | **13.90** | **25.13** |
| LaMa [51] | 3.98 | 8.82 | 22.57 | 8.75 | 2.34 | 8.77 |
| ICT [55] | 5.24 | 4.51 | 17.39 | 10.92 | 0.90 | 5.23 |
| MADF [79] | 10.43 | 6.25 | 14.62 | 23.59 | 0.50 | 1.44 |
| AOT GAN [70] | 9.64 | 5.61 | 14.62 | 22.91 | 0.47 | 1.65 |
| DeepFill v2 [67] | 5.69 | 6.62 | 16.82 | 13.23 | 0.84 | 2.62 |
| EdgeConnect [40] | 5.24 | 5.61 | 15.65 | 12.16 | 0.84 | 2.31 |

Table F.3. Quantitative results on CelebA-HQ at $256 \times 256$ size. The results of P-IDS and U-IDS are shown in percentage (%).

| Method | #Param. ×10⁶ | CelebA-HQ | | Places | |
|---|---|---|---|---|---|
| | | Small | Large | Small | Large |
| MAT (Ours) | 60 | **0.065** | **0.125** | 0.099 | 0.189 |
| CoModGAN [75]† | 109 | 0.073 | 0.140 | 0.101 | 0.192 |
| LaMa [51]† | 27/51 | 0.075 | 0.143 | **0.086** | **0.166** |
| ICT [55] | 150 | 0.105 | 0.195 | - | - |
| MADF [79] | 85 | 0.068 | 0.130 | 0.095 | 0.181 |
| AOT GAN [70] | 15 | 0.074 | 0.145 | 0.101 | 0.195 |
| HFill [65] | 3 | - | - | 0.148 | 0.284 |
| DeepFill v2 [67] | 4 | 0.117 | 0.221 | 0.113 | 0.213 |
| EdgeConnect [40] | 22 | 0.101 | 0.208 | 0.114 | 0.275 |

Table F.4. LPIPS [73] comparison on $512 \times 512$ CelebA-HQ [25] and Places [78] datasets. "†": CoModGAN [75] and LaMa [51] use 8M and 4.5M Places images to train their models, while our model is only trained on Places365-Standard (1.8M images). The LaMa models on CelebA-HQ and Places are different in size.

## G. Generalization to A Higher Resolution

Though trained on $512 \times 512$ images, our model generalizes well to larger resolutions. For example, we transfer our model and Big LaMa [51] trained at $512 \times 512$ resolution to $1024 \times 1024$. Compared to Big LaMa (FID↓ 6.31, PIDS↑ 4.98%), our model (FID↓ 5.83, P-IDS↑ 9.51%) obtains superior results on Places under the large mask setting. We suggest that maintaining a resolution consistency during training and testing yields better visual quality.

## H. Diversity-Fidelity Tradeoff

To evaluate the fidelity and diversity, apart from FID (depending on both diversity and fidelity), we also follow [**?**, **?**] to use Improved Precision and Recall to separately measure sample fidelity (precision) and diversity (recall). As shown in Table H.5, our method obtains better FID, higher recall yet slightly lower precision compared to CoModGAN on Places. It is noted that we use much less training data.

| Method | Training Data | FID↓ | Precision↑ | Recall↑ |
|---|---|---|---|---|
| **MAT (Ours)** | **1.8M** | **2.90** | 0.925 | **0.951** |
| CoModGAN | 8M | 2.92 | **0.929** | 0.942 |

Table H.5. Precision and Recall results of our MAT and Co-ModGAN on Places.

## I. Additional Qualitative Results

We present more visual comparisons on the Places [78] dataset between our MAT and other state-of-the-art methods. As shown in Fig J.4 and Fig J.5, our method generates more photo-realistic results with few artifacts, manifesting the effectiveness of MAT. Due to potential copyright issues with CelebA-HQ [25], we do not provide visual comparisons on this dataset. If necessary, you can process CelebA-HQ images with the provided code and model, or contact the authors.

## J. Licenses of Face Images

All face images used in the paper and supplementary material are from the FFHQ [26] dataset. Here we provide the detailed information on source and license.

- Face image in Fig.1 of main paper, source: https://www.flickr.com/photos/v63/5876049365/, license: CC BY-NC 2.0 (https://creativecommons.org/licenses/by-nc/2.0/).

- Face image in Fig.2 of main paper, source: https://www.flickr.com/photos/tbisaacs/4089001580/, license: CC BY 2.0 (https://creativecommons.org/licenses/by/2.0/).

- The first face image in Fig.6 of main paper, source: https://www.flickr.com/photos/southlanarkshirecouncil/8341157963/, license: CC BY-NC 2.0 (https://creativecommons.org/licenses/by-nc/2.0/).

- The second face image in Fig.6 of main paper, source: https://www.flickr.com/photos/afge/34804627253/, license: CC BY 2.0 (https://creativecommons.org/licenses/by/2.0/).

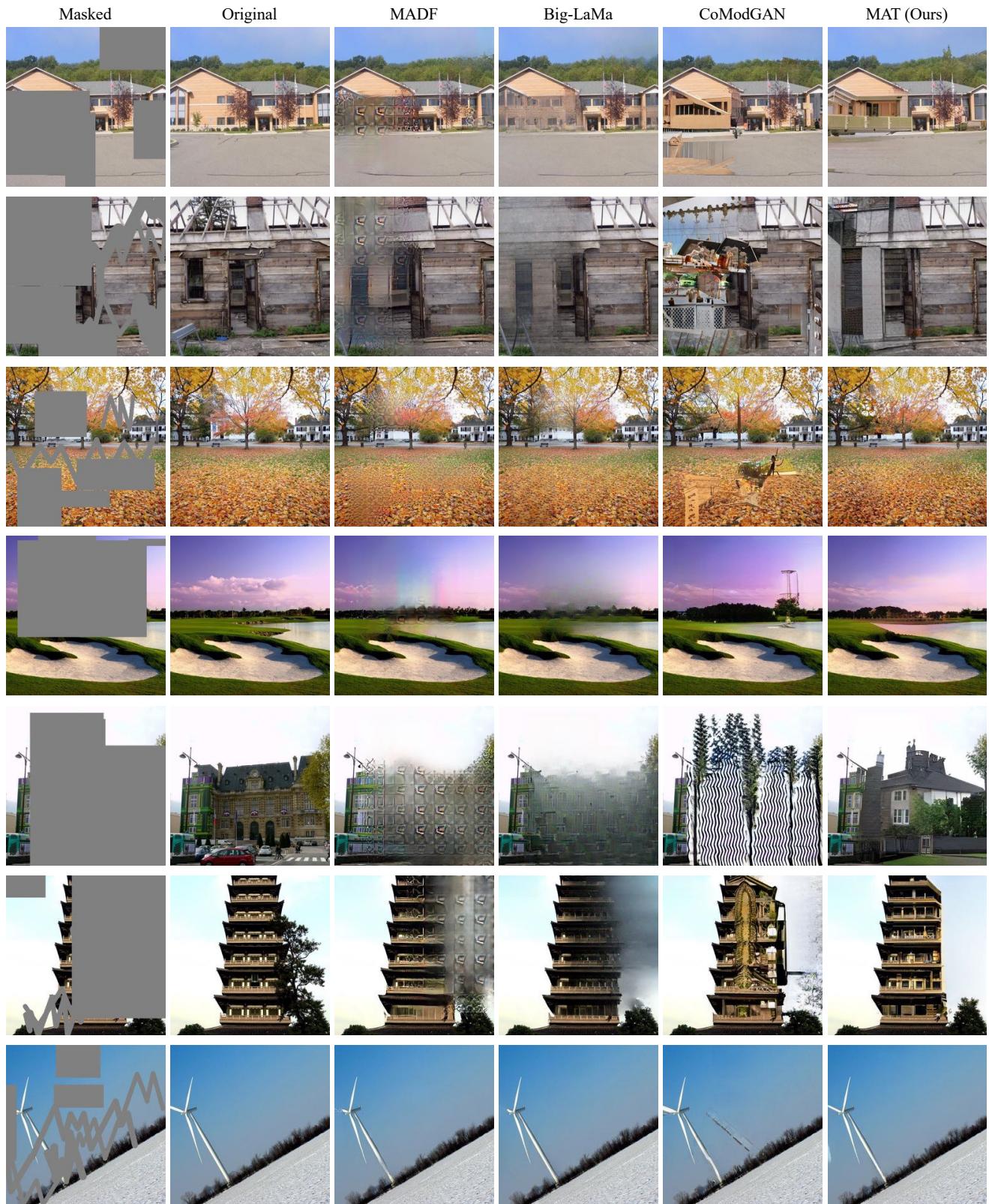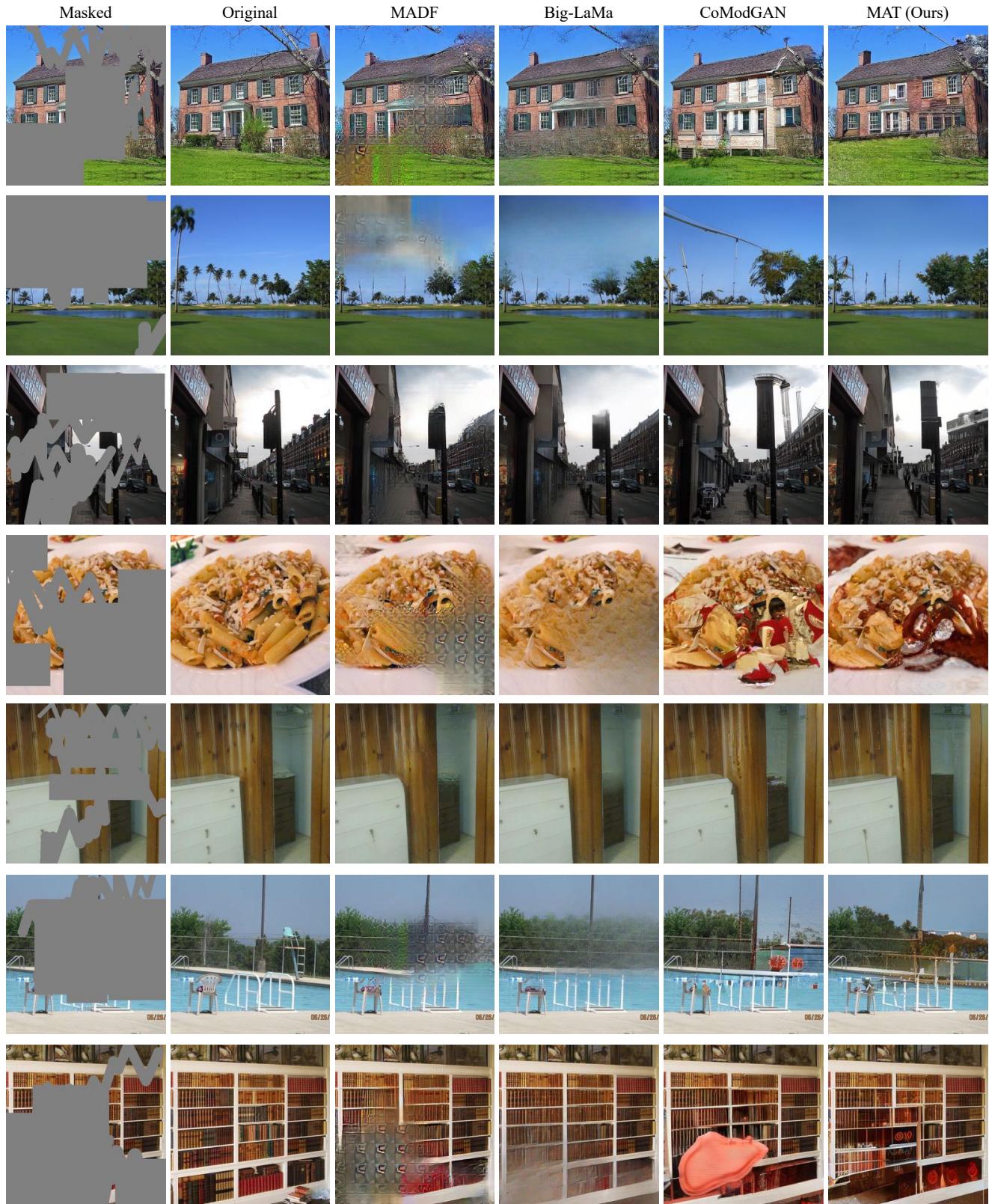| Masked | Original | MADF | Big-LaMa | CoModGAN | MAT (Ours) |
|--------|----------|------|----------|----------|------------|

Figure J.4. Qualitative comparison ($512 \times 512$) with state-of-the-art methods on the Places dataset. Zoom in for a better view.

| Masked | Original | MADF | Big-LaMa | CoModGAN | MAT (Ours) |

Figure J.5. Qualitative comparison ($512 \times 512$) with state-of-the-art methods on the Places dataset. Zoom in for a better view.