# PixelHacker: Image Inpainting with Structural and Semantic Consistency

Ziyang Xu[1,◇]    Kangsheng Duan[1,◇]    Xiaolei Shen[2]    Zhifeng Ding[2]    Wenyu Liu[1]

Xiaohu Ruan[2]    Xiaoxin Chen[2]    Xinggang Wang[1,✉]

[1]Huazhong University of Science and Technology    [2]VIVO AI Lab

**(a) Example 1**

**(b) Example 2**

Figure 1. **Qualitative comparison of our PixelHacker vs. other state-of-the-art methods.** While other methods struggle with complex structure and semantics, PixelHacker demonstrates remarkable semantic consistency and superior preservation of texture and color, as shown in Example 1. Moreover, PixelHacker maintains contextual structural consistency and logical coherence, as shown in Example 2.

## Abstract

*Image inpainting is a fundamental research area between image editing and image generation. Recent state-of-the-art (SOTA) methods have explored novel attention mechanisms, lightweight architectures, and context-aware modeling, demonstrating impressive performance. However, they often struggle with complex structure (e.g., texture, shape, spatial relations) and semantics (e.g., color consistency, object restoration, and logical correctness), leading to artifacts and inappropriate generation. To address this challenge, we design a simple yet effective inpainting paradigm called latent categories guidance, and further propose a diffusion-based model named PixelHacker. Specifically, we first construct a large dataset containing 14 million image-mask pairs by annotating foreground and background (potential 116 and 21 categories, respectively). Then, we encode potential foreground and background representations separately through two fixed-size embeddings, and intermittently inject these features into the denoising process via*

---

◇ Intern of VIVO AI Lab.
✉ Corresponding Author:xgwang@hust.edu.cn

1

*linear attention. Finally, by pre-training on our dataset and fine-tuning on open-source benchmarks, we obtain PixelHacker. Extensive experiments show that PixelHacker comprehensively outperforms the SOTA on a wide range of datasets (Places2, CelebA-HQ, and FFHQ) and exhibits remarkable consistency in both structure and semantics. Project page at https://hustvl.github.io/PixelHacker.*

## 1. Introduction

Image inpainting, as a fundamental research in computer vision, has been widely applied in image editing and object removal [22, 28]. The goal is to generate visually plausible content within the masked region by leveraging contextual pixel information from a given image-mask pair [15].

With the widespread application of inpainting techniques [2], there is an increasing demand for higher visual quality in generated content, particularly in terms of structural and semantic consistency. Structural consistency ensures that the generated content exhibits natural texture transitions, shape coherence, and physically plausible spatial relationships with the surrounding pixels. Semantic consistency requires smooth color blending, faithful reconstruction of object features when the mask partially covers a target, and logical coherence within the contextual pixel environment. Recent state-of-the- art (SOTA) approaches have focused on novel attention mechanisms, lightweight architectures, and improved contextual awareness, demonstrating impressive inpainting capabilities [14, 22, 25]. However, they often struggle with complex structures and semantics, as shown in Fig. 1. For instance, in Fig. 1 (a), when dealing with a texture-rich tree trunk and a sunlit water surface, these methods produce semantic inconsistencies (e.g., generating irrelevant objects on the tree trunk), structural distortions, color discrepancies, and blurriness. Similarly, in Fig. 1(b), when handling an image with multiple foreground elements (ground and people), middle-ground elements (railings and tree trunk), and background elements (forest and pathway), these methods fail to maintain structural consistency (e.g., discontinuous railings behind people), logical coherence, semantic relevance (e.g., generating objects unrelated to the scene), and also blurriness.

To address these challenges, we design a simple yet effective inpainting paradigm, termed Latent Categories Guidance (LCG), and further propose a diffusion-based model named PixelHacker. Specifically, we first construct a large-scale dataset containing 14 million image-mask pairs by annotating "foreground" and "background" (potential 116 and 21 categories, respectively). Then, we employ two fixed-size embeddings to encode the latent foreground and background representations separately. By intermittently injecting these latent features into the denoising process via linear attention, we effectively guide the generation process toward structural and semantic interaction. Finally, by

pre-training on our dataset and fine-tuning on open-source benchmarks, we obtain PixelHacker, a model that exhibits surprisingly well structural and semantics consistency. As illustrated in Fig. 1, PixelHacker generates visually coherent content with natural texture transitions, smooth color blending, and logical consistency, significantly outperforming existing SOTA methods. Notably, during the entire training pipeline, we do not require explicit supervision on the exact object category within the masked region (e.g., distinguishing between a person, car, or chair). Instead, we represent each masked object as either foreground or background, which encourages focusing on foreground-background semantics, implicitly compressing diverse object representations while reducing implementation costs. Extensive experiments demonstrate that PixelHacker consistently outperforms SOTA methods across various benchmarks (Places2 [38], CelebA-HQ [26], and FFHQ [11]), exhibiting superior structural and semantic consistency.

In summary, we make the following contributions:

(1) We introduce a simple yet effective inpainting paradigm, Latent Categories Guidance (LCG). LCG first constructs image-mask pairs based on "foreground" and "background" label. During training, LCG uses two fixed-size embeddings to encode the latent foreground and background features, which are intermittently injected into the denoising process via linear attention, leading to a model that inpaints with structural and semantic consistency.

(2) We propose PixelHacker, a diffusion-based inpainting model trained with LCG on 14M image-mask pairs and fine-tuning on open-source benchmarks. We demonstrate that PixelHacker achieves excellent structural and semantic consistency, outperforming various SOTA methods.

(3) Extensive experiments show that PixelHacker consistently outperforms SOTA approaches across multiple benchmarks (Places2, CelebA-HQ, and FFHQ).

## 2. Related Work

**Image Inpainting.** GAN-based methods [3, 4, 10, 11, 22, 31, 36, 37] typically introduce a large number of random masks during training to adversarially optimize the generator, enabling effective utilization of contextual texture features and achieving reasonable structural consistency. However, these methods often fail to infer semantic information from the masked region, leading to semantically implausible objects, as exemplified by MI-GAN in Fig. 1(b). Convolution-based approaches [14, 15, 25, 34, 35] focus primarily on contextual coherence. While they expand the receptive field and mitigate the interference of irrelevant features, they lack explicit attention to semantic consistency, often resulting in abrupt color transitions and inappropriate feature reconstructions, as observed in the results of LaMa and MAT in Fig. 1(a). Diffusion-based methods benefit from high-quality pretrained text-to-image (T2I)
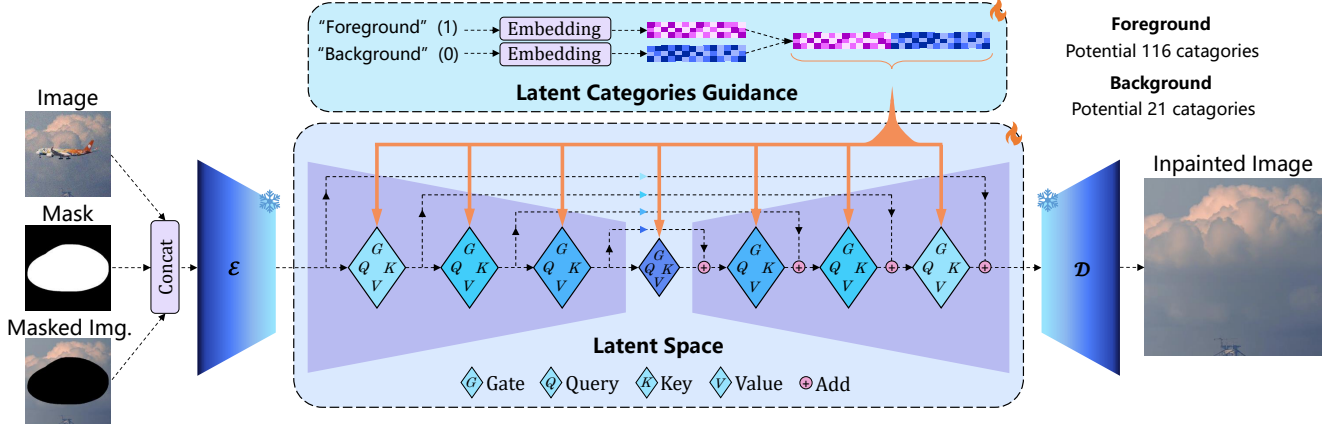
2

Figure 2. **Overall pipeline of our PixelHacker.** PixelHacker builds upon the latent diffusion architecture by introducing two fixed-size LCG embeddings to separately encode latent foreground and background features. We employ linear attention to inject these latent features into the denoising process, enabling intermittent structural and semantic multiple interactions. This design encourages the model to learn a data distribution that is both structurally and semantically consistent. We elaborate on the interaction details in Fig. 4 and Sec. 3.3.

models, allowing them to inject new semantic information via text prompts [19, 21, 28, 30, 40]. However, text prompts act as a double-edged sword in inpainting. On one hand, low-quality prompts degrade inpainting performance; on the other hand, introducing multiple text encoders to enhance prompt quality [19] results in computational redundancy. While text encoders improve semantic-aware generation, they often compromise structural consistency. For instance, as shown in Fig. 1(b), the results of SD-Inpainting [21], SDXL-Inpainting [19], and PowerPaint [40] exhibit discontinuous railings behind the people, which disrupts the original scene structure. Unlike these methods, Pixel-Hacker builds upon the latent diffusion architecture by introducing two fixed-size embeddings to separately encode latent foreground and background features. We employ linear attention to inject these latent features into the denoising process, enabling intermittent structural and semantic multiple interactions. This design encourages the model to learn a data distribution that is both structurally and semantically consistent. Similar to other Latent Diffusion Models (LDMs), we utilize Classifier-Free Guidance (CFG) inference [9]. Ultimately, PixelHacker achieves SOTA performance across multiple benchmarks [11, 26, 38].

## 3. Method

### 3.1. Overall Pipeline

The overall pipeline of our PixelHacker is illustrated in Fig. 2. First, following [21], PixelHacker takes a noised image, clean mask, and clean masked image as inputs, concatenates them, and feeds them into the encoder of a VAE [12], which transforms features from the pixel space to the latent space. Next, LCG constructs image-mask pairs based
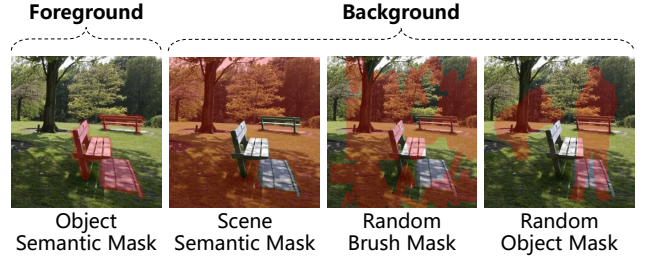


Figure 3. **Illustration of various masks we use to construct Latent Categories Guidance (LCG).** We assign object semantic masks to the foreground embedding and the other three masks to the background embedding. Details refer to Sec. 3.2.

on "foreground" and "background" label (potential 116 and 21 categories, respectively). Two fixed-size embeddings are then used to encode the latent foreground and background features separately. Then, in the latent space, we apply linear attention throughout both the downsampling and upsampling processes. By injecting the embeddings into the linear attention, we achieve intermittent structural and semantic consistency interactions. Finally, the encoded features in the latent space are passed through the decoder of the VAE to reconstruct the inpainted image.

### 3.2. Construction of Latent Categories Guidance

Previous studies have demonstrated that inpainting paradigms trained with random brush masks can achieve remarkable performance [18, 21, 34]. However, this strategy does not leverage segmentation masks, which contain rich semantic information. On the other hand, while we enumerate 116 foreground categories and 21 background categories when constructing the dataset, the fixed number
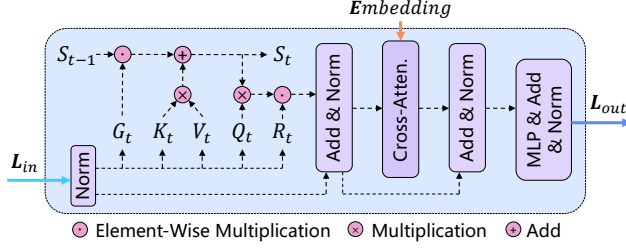
3

Figure 4. **The single interaction process between LCG embeddings and latent features.** We elaborate on the details in Sec. 3.3. Throughout the pipeline, multiple interactions are performed in a sequential manner, guiding the model to learn foreground semantics, background semantics, and contextual structures.

of categories inevitably limits the model's generalization and expansion capability for novel categories [6].

To integrate the advantages of both random masks and segmentation masks, while avoiding explicit reliance on precise category labels within foreground and background masks, we design and construct Latent Categories Guidance (LCG) to inject only two broad categories—"foreground" and "background"—into a conditioned diffusion model.

Specifically, we categorize masks into four types, as shown in Fig. 3, and assign them to either the foreground or background embedding. Object semantic masks are assigned to the foreground embedding, as illustrated in Fig. 3(a), where benches are masked. The goal is to enable the model to reconstruct foreground objects that are semantically aligned with the masked region [1], utilizing contextual background information. Scene semantic masks are assigned to the background embedding to enhance the model's ability to reconstruct background semantics. Similar to previous works [19, 21, 28, 30, 40], we also incorporate random brush masks. However, we assign them only to the background embedding, guiding the model to focus on structural information from the surrounding context. Finally, to prevent the model from overfitting object semantic masks by binding them too rigidly to foreground objects and ignoring embedding conditions, we assign random object masks to the background embedding as a regularization constraint. Overall, for the foreground mask $\mathcal{M}_{fg}$, we directly adopt the object semantic mask, i.e., $\mathcal{M}_{fg} = \mathcal{M}_{obj}$. For the background mask, we start with the scene semantic mask $\mathcal{M}scene$ and further incorporate random brush mask $\mathcal{M}rand$ and random object mask $\mathcal{M}obj'$ with probabilities $\mathcal{P}rand$ and $\mathcal{P}_{obj}$, respectively, as follows:

$$\mathcal{M}_{bg} = \mathcal{M}_{scene} + \mathcal{M}_{rand}\mathcal{P}_{rand} + \mathcal{M}'_{obj}\mathcal{P}_{obj}. \quad (1)$$

With this mask assignment strategy, our model learns to capture foreground object distributions under the "foreground" embedding condition, and background semantics and contextual structures under the "background" embed-

ding condition. During training, we do not explicitly provide category labels as textual prompts to guide generation. Instead, the model learns foreground semantics, background semantics, and contextual structures through learnable embedding weights, ultimately enabling the distribution of both categories to be injected into a single model. Despite adopting a different training paradigm, our approach remains compatible with other LDM inference schemes [21]. By applying Classifier-Free Guidance (CFG) [9], we can effectively merge the learned embeddings of the two categories to generate coherent outputs.

### 3.3. Structure & Semantic Consistency Interaction

As shown in Fig. 2, we perform multiple interactions between LCG embeddings and latent features during the denoising process. The detailed process of single interaction is illustrated in Fig. 4, where $\mathbf{L_{in}}$ and $\mathbf{L_{out}}$ represent the input and output features, respectively. The terms $\mathbf{G}_t$, $\mathbf{K}_t$, $\mathbf{V}_t$, $\mathbf{Q}_t$, $\mathbf{R}_t$, and $\mathbf{S}_t$ are all involved in computing the gated linear attention (GLA) [32], where $t$ denotes the index of token. We first compute self-attention on the normalized $\mathbf{L_{in}}$ using GLA to obtain self-decoded features. Then, following the standard transformer block architecture [8], we apply residual connections, normalization, cross-attention, and MLP to produce the final output features $\mathbf{L_{out}}$. Here, LCG embeddings are introduced via cross-attention, enabling the cross-decoding of self-decoded features with embeddings. Throughout the multiple interactions in the pipeline, self-decoding and cross-decoding alternate, and once the embeddings are first introduced via cross-decoding, all subsequent decoding steps incorporate LCG guidance.

Here, we elaborate on the process of computing the self-decoded features $\hat{\mathbf{L}}$ using GLA. We denote the linear projection weights and bias for a variable $x$ as $\mathbf{W}_x$ and $\mathbf{b}_x$, respectively. First, we compute the query ($\mathbf{Q}$), key ($\mathbf{K}$), and value ($\mathbf{V}$) based on $\mathbf{L_{in}} \in \mathbb{R}^{L \times d}$ (where $L$ is the sequence length and $d$ is the feature dimension) as follows:

$$\begin{aligned} \mathbf{Q} &= \mathbf{L_{in}}\mathbf{W}_Q \in \mathbb{R}^{L \times d_k}, \\ \mathbf{K} &= \mathbf{L_{in}}\mathbf{W}_K \in \mathbb{R}^{L \times d_k}, \quad (2) \\ \mathbf{V} &= \mathbf{L_{in}}\mathbf{W}_V \in \mathbb{R}^{L \times d_v}, \end{aligned}$$

here, $d_k$ represents the dimension of $\mathbf{Q}$ and $\mathbf{K}$, while $d_v$ represents that of $\mathbf{V}$. The subsequent computation is token-wise, and $t$ is the index of token. We use the 2D forget-gating matrix $\mathbf{G}_t$ to regulate the update of hidden states $\mathbf{S}_t$, enabling contextual interactions. The gating matrix $\mathbf{G}_t$ is computed as follows:

$$\begin{aligned} \mathbf{G}_t &= \boldsymbol{\alpha}_t^\top \boldsymbol{\beta}_t \in \mathbb{R}^{d_k \times d_v}, \\ \boldsymbol{\alpha}_t &= \sigma\left(\mathbf{L_{in}}\mathbf{W}_\alpha + \mathbf{b}_\alpha\right)^{\frac{1}{\tau}} \in \mathbb{R}^{L \times d_k}, \quad (3) \\ \boldsymbol{\beta}_t &= \sigma\left(\mathbf{L_{in}}\mathbf{W}_\beta + \mathbf{b}_\beta\right)^{\frac{1}{\tau}} \in \mathbb{R}^{L \times d_v}, \end{aligned}$$

4

where $\sigma$ represents the sigmoid function, $\beta$ is the bias term, and $\tau$ is the temperature term. Using $\mathbf{G}_t$ to update the hidden states $\mathbf{S}_t$, we obtain $\hat{\mathbf{L}}_\mathbf{t}$ as follows:

$$\begin{aligned}
\hat{\mathbf{L}}_\mathbf{t} &= (\mathbf{R}_t \odot \mathcal{L}(\mathbf{O}_t))\,\mathbf{W}_O \in \mathbb{R}^{1 \times d}, \\
\mathbf{R}_t &= \mathcal{S}(\mathbf{X}_t\mathbf{W}_r + \mathbf{b}_r) \in \mathbb{R}^{1 \times d_v}, \\
\mathbf{O}_t &= \mathbf{Q}_t^\top \mathbf{S}_t \in \mathbb{R}^{1 \times d_v}, \\
\mathbf{S}_t &= \mathbf{G}_t \odot \mathbf{S}_{t-1} + \mathbf{K}_t^\top V_t \in \mathbb{R}^{d_k \times d_v},
\end{aligned} \tag{4}$$

where $\odot$ denotes element-wise multiplication, $\mathcal{S}$ represents the Swish activation function [20], and $\mathcal{L}$ denotes Layer-Norm. After computing all $\hat{\mathbf{L}}_\mathbf{t}$, we obtain the final self-decoded features $\hat{\mathbf{L}}$.

# 4. Experiments

## 4.1. Datasets and Evaluation

**Construction of Our LCG Training Dataset.** We define 116 foreground categories and 21 background categories (detailed in the supplementary materials) and adopt an Auto-Labeling framework [24] by integrating AlphaCLIP and SAM [13] to obtain fine-grained segmentation masks for foreground and background across multiple datasets. Specifically, we utilize the following datasets: COCONut-Large [5] (0.36M images), Object365V2 [23] (2.02M images), GoogleLandmarkV2 [29] (4.13M images), and a natural scene dataset (7.49M images) that we curated and collected. In total, our dataset comprises 14 million images. Following the masking strategy described in Sec. 3.2, we construct a large-scale training dataset aligned with the LCG paradigm, consisting of 4.3M "foreground" image-mask pairs and 9.7M "background" image-mask pairs.

**Finetune and Evaluation.** Recent SOTA inpainting methods [14, 22, 25, 40] are typically evaluated on one or more public benchmarks, including Places2 [38], CelebA-HQ [26], and FFHQ [11]. To ensure a comprehensive comparison, we fine-tune PixelHacker on Places2, CelebA-HQ, and FFHQ while strictly following the evaluation protocols used in previous works. For Places2 (a natural scene dataset), we follow multiple evaluation settings: (1) Zhuang et al. [40] and Rombach et al. [21] sample 10K test images, apply center cropping to 512×512 resolution, and evaluate performance under 40-50% masked areas. (2) Li et al. [14] use 36.5K validation images, crop them to 512×512, and define two mask configurations (small and large masks) for evaluation. (3) Sargsyan et al. [22] resize the Places2 validation set to 256×256 and 512×512, applying highly randomized masking strategies for evaluation. To ensure fair comparisons, we strictly adhere to these evaluation protocols. For CelebA-HQ (a human face dataset), we follow Li et al. [14] and evaluate at 512×512 resolution. For FFHQ

(a human face dataset), we sample 10K images as the evaluation set and follow the masking strategy of Suvorov et al. [25], evaluating at 256×256 resolution for comparison.

## 4.2. Implementation Details

For training on our 14M-image dataset, we use 12 NVIDIA L40S GPUs with a batch size of 528 and a resolution of 512×512. The probabilities $\mathcal{P}_{rand}$ and $\mathcal{P}_{obj}$ are both set to 0.5, and the model is trained for 200K iterations. For fine-tuning on Places2, we use the 1.8M training set of Places2 and fine-tune the model for 120K iterations on 12 NVIDIA L40S GPUs with a batch size 528. For fine-tuning on CelebA-HQ, we follow the same training set split as Li et al. [14] and fine-tune the model for 90K iterations on 8 NVIDIA L40S GPUs with a batch size 352. For fine-tuning on FFHQ, we use 60K images as the training set and 10K images as the evaluation set. The model is fine-tuned for 60K iterations on 12 NVIDIA L40S GPUs with a batch size 528. In all experiments, we use the pretrained SDXL-VAE [19] and freeze its parameters. The input resolution is 512×512, the learning rate is set to 1e-5, and we adopt the AdamW optimizer [16] with betas (0.9, 0.999).

## 4.3. Comparison on Places2

We faithfully follow the three evaluation settings on Places2 and conduct fair comparisons with various SOTA methods. The results in Tab.1, Tab.2, and Tab. 3 correspond to evaluation settings (1), (2), and (3) described in Sec. 4.1, respectively. In particular, LDM [21] refers to a latent diffusion model fine-tuned for inpainting without text prompts. SD [21] denotes SDv1.5-Inpainting, a fine-tuned Stable Diffusion model trained with random masks and image captions for inpainting tasks. SDXL [19] represents SDXL-

| Method | Places2 (Test) | |
|---|---|---|
| | FID ↓ | LPIPS ↓ |
| LaMa [25] | 21.07 | 0.2133 |
| LDM [21] | 21.42 | 0.2317 |
| SD [21] | 19.73 | 0.2322 |
| SD("background") [21] | 19.21 | 0.2290 |
| SD("scenery") [21] | 18.93 | 0.2312 |
| SmartBrush("scenery") [30] | 87.21 | 0.2812 |
| MI-GAN [22] | 14.36 | 0.2390 |
| SDXL [19] | 8.66 | 0.2746 |
| PowerPaint [40] | 17.91 | 0.2225 |
| PixelHacker(Ours)‡ | 12.19 | 0.2100 |
| PixelHacker(Ours) | 8.59 | 0.2026 |

Table 1. **Quantitative comparison of PixelHacker with SOTA methods on 10k samples from Places2 [38] test set of 512 resolution with 40-50% masks.** Results of other methods except MI-GAN and SDXL are referred from [40]. "‡": our model without finetune on Places2 [38]. The best results are in red while the second best ones are in blue.

5

| Method | Places2 (Large Mask) | | | | Places2 (Small Mask) | | | |
|---|---|---|---|---|---|---|---|---|
| | FID ↓ | LPIPS ↓ | P-IDS(%) ↑ | U-IDS(%) ↑ | FID ↓ | LPIPS ↓ | P-IDS(%) ↑ | U-IDS(%) ↑ |
| CoModGAN [37]† | 2.92 | 0.192 | 19.64 | 35.78 | 1.10 | 0.101 | 26.95 | 41.88 |
| LaMa-Big [25]† | 2.97 | 0.166 | 13.09 | 32.29 | 0.99 | 0.086 | 22.79 | 40.58 |
| MADF [39] | 7.53 | 0.181 | 6.00 | 23.78 | 2.24 | 0.095 | 14.85 | 35.03 |
| AOT-GAN [36] | 10.64 | 0.195 | 3.07 | 19.92 | 3.19 | 0.101 | 8.07 | 30.94 |
| HiFill [33] | 28.92 | 0.284 | 1.24 | 11.24 | 7.94 | 0.148 | 3.98 | 23.60 |
| DeepFill v2 [35] | 9.27 | 0.213 | 4.01 | 21.32 | 3.02 | 0.113 | 9.17 | 32.56 |
| EdgeConnect [17] | 12.66 | 0.275 | 1.93 | 15.87 | 4.03 | 0.114 | 5.88 | 27.56 |
| MAT [14] | 2.90 | 0.189 | 19.03 | 35.36 | 1.07 | 0.099 | 27.42 | 41.93 |
| PixelHacker(Ours)‡ | 3.01 | 0.176 | 13.79 | 33.53 | 1.05 | 0.093 | 23.19 | 40.87 |
| PixelHacker(Ours) | 2.05 | 0.169 | 16.73 | 36.07 | 0.82 | 0.088 | 24.78 | 42.21 |

Table 2. **Quantitative comparison of PixelHacker with SOTA methods on Places2 [38] validation set of 512 resolution with large and small masks follow [14].** Results other than our method referred from [14]. "†": CoModGAN [37] use 8M training images on Places2, LaMa [25] use 4.5M training images, while other models (without "†") including our PixelHacker are trained on 1.8M training images from the standard part of Places2 [38]. Red, blue and "‡" are same to Tab. 1.
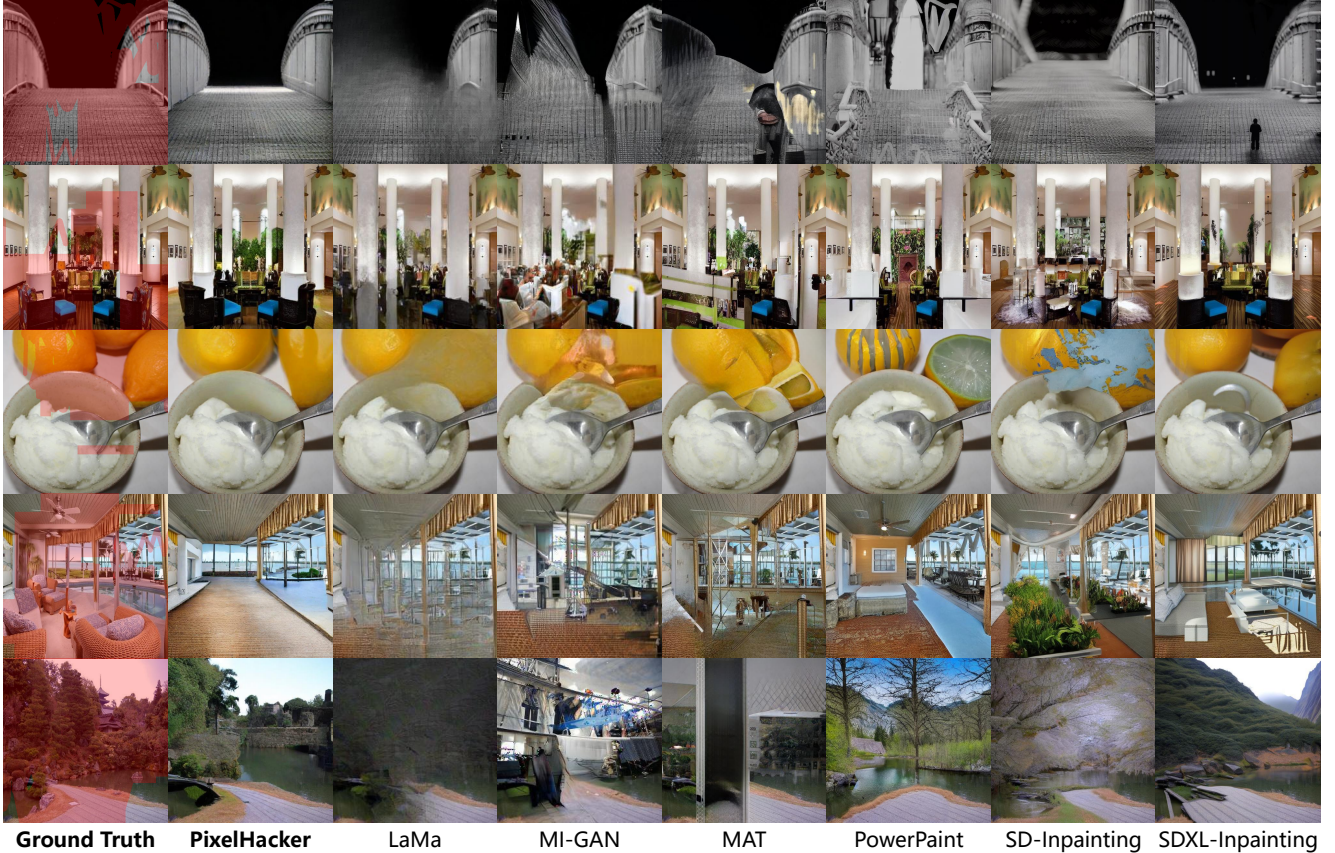


Figure 5. **Qualitative comparison of PixelHacker with SOTA methods on Places2.** Even under masks that cover almost the entire image, PixelHacker's generated results still exhibit remarkable structural and semantic consistency.

Inpainting, a latent text-to-image diffusion model fine-tuned from SDXL-Base [19] to enable mask-guided inpainting.

First, as shown in Tab. 1, we conduct quantitative comparisons on the Places2 test set using 512 resolution with 40-50% masked areas. The results for MI-GAN [22] and SDXL [19] are obtained using the official inference codes, while results for other models are taken from [40]. Our PixelHacker achieves the best performance, with FID **8.59** and LPIPS **0.2026**, surpassing the powerful SD [21] and SDXL [19]. Notably, even without fine-tuning, the zero-shot version of PixelHacker achieves the best LPIPS and second-best FID, trailing only SDXL. This strongly demonstrates the remarkable potential of our paradigm.

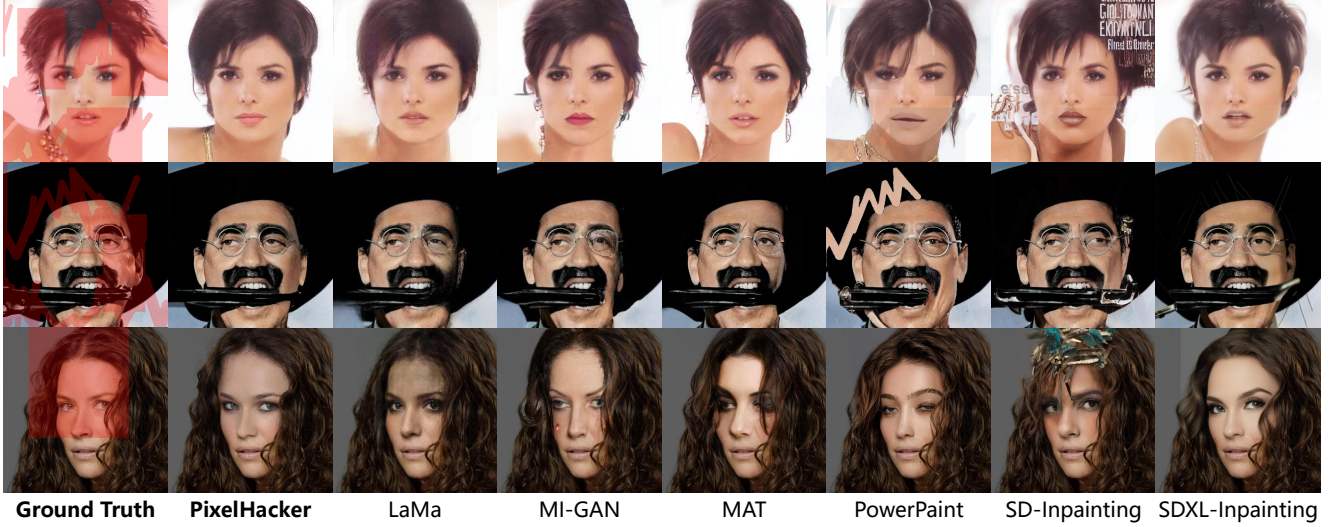Next, as shown in Tab. 2, we perform quantitative com-

**Figure 6. Qualitative comparison with SOTA methods on CelebA-HQ.** PixelHacker demonstrates exceptional robustness to mask shapes while maintaining highly consistent semantics, avoiding color discrepancies and artifacts commonly observed in other models.

| Method | Places2 (256) | | Places2 (512) | |
|---|---|---|---|---|
| | FID ↓ | LPIPS ↓ | FID ↓ | LPIPS ↓ |
| LaMa [25]† | 22.00 | 0.378 | 12.36 | 0.314 |
| CoModGAN [37]† | 9.32 | 0.397 | 8.05 | 0.343 |
| SH-GAN [31] | 7.40 | 0.392 | 7.03 | 0.339 |
| ZITS [7] | 16.78 | 0.356 | 12.94 | 0.310 |
| MAT [14] | 14.38 | 0.394 | 8.67 | 0.339 |
| LDM [21] | 13.40 | 0.385 | 8.46 | 0.342 |
| HiFill [33] | 81.27 | 0.488 | 64.07 | 0.438 |
| MI-GAN [22] | 11.83 | 0.394 | 10.00 | 0.345 |
| PixelHacker(Ours)‡ | 14.60 | 0.380 | 9.40 | 0.317 |
| PixelHacker(Ours) | 9.25 | 0.367 | 5.75 | 0.305 |

Table 3. **Quantitative comparison of PixelHacker with SOTA approaches on Places2 [38] validation set of 256 and 512 resolutions follow [22].** Results of other methods are referred from [22]. The images are resized from the High-Resolution version of Places2 dataset[38], into 256 and 512 respectively. Red, blue and "‡" are same to Tab. 1. "†" is same to Tab. 2

parisons on the Places2 validation set, following the large and small mask settings of [14], using 512 resolution. Our PixelHacker, finetuned with only 1.8M images from Places2, achieves the best FID and U-IDS, the second-best LPIPS, and the third-best P-IDS. It is worth noting that: LaMa-Big [25], which achieves the best LPIPS, was trained with 4.5M images from Places2. CoModGAN [37], which surpasses PixelHacker in P-IDS, was trained with the full 8M Places2 dataset. Both models utilize significantly more training images from Places2 than PixelHacker, highlighting the data efficiency of our paradigm.

Finally, as shown in Tab. 3, we follow [22] and conduct quantitative comparisons on the Places2 validation set at 256 and 512 resolutions. Our PixelHacker achieves SOTA

performance at 512 resolution and second-best results at 256 resolution. Notably, PixelHacker is trained and fine-tuned exclusively at 512 resolution, SH-GAN [31] is trained at both 256 and 512 resolutions, ZITS [7] undergoes extensive training across multiple resolutions between 256 and 512. Despite this, PixelHacker still achieves competitive results, demonstrating its robust generalization capability.

We present qualitative comparison results on Places2, as shown in Fig. 5. Even under masks that cover almost the entire image, PixelHacker's generated results still exhibit remarkable structural and semantic consistency. Unlike LaMa [25] and MI-GAN [22], which produce artifacts and ambiguous generations (e.g., the first row), PixelHacker does not suffer from such issues. Under medium-sized masks, even in semantically rich and structurally complex scenes, PixelHacker consistently maintains structural coherence (e.g., the second and third rows).

### 4.4. Comparison on CelebA-HQ and FFHQ

**CelebA-HQ.** As shown in Tab. 4, we follow [14] and conduct quantitative comparisons on CelebA-HQ using 512 resolution. PixelHacker consistently achieves SOTA performance, demonstrating its strong generalization ability when transferring from natural scenes to facial image domains under our LCG paradigm. We further provide qualitative comparisons, as shown in Fig. 6. PixelHacker generates clear and well-formed facial features without noticeable distortions, maintaining strong semantic consistency. Moreover, PixelHacker does not introduce abrupt or irrelevant textures, indicating superior structural consistency.

**FFHQ.** Tab. 5 presents PixelHacker's SOTA quantitative results on FFHQ [11], demonstrating its remarkable generalization ability to lower resolutions despite being exclu-

7

Figure 7. **Qualitative comparison with SOTA methods on FFHQ.** Our PixelHacker generates more realistic results compared to other methods and exhibits strong adaptability to scenes with complex hierarchies and challenging lighting conditions.

| Method | CelebA-HQ | |
| --- | --- | --- |
| | FID ↓ | LPIPS ↓ |
| CoModGAN [37] | 5.65 | 0.140 |
| LaMa [25] | 8.15 | 0.143 |
| ICT [27] | 12.84 | 0.195 |
| MADF [39] | 6.83 | 0.130 |
| AOT-GAN [36] | 10.82 | 0.145 |
| DeepFill v2 [35] | 24.42 | 0.221 |
| EdgeConnect [17] | 39.99 | 0.208 |
| MAT [14] | 4.86 | 0.125 |
| PixelHacker(Ours) | 4.75 | 0.115 |

Table 4. **Quantitative comparison with SOTA methods on CelebA-HQ in 512 resolution follow [14].** Results of other methods are referred from [14]. The best results are in red.

| Method | FFHQ | |
| --- | --- | --- |
| | FID ↓ | LPIPS ↓ |
| LaMa [25] | 52.17 | 0.366 |
| MI-GAN [22] | 27.65 | 0.358 |
| SD [21] | 40.24 | 0.359 |
| PowerPaint [40] | 38.25 | 0.409 |
| SDXL [19] | 12.36 | 0.278 |
| PixelHacker(Ours) | 6.35 | 0.229 |

Table 5. **Quantitative comparison with SOTA methods on FFHQ in 256 resolution.** The best results are in red.

sively trained at 512 resolution. Fig. 7 provides qualitative comparisons, illustrating that PixelHacker produces more realistic results than other methods while exhibiting strong adaptability to complex scene hierarchies and challenging lighting conditions.

## 4.5. Ablation Study

**Ablation of various masks in LCG.** We evaluate the impact of the mask construction strategy proposed in Sec. 3.2, as shown in Tab. 6. For each setting, we fairly initialize the model using weights trained for 90K iterations on our 14M dataset, then continue training 12K iterations before

| $\mathcal{M}_{obj}$ | $\mathcal{M}_{scene}$ | $\mathcal{M}_{rand}$ | $\mathcal{M}'_{obj}$ | FID ↓ | LPIPS ↓ |
| --- | --- | --- | --- | --- | --- |
| ✓ | ✓ | | | 12.69 | 0.1758 |
| ✓ | ✓ | ✓ | | 12.65 | 0.1735 |
| ✓ | ✓ | ✓ | ✓ | 12.62 | 0.1735 |

Table 6. **Ablation of various masks in LCG.** Best results in red.

| E.Dim | FID ↓ | LPIPS ↓ | E.Dim | FID ↓ | LPIPS ↓ |
| --- | --- | --- | --- | --- | --- |
| 20† | 0.755 | 0.264 | 256 | 0.754 | 0.258 |
| 64 | 0.735 | 0.257 | 1024 | 0.744 | 0.255 |

Table 7. **Ablation of embedding dims.** "†": we use 20 as the default size in PixelHacker. The best results are in red.

| Scale | 1.0 | 1.5 | 2.0 | 2.5 | 3.0 | 3.5 | 4.0 |
| --- | --- | --- | --- | --- | --- | --- | --- |
| FID ↓ | 9.70 | 8.60 | 8.44 | 8.47 | 8.48 | 8.51 | 8.73 |
| LPIPS ↓ | 0.124 | 0.119 | 0.118 | 0.119 | 0.120 | 0.121 | 0.123 |

Table 8. **Ablation of guidance scales.** The best scale is 2.0.

evaluating on a 3K-image subset of the Places2 validation set [38]. The results demonstrate that using $\mathcal{M}_{obj}$, $\mathcal{M}_{scene}$, $\mathcal{M}_{rand}$, and $\mathcal{M}'_{obj}$ yields the best performance.

**Ablation of embedding dims.** To investigate whether increasing the embedding dimension (E.Dim) leads to significant performance improvements, we compare various E.Dim values of 20, 64, 256, and 1024, as presented in Tab. 7. For each configuration, we fairly initialize the model using weights trained for 200K iterations on our 14M dataset, then continue training 36K iterations, evaluating on custom 10K images. The results indicate that scaling up E.Dim does not significantly enhance performance. This suggests that a smaller E.Dim is sufficient to represent latent "foreground" and "background" features. Based on these findings, we adopt E.Dim = 20 as the default setting.

**Ablation of guidance scales.** Similar to other Latent Diffusion Models, we employ Classifier-Free Guidance inference [9]. Here, to assess the impact of guidance scale on generation quality, we compare multiple scales ranging from 1.0 (no guidance) to 4.0, as shown in Tab. 8. For each configuration, we initialize the model using weights that fine-tune on CelebA-HQ for 21K iterations, and evaluate on the CelebA-HQ validation set used in Sec. 4.4. The results confirm that 2.0 is the best, which we adopt as default.

## 5. Conclusion

In this work, we introduce Latent Categories Guidance (LCG), a simple yet effective inpainting paradigm that guides the model toward structural and semantic consistency through latent foreground and background features. Then, we propose PixelHacker, a diffusion-based inpainting model trained with LCG on 14M image-mask pairs and fine-tuning on open-source benchmarks. Extensive experiments demonstrate that PixelHacker consistently achieves SOTA performance across various benchmarks.

# References

[1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18208–18218, 2022. 4

[2] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, page 417–424, USA, 2000. ACM Press/Addison-Wesley Publishing Co. 2

[3] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8789–8797, 2018. 2

[4] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8185–8194, 2020. 2

[5] Xueqing Deng, Qihang Yu, Peng Wang, Xiaohui Shen, and Liang-Chieh Chen. Coconut: Modernizing coco segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 5

[6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, pages 8780–8794. Curran Associates, Inc., 2021. 4

[7] Qiaole Dong, Chenjie Cao, and Yanwei Fu. Incremental transformer structure enhanced image inpainting with masking positional encoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 7

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 4

[9] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 3, 4, 8

[10] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. 2

[11] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(12):4217–4228, 2021. 2, 3, 5, 7, 12

[12] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. 3

[13] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026, 2023. 5

[14] Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. Mat: Mask-aware transformer for large hole image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2, 5, 6, 7, 8, 11

[15] 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 *The European Conference on Computer Vision (ECCV)*, 2018. 2

[16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 5

[17] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Qureshi, and Mehran Ebrahimi. Edgeconnect: Structure guided image inpainting using edge prediction. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*, 2019. 6, 8

[18] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2536–2544, 2016. 3

[19] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. 3, 4, 5, 6, 8, 11

[20] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017. 5

[21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 3, 4, 5, 6, 7, 8, 11

[22] Andranik Sargsyan, Shant Navasardyan, Xingqian Xu, and Humphrey Shi. Mi-gan: A simple baseline for image inpainting on mobile devices. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7335–7345, 2023. 2, 5, 6, 7, 8, 11

[23] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8429–8438, 2019. 5

[24] Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Alphaclip: A clip model focusing on wherever you want, 2023. 5

[25] 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. 2, 5, 6, 7, 8, 11

[26] Samuli Laine Tero Karras, Timo Aila and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability and variation. In *International Conference on Learning Representations (ICLR)*, 2018. 2, 3, 5, 12

[27] Ziyu Wan, Jingbo Zhang, Dongdong Chen, and Jing Liao. High-fidelity pluralistic image completion with transformers. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4672–4681, 2021. 8

[28] Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18359–18369, 2023. 2, 3, 4

[29] Tobias Weyand, André Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 – a large-scale benchmark for instance-level recognition and retrieval. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2572–2581, 2020. 5

[30] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22428–22437, 2023. 3, 4, 5

[31] Xingqian Xu, Shant Navasardyan, Vahram Tadevosyan, Andranik Sargsyan, Yadong Mu, and Humphrey Shi. Image completion with heterogeneously filtered spectral hints. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4591–4601, 2023. 2, 7

[32] Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. In *Proceedings of ICML*, 2024. 4

[33] Zili Yi, Qiang Tang, Shekoofeh Azizi, Daesik Jang, and Zhan Xu. Contextual residual aggregation for ultra high-resolution image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7508–7517, 2020. 6, 7

[34] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. Generative image inpainting with contextual attention. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018. 2, 3

[35] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas Huang. Free-form image inpainting with gated convolution. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4470–4479, 2019. 2, 6, 8

[36] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Aggregated contextual transformations for high-resolution image inpainting. *IEEE Transactions on Visualization and Computer Graphics*, 29(7):3266–3280, 2023. 2, 6, 8

[37] Shengyu Zhao, Jonathan Cui, Yilun Sheng, Yue Dong, Xiao Liang, Eric I Chang, and Yan Xu. Large scale image completion via co-modulated generative adversarial networks. In *International Conference on Learning Representations (ICLR)*, 2021. 2, 6, 7, 8

[38] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 2, 3, 5, 6, 7, 8, 11

[39] 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. *IEEE Transactions on Image Processing*, 30:4855–4866, 2021. 6, 8

[40] Junhao Zhuang, Yanhong Zeng, Wenran Liu, Chun Yuan, and Kai Chen. A task is worth one word: Learning with task prompts for high-quality versatile image inpainting, 2023. 3, 4, 5, 6, 8, 11

# Supplementary Materials

**Potential Categories.** (1) "foreground": person, bicycle, car, motorcycle, airplane, bus, train, truck, boat, traffic light, fire hydrant, stop sign, parking meter, bench, bird, cat, dog, horse, sheep, cow, elephant, bear, zebra, giraffe, backpack, umbrella, handbag, tie, suitcase, frisbee, skis, snowboard, sports ball, kite, baseball bat, baseball glove, skateboard, surfboard, tennis racket, bottle, wine glass, cup, fork, knife, spoon, bowl, banana, apple, sandwich, orange, broccoli, carrot, hot dog, pizza, donut, cake, chair, couch, potted plant, bed, dining table, toilet, tv, laptop, mouse, remote, keyboard, cell phone, microwave, oven, toaster, sink, refrigera-tor, book, clock, vase, scissors, teddy bear, hair drier, toothbrush, banner, blanket, bridge, cardboard, counter, curtain, door-stuff, flower, fruit, gravel, house, light, mirror-stuff, net, pillow, railroad, roof, shelf, stairs, tent, towel, window-blind, window-other, tree-merged, fence-merged, cabinet-merged, mountain-merged, dirt-merged, paper-merged, food-other-merged, building-other-merged, rock-merged, rug-merged, trash can, clothes, shadow. (2) "background": background, floor-wood, platform, playing-field, river, road, sand, sea, snow, wall-brick, wall-stone, wall-tile, wall-wood, water-other, ceiling-merged, sky-other-merged, table-merged, floor-other-merged, pavement-merged, grass-merged, wall-other-merged.
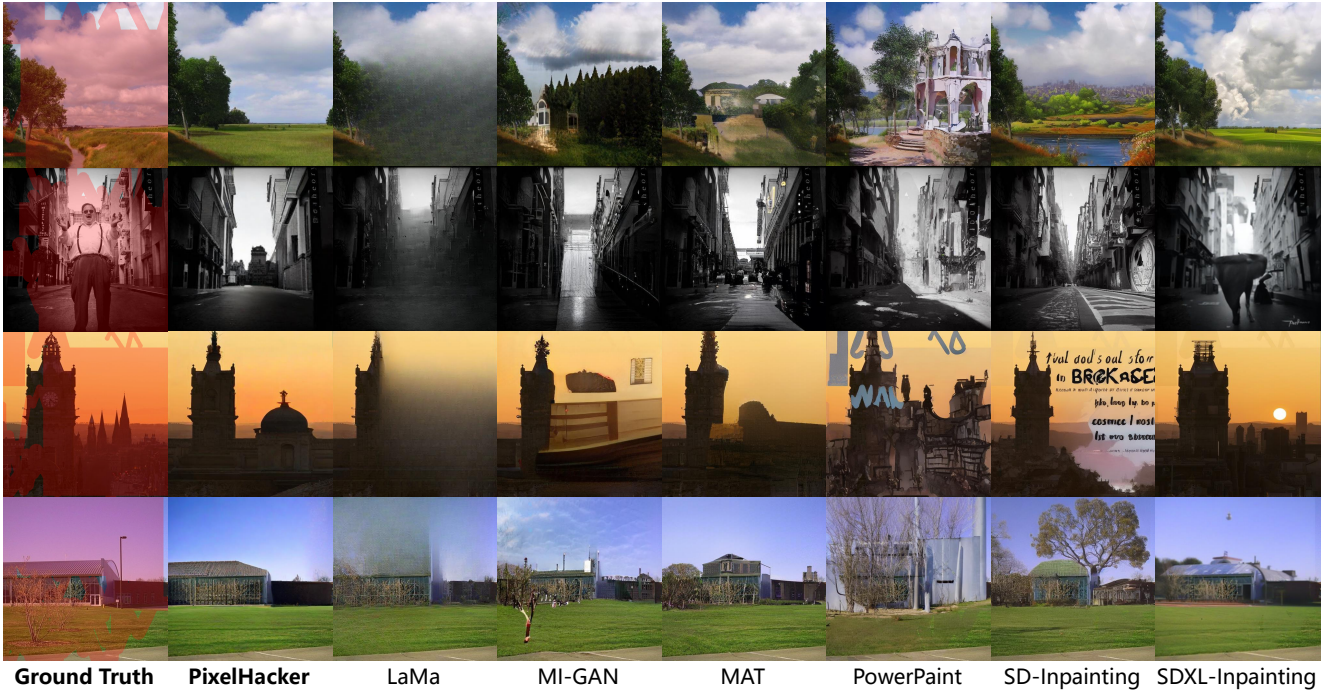


Figure 8. **More qualitative results of our PixelHacker vs. other SOTA methods [14, 19, 21, 22, 25, 40] on Places2 [38].** Even under masks that obscure almost the entire image, PixelHacker consistently maintains remarkable semantic and structural coherence.



Figure 9. **Failure cases on natural and human-face scenes.** In the first row, our PixelHacker struggles with generating fine details for the small black humanoid silhouette; in the second row, PixelHacker exhibits suboptimal reconstruction of the fingers. Nevertheless, it still maintains significantly better semantic and structural coherence compared to all other SOTA methods.

11

(a) More qualitative results of our PixelHacker vs. other SOTA methods on CelebA-HQ



(b) More qualitative results of our PixelHacker vs. other SOTA methods on FFHQ

Figure 10. **More qualitative results of our PixelHacker vs. other SOTA methods on CelebA-HQ [26] and FFHQ [11].** PixelHacker preserves highly detailed facial textures while avoiding perceptible color discrepancies or artifacts, producing more realistic results than other methods and demonstrating strong adaptability to complex scene structures and challenging lighting conditions.

**More Qualitative Results.** We provide more qualitative results of our PixelHacker with other SOTA methods on Places2, CelebA-HQ, and FFHQ, in Fig. 8 and Fig. 10.

**Failure Cases.** We provide some failure cases as shown in Fig. 9. PixelHacker works imperfectly on these examples, but is still significantly better than other methods.