# RealFill: Reference-Driven Generation for Authentic Image Completion

LUMING TANG, Cornell University, US
NATANIEL RUIZ, Google Research, US
QINGHAO CHU, Google Research, US
YUANZHEN LI, Google Research, US
ALEKSANDER HOŁYŃSKI, Google Research, US
DAVID E. JACOBS, Google Research, US
BHARATH HARIHARAN, Cornell University, US
YAEL PRITCH, Google Research, Israel
NEAL WADHWA, Google Research, US
KFIR ABERMAN, Snap Research, US
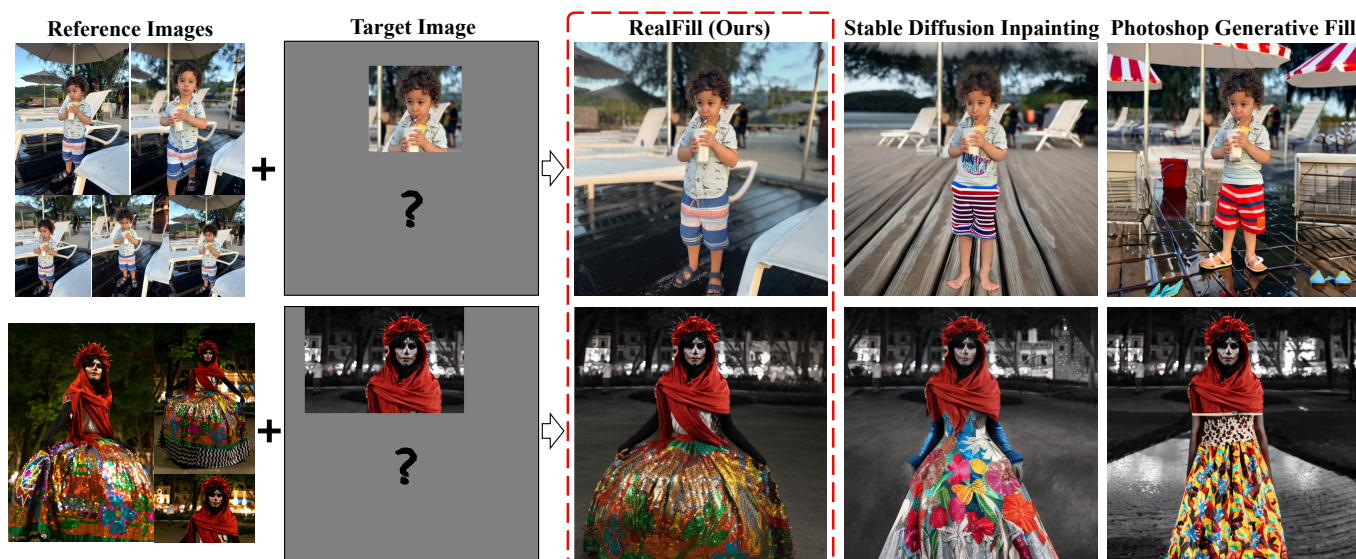MICHAEL RUBINSTEIN, Google Research, US

Fig. 1. Given a few reference images that roughly capture the same scene, and a target image with a missing region, **RealFill** is able to complete the target image with image content that is faithful to the true scene. In contrast, standard prompt-based inpainting methods hallucinate plausible but inauthentic content due to their lack of knowledge of the original scene.

Recent advances in generative imagery have brought forth outpainting and inpainting models that can produce high-quality, plausible image content in unknown regions. However, the content these models hallucinate is necessarily inauthentic, since they are unaware of the true scene. In this work, we propose **RealFill**, a novel generative approach for image completion that fills in missing regions of an image with the content that should have been there. RealFill is a generative inpainting model that is personalized using only a few reference images of a scene. These reference images do not have to be aligned with the target image, and can be taken with drastically varying viewpoints, lighting conditions, camera apertures, or image styles. Once personalized, RealFill is able to complete a target image with visually compelling contents that are faithful to the original scene. We evaluate RealFill on a new image completion benchmark that covers a set of diverse and challenging scenarios, and find that it outperforms existing approaches by a large margin. Project page: https://realfill.github.io.

CCS Concepts: • **Computing methodologies** → **Computational photography**; **Image processing**; **Computer vision**.

Additional Key Words and Phrases: Image Completion, Diffusion Model

## 1 INTRODUCTION

Photographs capture ephemeral and invaluable experiences in our lives, but can sometimes fail to do these memories justice. In many cases, no single shot may have captured the perfect angle, framing, timing, or composition, and unfortunately, just as the experiences themselves cannot be revisited, these elements of the captured images are also unalterable. We show one such example in Fig. 2: imagine having taken a nearly perfect photo of your daughter dancing on stage, but her unique and intricate crown is partially cut out of the frame. Of course, there are many other pictures from the performance that showcase her crown, but they all fail to capture that precise special moment: her pose mid-dance, her facial expression, and the perfect lighting. Given this collection of imperfect photos, you can certainly imagine the missing parts of this perfect shot, but actually creating a complete, shareable version of this image is much harder.

In this paper, we focus on this problem, which we call *Authentic Image Completion*. Given a few reference images (up to five) and one target image that captures roughly the same scene (but in a different arrangement or appearance), we aim to fill missing regions of the target image with high-quality image content that is faithful to the originally captured scene. Note that for the sake of practical benefit, we focus particularly on the more challenging, unconstrained setting in which the target and reference images may have very different viewpoints, environmental conditions, camera apertures, image styles, or even moving objects.

Approaches to solving variants of this problem have been proposed using classical geometry-based pipelines [Shan et al. 2014; Zhao et al. 2023a; Zhou et al. 2021] that rely on correspondence matching, depth estimation, and 3D transformations, followed by patch fusion and image harmonization. These methods tend to encounter catastrophic failure when the scene's structure cannot be accurately estimated, e.g., when the scene geometry is too complex or contains dynamic objects. On the other hand, recent generative models [Chang et al. 2023, 2022; Yu et al. 2018], and in particular diffusion models [Ho et al. 2020; Rombach et al. 2022; Song et al. 2021], have demonstrated strong performance on the tasks of image inpainting and outpainting [Ramesh et al. 2022; Stability AI 2022; Wang et al. 2023]. These methods, however, struggle to recover genuine scene structure and fine details, since they are only guided by text prompts, and can't use reference image content.

To this end, we present a simple yet effective reference-driven image completion framework called *RealFill*. For a given scene, we first create a personalized generative model by finetuning a pretrained inpainting diffusion model [Stability AI 2022] on the reference and target images. This finetuning process is designed such that the adapted model not only maintains a good image prior, but also learns the contents, lighting, and style of the scene in the input images. We then use this finetuned model to fill the missing regions

in the target image through a standard diffusion sampling process. Given the stochastic nature of generative inference, we propose *Correspondence-Based Seed Selection*, to automatically select a small set of high-quality generations by exploiting a special property of our completion task: that there exists true correspondences between generated contents and reference images. Specifically, we filter out samples that have too few keypoint matches with references, which greatly reduces human labor to select high-quality model outputs.

As shown in Figs. 1 to 4, RealFill is able to effectively inpaint and outpaint a target image with its *genuine* scene content. Most importantly, our method is able to handle large differences between reference and target images, e.g., viewpoint, lighting, aperture, style or dynamic deformations — differences which are very difficult for previous geometry-based approaches. Existing benchmarks for image completion [Zhou et al. 2021] mainly focus on small inpainting tasks and minimal changes between reference and target images. In order to quantitatively evaluate the aforementioned challenging use-case, we collect a dataset containing 10 inpainting and 23 outpainting examples along with corresponding ground-truth, and show that RealFill outperforms baselines by a large margin across multiple image similarity metrics.

In summary, our contributions are as follows:

- We define a new problem named *Authentic Image Completion*, i.e., given a set of reference images and a target image with missing regions, we seek to complete those missing regions with content that is faithful to the scene as captured in the references. In essence, the goal is to complete the target image with what "should have been there" rather than what "could have been there", like in typical generative inpainting.
- We introduce *RealFill*, a method that aims to solve this problem by finetuning an inpainting diffusion model on reference and target images. This model is sampled with *Correspondence-Based Seed Selection* to filter outputs with low fidelity to the reference images. *RealFill* is the first method that expands the expressive power of generative inpainting models by conditioning the process beyond text, enabling extra conditioning on reference images.
- We propose *RealBench*, a dataset for quantitative evaluation of authentic image completion, composed of 33 scenes spanning both inpainting and outpainting tasks.

## 2 RELATED WORK

**Adapting Pretrained Diffusion Models**. Diffusion models [Dhariwal and Nichol 2021; Ho et al. 2020; Song et al. 2021] have shown strong performance in text-to-image (T2I) generation [Ramesh et al. 2022; Rombach et al. 2022; Saharia et al. 2022]. Recent works make use of this pretrained image prior by finetuning them for various tasks. Personalization methods propose to finetune the T2I model [Avrahami et al. 2023; Chen et al. 2023; Ruiz et al. 2023a,b] or text embedding [Gal et al. 2022; Voynov et al. 2023], on a few images to achieve arbitrary text-driven generation of a given object or style. Other techniques instead finetune a T2I model to add new conditioning signals, either for image editing [Brooks et al. 2023; Kawar et al. 2023; Wang et al. 2023] or more controllable generation [Mou et al. 2023; Sohn et al. 2023; Zhang et al. 2023b]. The same

approach can be also applied to specialized tasks [Liu et al. 2023; Raj et al. 2023; Wu et al. 2022; Zhao et al. 2023b] such as converting a T2I model into a 3D or video generation model. Our method shows that a pretrained T2I inpainting diffusion model can be adapted to perform reference-driven image completion.

**Image Completion**. As an enduring challenge in computer vision, image completion aims to fill missing parts of an image with plausible content, i.e., inpainting and outpainting. Traditional approaches [Barnes et al. 2009; Bertalmio et al. 2000; Criminisi et al. 2003; Hays and Efros 2007] rely on handcrafted heuristics while more recent deep learning-based methods [Iizuka et al. 2017; Kim et al. 2022; Liu et al. 2018; Suvorov et al. 2022] instead directly train end-to-end neural networks that take original image and mask as inputs and generate the completed image. Given the challenging nature of this problem [Zheng et al. 2019], many works [Chang et al. 2023, 2022; Lugmayr et al. 2022; Yeh et al. 2017] propose to leverage the image prior from a pretrained generative model for this task. Built upon powerful T2I diffusion models, recent solutions [Adobe Inc. 2023; Ramesh et al. 2022; Stability AI 2022] demonstrate strong text-driven image completion capabilities. However, due to their sole dependence on a text prompt (which has limited descriptive power), generated image content can often be hard to control, resulting in tedious prompt engineering, especially when a particular or otherwise true scene content is desired. This is one of the main issues we aim to tackle in our work.

**Reference-Based Image Inpainting**. Existing works for reference-based inpainting [Zhao et al. 2023a; Zhou et al. 2021] or outpainting [Shan et al. 2014] usually make use of carefully tuned pipelines containing many individual components like depth and pose estimation, image warping, and harmonization. Each of these modules usually tackles a moderately challenging problem itself and the resulting prediction error can, and often does, propagate and accumulate through the pipeline. This can lead to catastrophic failure especially in challenging cases with complex scene geometry, changes in appearance, or scene deformation. Paint-by-Example [Yang et al. 2023] proposes a novel latent diffusion model [Rombach et al. 2022] whose generation is conditioned on both a reference and target image. However, the conditioning is based on CLIP embedding [Radford et al. 2021] of a single reference image, therefore is only able to capture high-level semantics of the reference object. In contrast, our method is the first to demonstrate multiple reference image-driven inpainting and outpainting that is both visually compelling and faithful to the original scene, even in cases where there are large appearance changes between reference and target images.

## 3 METHOD

### 3.1 Reference-Based Image Completion

Given a set of casually captured reference images (up to five), our goal is to complete (either outpaint or inpaint) a target image of roughly the same scene. The output image is expected to be not only plausible and photorealistic, but also faithful to the reference images, i.e., recovering content and scene detail that were present in the actual scene. In essence, we want to achieve *authentic image completion*, where we generate what "should have been there" instead of what "could have been there". We purposefully pose this

as a broad and challenging problem with few constraints on the inputs. For example, the images could be taken from very different viewpoints with unknown camera poses. They could also have different lighting conditions or styles, and the scene could potentially be non-static and have significantly varying layout across images.

In this section, we first provide background knowledge on diffusion models and subject-driven generation (Sec. 3.2). Then, we formally define the problem of authentic image completion (Sec. 3.3). Finally, we present RealFill, our method to perform reference-based image completion with a pretrained diffusion image prior (Sec. 3.4).

### 3.2 Preliminaries

**Diffusion models** are generative models that aim to transform a Normal distribution into an arbitrary target data distribution. During training, different magnitudes of Gaussian noise are added to a data point $x_0$ to obtain noisy $x_t$:

$$x_t = \sqrt{\alpha_t} x_0 + (\sqrt{1 - \alpha_t})\epsilon \tag{1}$$

where the noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$, and $\{\alpha_t\}_{t=1}^{T}$ define a fixed noise schedule with larger $t$ corresponding to more noise. Then, a neural network $\epsilon_\theta$ is trained to predict the noise using the following loss function:

$$\mathcal{L} = \mathbb{E}_{x,t,\epsilon} \|\epsilon_\theta(x_t, t, c) - \epsilon\|_2^2 \tag{2}$$

where $\epsilon_\theta$ is conditioned on some signal $c$, e.g., a language prompt for a T2I model, a masked image for an inpainting model. During inference, starting from $x_T \sim \mathcal{N}(0, \mathbf{I})$, $\epsilon_\theta$ is used to iteratively remove noise from $x_t$ to get a less noisy $x_{t-1}$, eventually leading to a sample $x_0$ from the target data distribution.

**DreamBooth** [Ruiz et al. 2023a] enables T2I diffusion models to generate images of a specific subject with semantic modifications. The core idea is to finetune the model $\epsilon_\theta$ on a few subject images using the loss in Eq. 2. Instead of finetuning all the network weights, it is possible to combine DreamBooth with Low Rank Adaptations (LoRA) [Hu et al. 2022; Ryu 2023], for a more memory-efficient alternative, by injecting learnable residual modules $\Delta W$ to each network weight matrix $W$. $\Delta W$ is a composition of low-rank matrices, i.e., $W + \Delta W = W + AB$ where $W \in \mathbb{R}^{n \times n}$, $A \in \mathbb{R}^{n \times r}$, $B \in \mathbb{R}^{r \times n}$, $r \ll n$, and only the added $\Delta W$ is being updated during training while model's original parameters $W$ stay frozen.

### 3.3 Problem Setup

Formally, the model is given $n$ ($n \leq 5$) reference images $X_{ref} := \{I_{ref}^k\}_{k=1}^n$, a target image $I_{tgt} \in \mathbb{R}^{H \times W \times 3}$ and its associated binary mask $M_{tgt} \in \{0, 1\}^{H \times W}$, in which 1 denotes the region to fill and 0 denotes the existing area in $I_{tgt}$. The model is expected to generate a harmonized image $I_{out} \in \mathbb{R}^{H \times W \times 3}$ whose pixels should stay as similar as possible to $I_{tgt}$ where $M_{tgt}$ equals 0, while staying faithful to the corresponding contents in $X_{ref}$ where $M_{tgt}$ equals 1. We assume there is enough overlap between $X_{ref}$ and $I_{tgt}$ such that a human could imagine a plausible $I_{out}$.

### 3.4 RealFill

This task is challenging for both geometry-based [Zhao et al. 2023a; Zhou et al. 2021] and reconstruction-based approaches [Mildenhall et al. 2020] because there are barely any geometric constraints
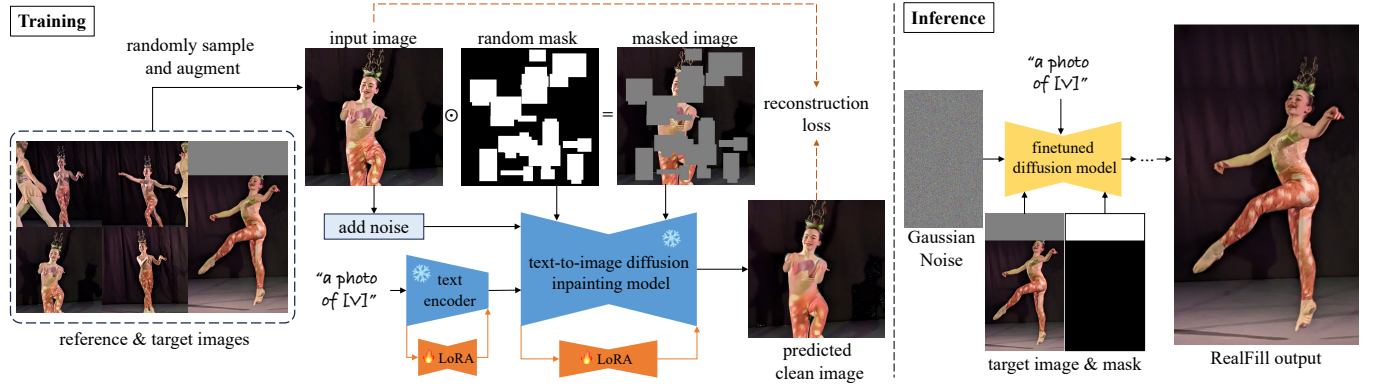
Fig. 2. **Training and inference pipelines of RealFill**. RealFill's inputs are a target image to be filled and a few reference images of the same scene. We first finetune LoRA weights of a pretrained inpainting diffusion model on the reference and target images (with random patches masked out). Then, we use the adapted model to fill the desired region of the target image, resulting in a faithful, high-quality output. For example, the girl's crown is recovered in the target image, despite the girl being in very different poses in the reference images.

between $\mathcal{X}_{ref}$ and $I_{tgt}$, there are only a few images available as inputs, and the reference images may have different styles, lighting conditions, and subject poses from the target. One alternative is to use a controllable inpainting or outpainting model, however, these models are either prompt-based [Adobe Inc. 2023; Rombach et al. 2022] or single-image object-driven [Yang et al. 2023], which makes them hard to use for recovering complex scene-level structure and details.

Therefore, we propose to first inject knowledge of the scene into a pretrained generative model by finetuning it on a set of reference images, then use the adapted model to generate $I_{out}$ conditioned on $I_{tgt}$ and $M_{tgt}$, such that it is aware of the scene's contents.

**Training**. Starting from a state-of-the-art T2I diffusion inpainting model [Rombach et al. 2022], we inject LoRA weights and finetune it on both $\mathcal{X}_{ref}$ and $I_{tgt}$ with randomly generated binary masks $m \in \{0, 1\}^{H \times W}$. The loss function is

$$\mathcal{L} = \mathbb{E}_{x,t,\epsilon,m} \|\epsilon_\theta(x_t, t, p, m, (1-m) \odot x) - \epsilon\|_2^2 \quad (3)$$

where $x \in \mathcal{X}_{ref} \cup \{I_{tgt}\}$, $p$ is a fixed language prompt, $\odot$ denotes the element-wise product and therefore $(1-m) \odot x$ is the masked clean image. For $I_{tgt}$, the loss is only calculated on the existing region, i.e., where $M_{tgt}$'s entry equals 0. Specifically, we use the open-sourced Stable Diffusion v2 inpainting model [Stability AI 2022] and inject LoRA layers into its text encoder and U-Net for finetuning. Following [Ruiz et al. 2023a], we fix $p$ to be a sentence containing a rare token, i.e., "a photo of [V]". For each training example, similar to [Suvorov et al. 2022], we generate multiple random rectangles and take either their union or the complement of the union to get the final random mask $m$. Fig. 2 illustrates the whole pipeline.

**Inference**. After training, we use the DDPM [Ho et al. 2020] sampler to generate an image $I_{gen}$, conditioning the model on $p$, $I_{tgt}$ and $M_{tgt}$. However, similar to the observation in [Zhu et al. 2023], we notice that the existing region in $I_{tgt}$ is distorted in $I_{gen}$. To resolve this, we first feather the mask $M_{tgt}$, then use it to alpha composite $I_{gen}$ and $I_{tgt}$, leading to the final $I_{out}$ with full recovery on the existing area and a smooth transition at the boundary of the generated region.

**Correspondence-Based Seed Selection**. The diffusion inference process is stochastic, i.e., the same inputs may produce any number of generated images depending on the random seeds of the sampling process. The generation quality can vary due to this stochasticity, thus requiring humans to select high-quality samples. While there is work to identify good samples from a collection of generated outputs [Samuel et al. 2023], this remains an open problem. In our case, the reference images actually provide a grounding signal for the true scene content, thus can be used to help identify high-quality outputs.

Specifically, we use the number of image feature correspondences between $I_{out}$ and $\mathcal{X}_{ref}$ as a metric to roughly quantify whether the result is faithful to the reference images. During inference, we first generate a batch of outputs, i.e., $\{I_{out}\}$, then extract a set of correspondences (e.g., using LoFTR [Sun et al. 2021]) between $\mathcal{X}_{ref}$ and the filled region of each $I_{out}$, and finally rank the generated results $\{I_{out}\}$ by the number of correspondences. This allows us to automatically filter generations to a small set of high-quality results.

## 4 EXPERIMENTS

**Qualitative Results**. In Figs. 3 and 4, we show that RealFill is able to convincingly outpaint and inpaint image content that is faithful to the reference images. Notably, it is able to handle dramatic differences in camera pose, lighting, defocus blur, image style, and even subject pose. This is because RealFill has both a good image prior (from the pretrained diffusion model) and knowledge of the scene (from finetuning on the input images).

**Evaluation Dataset**. Existing benchmarks for reference-based image completion [Zhou et al. 2021] primarily focus on inpainting small regions, and assume at most very minor changes between the reference and target images. To better evaluate our target use cases, we create our own dataset, *RealBench*. It consists of 33 scenes (23 outpainting and 10 inpainting), where each scene has a set of reference images $\mathcal{X}_{ref}$, a target image $I_{tgt}$ to fill, a binary mask $M_{tgt}$ indicating the missing region, and the ground-truth result $I_{gt}$. The number of reference images in each scene varies from 1 to 5.

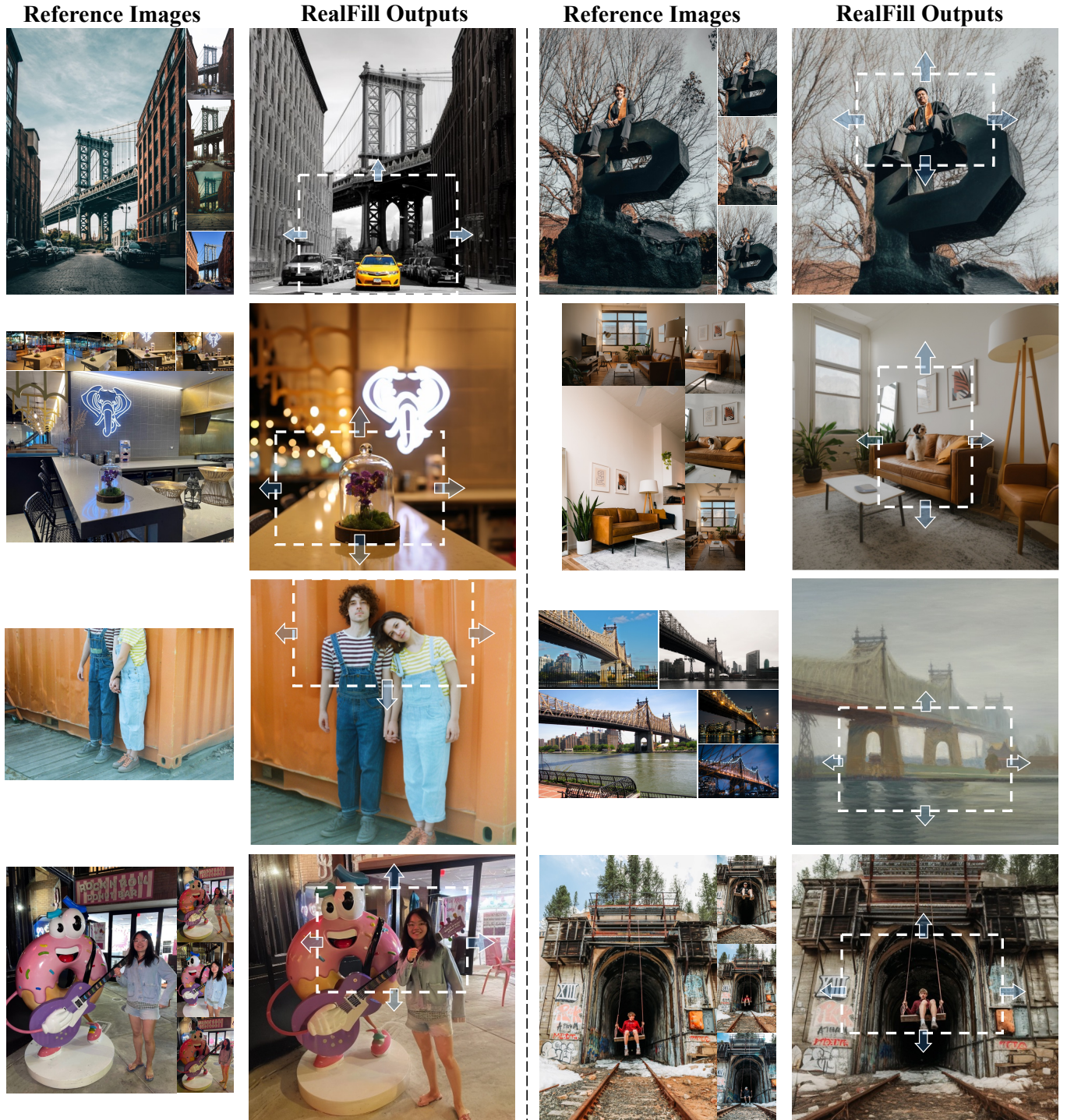Reference Images | RealFill Outputs | Reference Images | RealFill Outputs



Fig. 3. **Reference-based outpainting with RealFill**. Given the reference images on the left, RealFill outpaints the corresponding target images on the right. The region inside the white box is provided to the network as known pixels, and the region outside the white box is generated. RealFill produces high-quality images that are faithful to the references, even when there are dramatic differences between the references and targets such as changes in viewpoint, aperture, lighting, image style, and object motion.
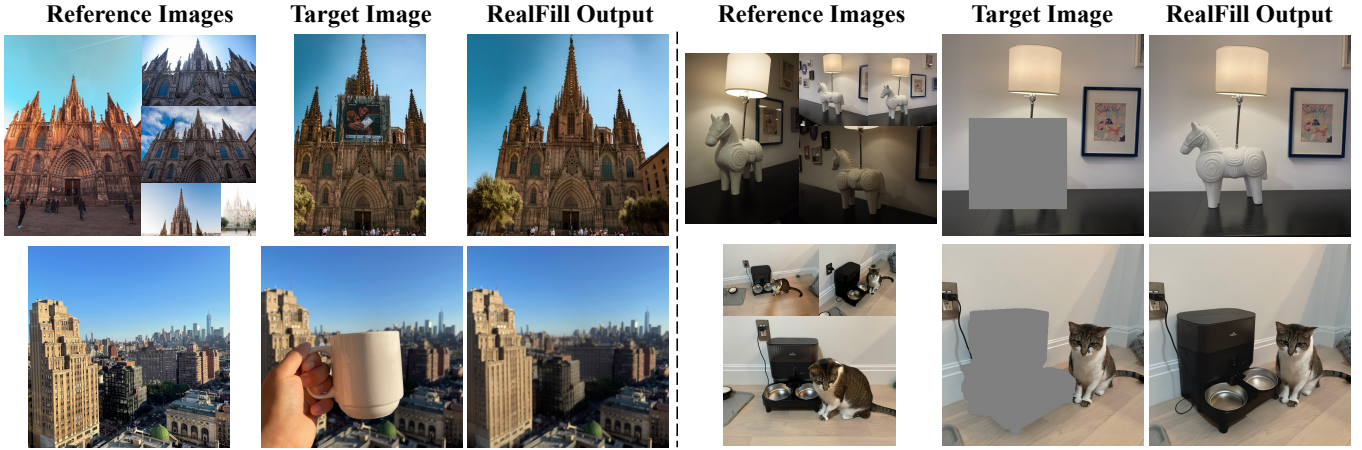
| Reference Images | Target Image | RealFill Output | Reference Images | Target Image | RealFill Output |



**Fig. 4. Reference-based inpainting with RealFill**. Given the references on the left, RealFill can not only remove undesired objects in the target image and reveal the occluded contents faithfully (left column), but also insert objects into the scene despite significant viewpoint changes between reference and target images (right column). In the bottom left example, the reference and target images have different defocus blurs. RealFill not only recovers the buildings behind the mug, but also keeps the same amount of blur as in the target image.

The dataset contains diverse, challenging scenarios with significant variations between the reference and target images, such as changes in viewpoint, defocus blur, lighting, style, and subject pose.

**Evaluation Metrics**. We use multiple metrics to evaluate the quality and fidelity of our model outputs. We compare the generated images with the ground-truth target image at multiple levels of image similarity, including PSNR, SSIM, LPIPS [Zhang et al. 2018] for low-level, DreamSim [Fu et al. 2023] for mid-level, and DINO [Caron et al. 2021], CLIP [Radford et al. 2021] for high-level. For low-level metrics, we only do calculation inside the filled-in region, i.e., where $M_{tgt}$ is 1. For high-level image similarity, we use the cosine distance between the full image embeddings from CLIP and DINO. DreamSim [Fu et al. 2023] is a mid-level similarity between two full images, emphasizing differences in image layouts, object poses, and semantic contents.

**Baseline Methods**. We compare to two groups of baselines: models that take reference image as conditioning input, i.e., TransFill [Zhou et al. 2021] and Paint-by-Example [Yang et al. 2023]; and prompt-based image filling approaches including Stable Diffusion Inpainting [Stability AI 2022] and Photoshop Generative Fill [Adobe Inc. 2023]. Since TransFill and Paint-by-Example can only use one reference image during inference, we randomly sample one from $X_{ref}$ as reference for each run of them. Choosing an appropriate prompt for prompt-based filling methods is a necessary component of getting a high quality result. So, for a fair comparison, instead of using a generic prompt like "a beautiful photo", for each scene, we manually design a long prompt that describe the true scene in detail with the help of ChatGPT [OpenAI 2023].

**Quantitative Comparison**. We quantitatively evaluate all methods on RealBench. For each approach, we report average metrics across all 33 scenes. Specifically, for generative methods, each scene's metric is itself computed from an average of 64 randomly generated samples (18 for Generative Fill due to Photoshop UI-only limits); for TransFill, the metrics are averaged over different choices of

**Table 1. Quantitative comparison of RealFill and baselines**. On Real-Bench, our evaluation set of 33 diverse challenging scenes, RealFill outperforms both prompt-based and reference-based baselines by a large margin on all types of metrics.

| | Method | low-level | | | mid-level | high-level | |
|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | DreamSim↓ | DINO↑ | CLIP↑ |
| prompt based | SD Inpaint | 10.63 | 0.282 | 0.605 | 0.213 | 0.831 | 0.874 |
| | Generative Fill | 10.92 | 0.311 | 0.598 | 0.212 | 0.851 | 0.898 |
| reference based | Paint-by-Example | 10.13 | 0.244 | 0.642 | 0.237 | 0.797 | 0.859 |
| | TransFill | 13.28 | 0.404 | 0.542 | 0.192 | 0.860 | 0.866 |
| | **RealFill (ours)** | **14.78** | **0.424** | **0.431** | **0.077** | **0.948** | **0.962** |

single reference image. As shown in Tab. 1, RealFill outperforms all baselines by a large margin across all metrics.

**Qualitative Comparison**. In Fig. 5, we present a visual comparison between RealFill and the baselines. We also show the ground-truth and input images for each example. In order to better highlight the regions which are being generated, we overlay a semi-transparent white mask on the ground-truth and output images, covering the known regions of the target image. RealFill not only generates high-quality images, but also more faithfully reproduces the scene than the baseline methods. Paint-by-Example relies on the CLIP embedding of the reference images as the condition. This poses a challenge when dealing with complex scenes or attempting to restore object details, since CLIP embeddings only capture high-level semantic information. Although the geometric-based TransFill has decent numbers in terms of low-level metrics like PSNR, the outputs have much lower quality due to the lack of a good image prior, especially when the scene structure has complex depth variants beyond a planar surface, which is hard for homography transformations to approximate. The generated results from Generative Fill are plausible on their own. However, because natural language is limited in conveying complex visual information, they often exhibit substantial deviations from the original scenes depicted in the references.
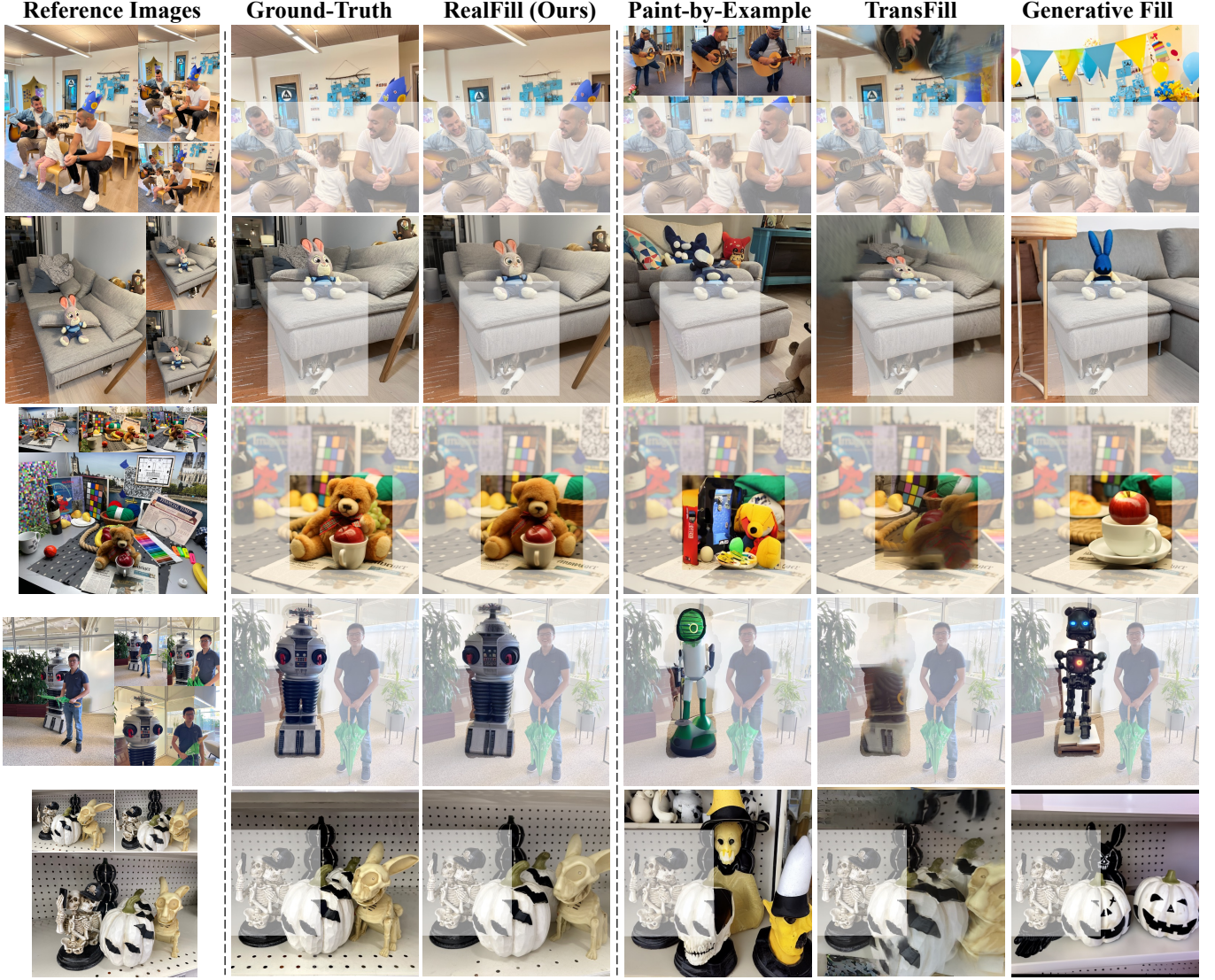
**Fig. 5. Qualitative comparison of RealFill and baselines**. Transparent white masks are overlaid on the unaltered known regions of the target images. Paint-by-Example loses fidelity with the reference images because it relies on CLIP embeddings, which only capture high-level semantic information. TransFill outputs low quality images due to the lack of a good image prior and the limitations of its geometry-based pipeline. While Generative Fill produces plausible results, they are inconsistent with the reference images because prompts have limited expressiveness. In contrast, RealFill generates high-quality results that have high fidelity with respect to the reference images.

Table 2. **Effect of correspondence-based seed selection**. It helps RealFill output higher quality results, i.e., filtering out samples with fewer matches results in better quantitative scores.

| Filtering Rate | PSNR↑ | SSIM↑ | LPIPS↓ | DreamSim↓ | DINO↑ | CLIP↑ |
|---|---|---|---|---|---|---|
| 0% | 14.78 | 0.424 | 0.431 | 0.077 | 0.948 | 0.962 |
| 25% | 15.01 | 0.427 | 0.421 | 0.066 | 0.955 | 0.967 |
| 50% | 15.05 | 0.427 | 0.418 | 0.063 | 0.958 | 0.969 |
| **75%** | **15.10** | **0.427** | **0.417** | **0.060** | **0.961** | **0.970** |

Table 3. **User study results**. Here we show the preference percentages for the most realistic and most faithful image completions across 58 scenes by 44 participants. RealFill significantly surpasses baselines, especially on the faithfulness criterion.

| Method | Most Realistic ↑ | Most Faithful ↑ |
|---|---|---|
| TransFill [Zhou et al. 2021] | 2.0% | 3.7% |
| Paint-by-Example [Yang et al. 2023] | 11.0% | 1.9% |
| Generative Fill [Adobe Inc. 2023] | 23.4% | 7.2% |
| **RealFill (ours)** | **63.7%** | **87.2%** |

**Reference Images**   **Target Image**   **Ground-Truth**   **RealFill Outputs**

# of matches: 1111    # of matches: 692    # of matches: 83

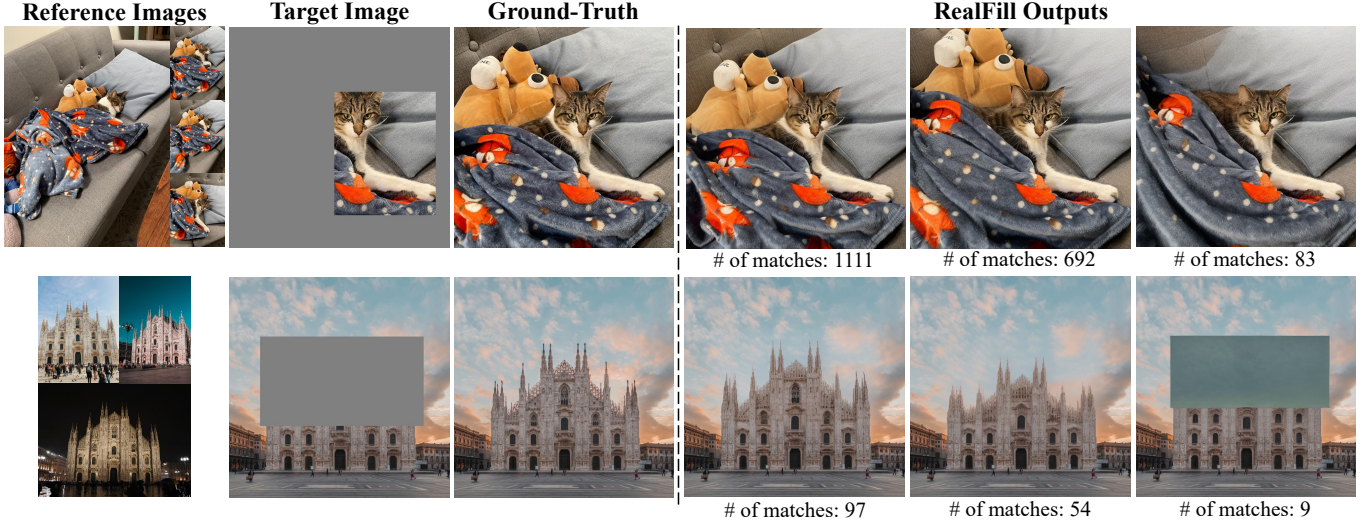# of matches: 97    # of matches: 54    # of matches: 9

Fig. 6. **Correspondence-based seed selection**. Given the reference images on the left, we show multiple RealFill outputs on the right along with the number of matched key points. We can see that fewer matches correlate with lower-quality outputs that are more divergent from the ground-truth.
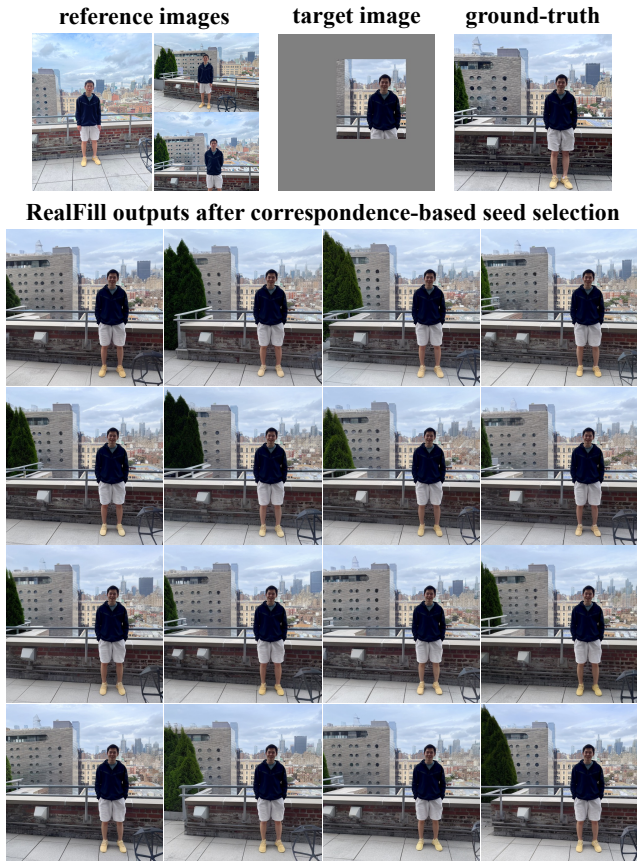
**reference images**   **target image**   **ground-truth**



**RealFill outputs after correspondence-based seed selection**

Fig. 7. **Top 16 RealFill outputs after correspondence-based selection**. Model first generates 64 examples in a batch with different sampled noises, then the top 16 images are automatically selected based on predicted correspondences. The outputs are mostly high-quality, so it's easy for users to pick the final image based on their own preference.

**Correspondence-Based Seed Selection**. We evaluate the effect of our proposed correspondence-based seed selection described in Sec. 3.4. To measure the correlation between our seed selection mechanism and high-quality results, for each scene, we rank Real-Fill's outputs $\{I_{out}\}$ according to the number of matched keypoints, and then filter out a certain percent of the lowest-ranked samples. We then average the evaluation metrics only across the remaining samples. Higher filtering rates like 75% are quantitatively better than no filtering (Tab. 2). In Fig 6, we show multiple RealFill outputs with the corresponding number of matched keypoints. These demonstrate a clear trend, where fewer matches usually indicate lower-quality results.

Therefore, we followed such strategy to select RealFill outputs in Figs. 1 to 5: for each scene, model first generates 64 examples in a batch with different sampled noises, then we only keep the top 16 images based on predicted correspondences, and manually pick one from them. For fair comparison, we also manually pick the best outputs for each baseline as in Figs. 1 and 5. As shown in Fig. 7, after correspondence-based selection, the outputs are mostly high-quality with small variations. Users can make choices based on their own preference, therefore the involved human labor is very light.

**User study**. In our study, 44 users evaluated the realism and faithfulness of image completions from four different methods across 58 scenes. Participants first review an incomplete target image alongside the outputs of these methods, choosing the one they found most realistic. Then, with reference images provided, they select the most faithful completion. The study includes a mix of Real-Bench and 25 additional challenging scenes. All method outputs are randomly sampled to avoid bias. Overall, 2552 votes per criterion were collected. The results, as summarized in Tab. 3, show Real-Fill's superiority compared to the baseline methods especially in faithfulness.

**Target Image**  **TransFill Output**  **Ground-Truth**
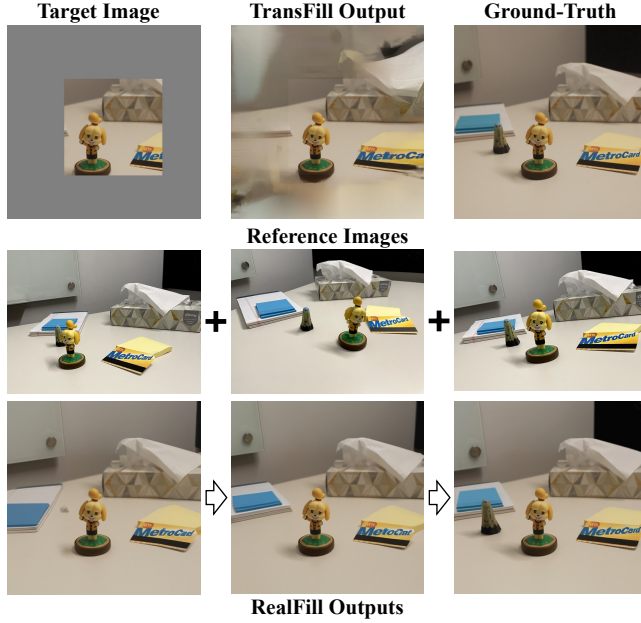
**Reference Images**

**RealFill Outputs**

Fig. 8. **Influence of reference image choice on RealFill**. For the given target image, we show different RealFill outputs by gradually increasing number of reference images from left to right. We can see outputs are getting better when having more references, e.g., the ghost gets recovered when having all three reference images. Note that all RealFill outputs significantly outperform the TransFill baseline.

## 5  DISCUSSION

**How does reference image choice affect RealFill?** Empirically, when there are more reference images, or when the references have smaller variations from the target in terms of viewpoint and lighting, RealFill gives better results, as shown in Fig. 8.

**Would other baselines work?**

*Image Stitching*. It is possible to stitch the reference and target images together using correspondences. However, we find that even strong commercial image stitching software doesn't work when there are dramatic lighting changes or object motion. Taking the two scenes in Fig. 9 for example, multiple commercial software solutions produce no output, asserting that there are insufficient correspondences. In contrast, RealFill faithfully recovers these scenes.

*DreamBooth*. Instead of finetuning an inpainting model, an alternative is to finetune a standard Stable Diffusion model on the reference images, i.e., DreamBooth, then use the finetuned T2I model to inpaint the target image [Lugmayr et al. 2022], as implemented in the popular Diffusers library [von Platen et al. 2022][1]. However, because this model is never trained with a masked prediction objective, it is much worse than RealFill, as shown in Fig. 10.

**What makes RealFill work?** To explore why RealFill leads to strong results, especially on complex scenes, we make the following two hypotheses:

*RealFill relates multiple elements in a scene*. If we make the conditioning image a blank canvas during inference, i.e., all entries of

[1]Diffusers' Stable Diffusion inpainting pipeline code.

**Reference Images**  **Target Image**  **RealFill Outputs**

Fig. 9. **RealFill vs. image stitching**. Commercial image stitching softwares fail to produce any outputs when there are significant variations between reference and target images, such as disparate lighting conditions. In contrast, RealFill excels by producing accurate and high-quality results. It effectively recovers elements like rooftop tanks and balloons, maintaining fidelity even under varying lighting scenarios between the compared images.
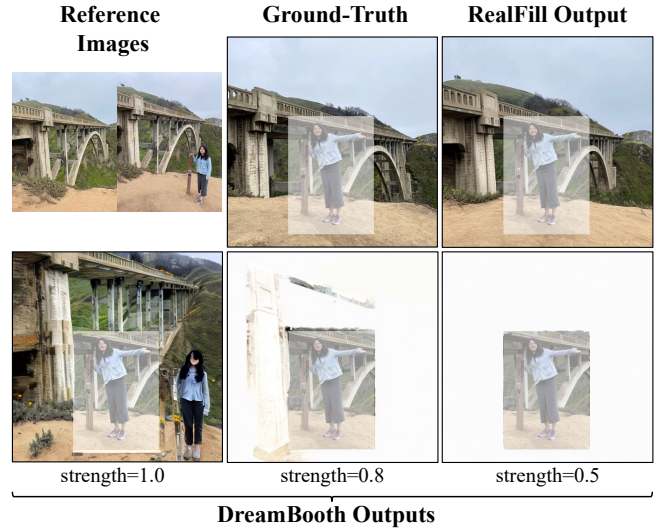
**Reference Images**  **Ground-Truth**  **RealFill Output**

strength=1.0    strength=0.8    strength=0.5

**DreamBooth Outputs**

Fig. 10. **RealFill vs. DreamBooth**. Finetuning a standard Stable Diffusion model on reference images and using it to fill missing regions, leads to drastically worse results compared to RealFill. We show samples for various levels of the strength hyper-parameter.

$M_{tgt}$ equal 1, we can see in Fig. 11 that the finetuned model is able to generate multiple scene variants with different structures, e.g., removing the foreground or background object, or manipulating the object layouts. This suggest that RealFill may understand the scene composition.

*RealFill captures correspondences among input images*. Even if the reference and target images do not depict the same scene, the finetuned model is still able to fuse the corresponding contents of

**Fig. 11.** RealFill is able to generate multiple scene variants when conditioned on a blank image as input, e.g., people are added or removed in the first and second rows. This suggests that the finetuned model can relate elements inside the scene in a compositional manner.
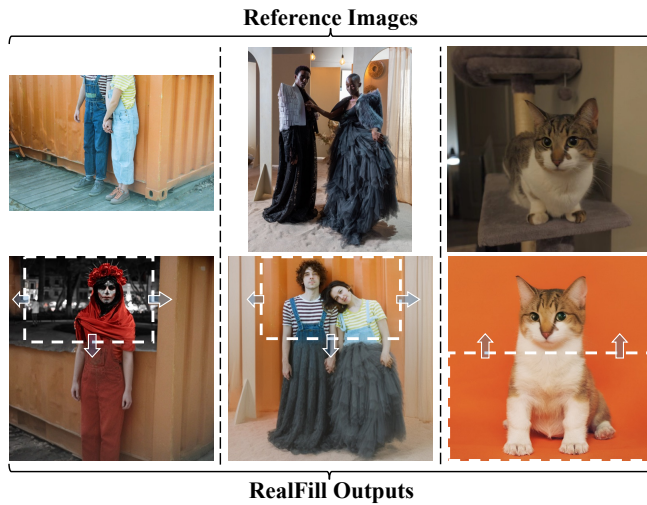


**Fig. 12.** When the reference and target images do not depict the same scene, the finetuned model is still able to fuse the reference contents into the target image in a semantically-reasonable way, suggesting that it captures both real or invented correspondences between input images.

the reference images into the target area seamlessly, as shown in Fig. 12. This suggests that RealFill is able to capture and utilize real or invented correspondences between reference and target images to do generation. Previous works [Luo et al. 2023; Tang et al. 2023;



**Fig. 13.** **Failure cases of RealFill**. (Top) RealFill fails to recover the precise 3D scene structure, e.g., the output husky plush has different pose compared to the reference; (Bottom) RealFill fails to handle cases that are also challenging for the base T2I diffusion model, e.g., the words on the store sign are wrongly spelled.

Zhang et al. 2023a] also found similar emergent correspondence inside pretrained diffusion models.

**Limitations**. Because RealFill requires a gradient-based finetuning process on input images, it is relatively slow. Empirically, we also find that, when there is a large viewpoint change between reference and target images, RealFill fails to recover the 3D scene faithfully, especially when there's only a single reference image, as in Fig. 13 top. In addition, because RealFill relies on the base pretrained model's image prior, it also fails to handle cases that are challenging for the base model, e.g., Stable Diffusion is not good at generating fine details, such as text, human faces, or body parts, as in Fig. 13 bottom.

Lastly, similar to the user study in Tab. 3, we conducted another study on RealBench where participants compare randomly sampled RealFill output vs. ground-truth image for each scene. Among the collected 396 pairwise comparisons, RealFill only gets 23.7% of the votes for realism and 22.0% for faithfulness, vs. ground-truth's 76.3% and 78.0% respectively. This is reasonable because it's easy for human to spot artifacts especially for side-by-side comparison, but it also shows that more future improvements are needed to make RealFill achieve perfect authentic image completion.

**Societal Impact**. This research aims to create a tool that can help users express their creativity and improve the quality of their personal photographs through image generation. However, advanced image generation methods can have complex impacts on society. Our proposed method inherits some of the concerns that are associated with this class of technology, such as the potential to alter sensitive personal characteristics. The open-source pretrained model that we use in our work, Stable Diffusion, exhibits some of these concerns. However, we have not found any evidence that our method is more likely to produce biased or harmful content than previous work. Despite these findings, it is important to continue investigating the potential risks of image generation technology. Future research should focus on developing methods to mitigate bias and harmful

content, and to ensure that image generation tools are used in a responsible manner.

## 6  CONCLUSION

In this work, we introduce the problem of *Authentic Image Completion*, where given a few reference images, we intend to complete some missing regions of a target image with the content that "*should* have been there" rather than "what *could* have been there". To tackle this problem, we proposed a simple yet effective approach called RealFill, which first finetunes a T2I inpainting diffusion model on the reference and target images, and then uses the adapted model to fill the missing regions. We show that RealFill produces high-quality image completions that are faithful to the content in the reference images, even when there are large differences between reference and target images such as viewpoint, aperture, lighting, image style, and object pose.

## ACKNOWLEDGMENTS

## A  IMPLEMENTATION DETAILS

### A.1  RealFill

For each scene, we finetune the Stable Diffusion inpainting model [Stability AI 2022] for 2,000 iterations with a batch size of 16 on a single NVIDIA A100 GPU with LoRA rank 8. With a probability of 0.1, we randomly dropout prompt $p$, mask $m$ and LoRA layers independently during training. The learning rate is set to 2e-4 for the U-Net and 4e-5 for the text encoder. The whole finetuning process takes around one hour but usually 20 minutes would give pretty good results already in many cases. Note that these hyper-parameters could be further tuned for each scene to get better performance, e.g., some scenes converge more quickly may overfit if trained for too long. However, for the sake of fair comparison, we use a constant set of hyper-parameters for all results shown in the paper. During inference, we use DDPM [Ho et al. 2020] sampler with step 200 and guidance weight 1.0, i.e., without classifier-free guidance.

### A.2  User Study

In the user study, participants are asked to evaluated the realism and faithfulness of image completions across 58 scenes from four different methods including RealFill, Paint-by-Example [Yang et al. 2023], TransFill [Zhou et al. 2021], and Photoshop Generative Fill [Adobe Inc. 2023]. All method outputs are randomly sampled without human intervention. For each scene, the placement organization of different method outputs are also randomly shuffled to avoid user bias.
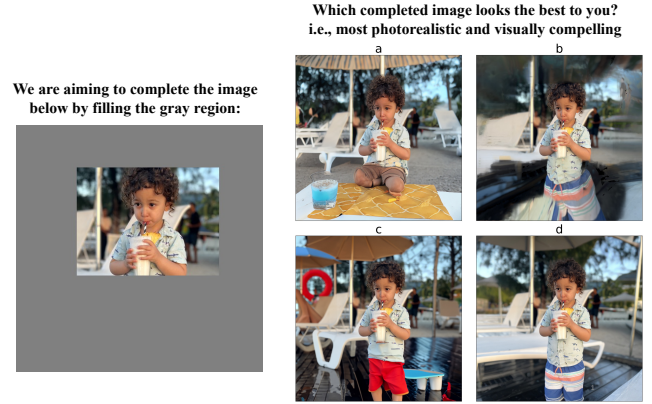


Fig. 14. An example of the realism question in the user study. Users first review an incomplete target image alongside the shuffled outputs of four methods, then choose the one they find most realistic.
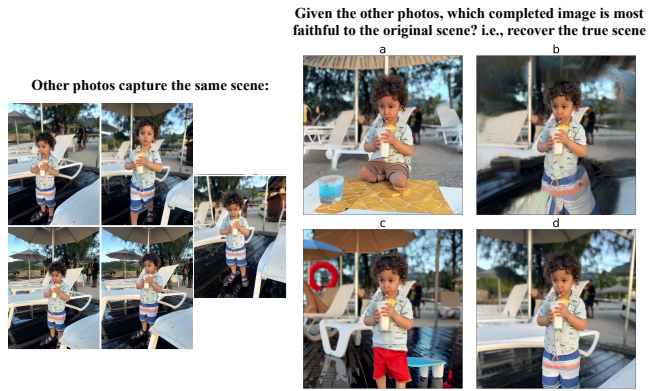


Fig. 15. An example of the faithfulness question in the user study. After the user answered the realism question, the reference images are provided. Then, they are asked to select the most faithful completion.

Specifically, participants first review an incomplete target image alongside the outputs of these methods, choosing the one they found most realistic, as shown in Fig. 14. Then, with reference images provided, they select the most faithful completion, as shown in Fig. 15.

## REFERENCES

Adobe Inc. 2023. Adobe Photoshop.  https://www.adobe.com/products/photoshop.html

Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. 2023. Break-A-Scene: Extracting Multiple Concepts from a Single Image. *ArXiv preprint* abs/2305.16311 (2023).  https://arxiv.org/abs/2305.16311

Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. Patch-Match: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.* 28, 3 (2009), 24.

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

Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2023. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18392–18402.

Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging Properties in Self-Supervised Vision Transformers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 9630–9640. https://doi.org/10.1109/ICCV48922.2021.00951

Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. 2023. Muse: Text-to-image generation via masked generative transformers. *ArXiv preprint* abs/2301.00704 (2023). https://arxiv.org/abs/2301.00704

Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. 2022. MaskGIT: Masked Generative Image Transformer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 11305–11315. https://doi.org/10.1109/CVPR52688.2022.01103

Wenhu Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W Cohen. 2023. Subject-driven text-to-image generation via apprenticeship learning. *ArXiv preprint* abs/2304.00186 (2023). https://arxiv.org/abs/2304.00186

Antonio Criminisi, Patrick Perez, and Kentaro Toyama. 2003. Object removal by exemplar-based inpainting. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, Vol. 2. IEEE, II–II.

Prafulla Dhariwal and Alexander Quinn Nichol. 2021. Diffusion Models Beat GANs on Image Synthesis. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 8780–8794. https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html

Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. 2023. DreamSim: Learning New Dimensions of Human Visual Similarity using Synthetic Data. *ArXiv preprint* abs/2306.09344 (2023). https://arxiv.org/abs/2306.09344

Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *ArXiv preprint* abs/2208.01618 (2022). https://arxiv.org/abs/2208.01618

James Hays and Alexei A Efros. 2007. Scene completion using millions of photographs. *ACM Transactions on Graphics (ToG)* 26, 3 (2007), 4–es.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net. https://openreview.net/forum?id=nZeVKeeFYf9

Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. 2017. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–14.

Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. 2023. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6007–6017.

Soo Ye Kim, Kfir Aberman, Nori Kanazawa, Rahul Garg, Neal Wadhwa, Huiwen Chang, Nikhil Karnad, Munchurl Kim, and Orly Liba. 2022. Zoom-to-Inpaint: Image Inpainting with High-Frequency Details. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2022, New Orleans, LA, USA, June 19-20, 2022*. IEEE, 476–486. https://doi.org/10.1109/CVPRW56347.2022.00063

Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. 2018. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European conference on computer vision (ECCV)*. 85–100.

Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 2023. Zero-1-to-3: Zero-shot One Image to 3D Object. arXiv:2303.11328 [cs.CV]

Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. 2022. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11461–11471.

Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. 2023. Diffusion Hyperfeatures: Searching Through Time and Space for Semantic Correspondence. *ArXiv preprint* abs/2305.14334 (2023). https://arxiv.org/abs/2305.14334

Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.

Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. 2023. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *ArXiv preprint* abs/2302.08453 (2023). https://arxiv.org/abs/2302.08453

OpenAI. 2023. ChatGPT: Optimizing Language Models for Dialogue. https://openai.com/blog/chatgpt. Accessed: November 17, 2023.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. http://proceedings.mlr.press/v139/radford21a.html

Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Nataniel Ruiz, Ben Mildenhall, Shiran Zada, Kfir Aberman, Michael Rubinstein, Jonathan Barron, Yuanzhen Li, and Varun Jampani. 2023. DreamBooth3D: Subject-Driven Text-to-3D Generation. *ArXiv preprint* abs/2303.13508 (2023). https://arxiv.org/abs/2303.13508

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *ArXiv preprint* abs/2204.06125 (2022). https://arxiv.org/abs/2204.06125

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.

Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023a. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22500–22510.

Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, Michael Rubinstein, and Kfir Aberman. 2023b. HyperDreamBooth: HyperNetworks for Fast Personalization of Text-to-Image Models. *ArXiv preprint* abs/2307.06949 (2023). https://arxiv.org/abs/2307.06949

Simo Ryu. 2023. Low-rank Adaptation for Fast Text-to-Image Diffusion Fine-tuning. https://github.com/cloneofsimo/lora.

Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* 35 (2022), 36479–36494.

Dvir Samuel, Rami Ben-Ari, Simon Raviv, Nir Darshan, and Gal Chechik. 2023. It is all about where you start: Text-to-image generation with seed selection. *ArXiv preprint* abs/2304.14530 (2023). https://arxiv.org/abs/2304.14530

Qi Shan, Brian Curless, Yasutaka Furukawa, Carlos Hernandez, and Steven M Seitz. 2014. Photo uncrop. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI 13*. Springer, 16–31.

Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. 2023. StyleDrop: Text-to-Image Generation in Any Style. *ArXiv preprint* abs/2306.00983 (2023). https://arxiv.org/abs/2306.00983

Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2021. Score-Based Generative Modeling through Stochastic Differential Equations. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=PxTIG12RRHS

Stability AI. 2022. Stable-Diffusion-2-Inpainting. https://huggingface.co/stabilityai/stable-diffusion-2-inpainting.

Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. 2021. LoFTR: Detector-Free Local Feature Matching With Transformers. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 8922–8931. https://doi.org/10.1109/CVPR46437.2021.00881

Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. 2022. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 2149–2159.

Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. 2023. Emergent Correspondence from Image Diffusion. In *Thirty-seventh Conference on Neural Information Processing Systems*. https://openreview.net/forum?id=ypOiXjdfnU

Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. 2022. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers.

Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. 2023. *P+*: Extended Textual Conditioning in Text-to-Image Generation. *ArXiv preprint* abs/2303.09522 (2023). https://arxiv.org/abs/2303.09522

Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. 2023. Imagen

editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18359–18369.

Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. 2022. Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation. *ArXiv preprint* abs/2212.11565 (2022). https://arxiv.org/abs/2212.11565

Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. 2023. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18381–18391.

Raymond A. Yeh, Chen Chen, Teck-Yian Lim, Alexander G. Schwing, Mark Hasegawa-Johnson, and Minh N. Do. 2017. Semantic Image Inpainting with Deep Generative Models. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 6882–6890. https://doi.org/10.1109/CVPR.2017.728

Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. 2018. Generative Image Inpainting With Contextual Attention. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 5505–5514. https://doi.org/10.1109/CVPR.2018.00577

Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. 2023a. A Tale of Two Features: Stable Diffusion Complements DINO for Zero-Shot Semantic Correspondence. *ArXiv preprint* abs/2305.15347 (2023). https://arxiv.org/abs/2305.15347

Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023b. Adding Conditional Control to Text-to-Image Diffusion Models.

Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 586–595. https://doi.org/10.1109/CVPR.2018.00068

Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. 2023b. Unleashing Text-to-Image Diffusion Models for Visual Perception. *ICCV* (2023).

Yunhan Zhao, Connelly Barnes, Yuqian Zhou, Eli Shechtman, Sohrab Amirghodsi, and Charless Fowlkes. 2023a. Geofill: Reference-based image inpainting with better geometric understanding. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1776–1786.

Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. 2019. Pluralistic Image Completion. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 1438–1447. https://doi.org/10.1109/CVPR.2019.00153

Yuqian Zhou, Connelly Barnes, Eli Shechtman, and Sohrab Amirghodsi. 2021. Trans-Fill: Reference-Guided Image Inpainting by Merging Multiple Color and Spatial Transformations. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2266–2276. https://doi.org/10.1109/CVPR46437.2021.00230

Zixin Zhu, Xuelu Feng, Dongdong Chen, Jianmin Bao, Le Wang, Yinpeng Chen, Lu Yuan, and Gang Hua. 2023. Designing a Better Asymmetric VQGAN for StableDiffusion. *ArXiv preprint* abs/2306.04632 (2023). https://arxiv.org/abs/2306.04632