

CorrFill: Enhancing Faithfulness in Reference-based Inpainting with Correspondence Guidance in Diffusion Models

Kuan-Hung Liu¹ Cheng-Kun Yang^{2*} Min-Hung Chen³ Yu-Lun Liu¹ Yen-Yu Lin¹

¹National Yang Ming Chiao Tung University ²National Taiwan University ³NVIDIA

<https://corrfill.github.io/>

Abstract

In the task of reference-based image inpainting, an additional reference image is provided to restore a damaged target image to its original state. The advancement of diffusion models, particularly Stable Diffusion, allows for simple formulations in this task. However, existing diffusion-based methods often lack explicit constraints on the correlation between the reference and damaged images, resulting in lower faithfulness to the reference images in the inpainting results. In this work, we propose CorrFill, a training-free module designed to enhance the awareness of geometric correlations between the reference and target images. This enhancement is achieved by guiding the inpainting process with correspondence constraints estimated during inpainting, utilizing attention masking in self-attention layers and an objective function to update the input tensor according to the constraints. Experimental results demonstrate that CorrFill significantly enhances the performance of multiple baseline diffusion-based methods, including state-of-the-art approaches, by emphasizing faithfulness to the reference images.

1. Introduction

Image inpainting aims to restore damaged regions of a target image. This task is inherently ill-posed, as any plausible outcome could be considered valid. Consequently, general image inpainting approaches are insufficient for faithfully recovering the original content of the images. To address this issue, reference-based image inpainting introduces supplementary images, known as reference images, to guide the recovery process for damaged regions [15]. These reference images can be photographs of the same scene with the target image, taken from different viewpoints or at different time slots. With the guidance of reference images, it becomes more practical to restore the target image to its original state.

Denosing diffusion probabilistic models [9] excel as

*Now at MediaTek Inc., Taiwan.

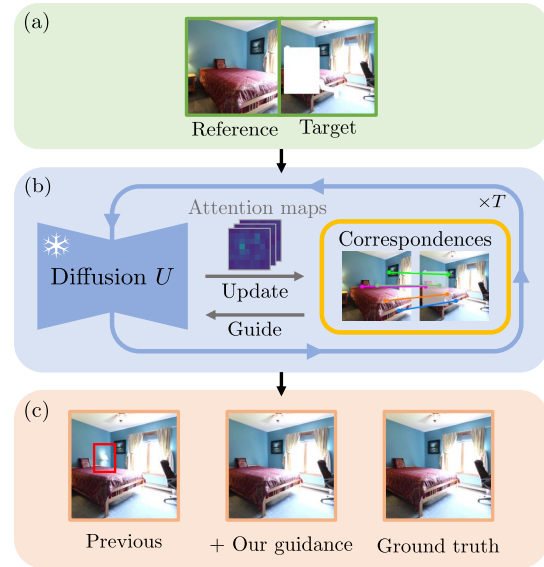


Figure 1. **Overview.** (a) The reference and target images are stitched side by side, serving as inputs to the model. (b) Reference-based inpainting using an inpainting fine-tuned Stable Diffusion [19] that employs our training-free correspondence guidance. (c) Our method captures more reliable correlations between references and targets than previous methods [2], thereby avoiding incorrect geometry and unwanted objects.

generative models, producing high-quality and diverse images [5], and showing significant potential in reference-based inpainting [2, 22, 24–26]. Existing diffusion-based methods for reference-based inpainting [24, 25] focus on training or fine-tuning an image-conditioned model to fill damaged regions based on reference images. However, they lack direct awareness of the relationships between targets and references, which is crucial for earlier approaches based on geometry matching [30, 32]. Without this awareness, diffusion models merely conditioned on reference images fail to ensure correct reference-target geometric correlation, leading to inpainting results that do not fully adhere to the content of the references, thus losing faithfulness. As shown in Figure 1, methods lacking direct reference-target aware-

ness suffer from unwanted objects in the generated results and can lead to incorrect scene layouts, or geometry as well.

In this work, we propose **CorrFill**, a plug-in module for reference-based inpainting diffusion models. It estimates image correspondences between targets and references, guiding the inpainting process with these correspondences as constraints, thereby preserving reference-target geometric relationships. To condition the inpainting process on the reference image, we stitch the reference and the target into a single image as input, allowing the self-attention layers in the diffusion models to attend across the reference and the target [2]. As demonstrated in Figure 1, our method collects self-attention scores at each denoising step, computes the correspondence, and guides the subsequent denoising step using this correspondence. It is worth noting that, in the context of reference-based inpainting, the application of existing image correspondence methods is hindered by the presence of damaged regions within the target image, rendering these methods ineffective for obtaining accurate correspondences.

Based on the observation that self-attention scores of inpainting diffusion models on the stitched image present the primitive approximations of correspondence [2] even in the damaged regions, we propose utilizing the correspondences derived from attention scores as the explicit constraints for guidance. As illustrated in Figure 1, the derived correspondence approximation is updated with the newly produced attention scores at each iteration, which is then used to guide the next iteration of denoising. This cyclic interaction between correspondence approximation and inpainting enables joint improvement, progressively refining the inpainting process to achieve a more faithful result. Furthermore, this method does not introduce additional learnable modules, making it a general strategy for reference-based inpainting diffusion models to improve the faithfulness to the reference images.

Our key contributions are summarized as follows: First, we propose CorrFill, a plug-in module for diffusion models, which utilizes correspondences as the explicit constraints to enhance the faithfulness to the reference images in reference-based inpainting. Second, the cyclic enhancement strategy employed in CorrFill facilitates the derivation of correspondence approximations for damaged image pairs without the need for additional training. This is achieved through the joint refinement of the inpainting process and the correspondence approximations. Third, our approach increases the performances of multiple diffusion model-based approaches on datasets RealEstate10K and MegaDepth.

2. Related Work

Reference-based Image Inpainting. Reference-based image inpainting aims to restore images to their original

states with additional knowledge provided by reference images. Geometry-based methods [30, 32] rely on geometric estimations, resulting in complex pipelines and a tendency for error propagation, particularly when dealing with large damaged regions and insufficient overlapping areas.

The profound capability of the diffusion models [9] in image generation reveals their potential to perform reference-based inpainting without the need for complex pipelines. Specifically, Stable Diffusion [19], notable for its capability of high-quality text-to-image generation, spurs active development of various downstream applications, including image-conditioned generation. Excelling in controllable generations with fine-grained guidance, image-conditioned variants of Stable Diffusion [14, 26, 29] demonstrate significant potential for the reference-based image inpainting task. For instance, [2, 22, 25, 26] address reference-based inpainting in an end-to-end manner based on Stable Diffusion.

Yang *et al.* [25] retrain a Stable Diffusion model to condition the generation process on the CLIP embedding [17] of the reference image. LeftRefill [2] employs prompt tuning techniques on a pre-trained text-to-image model, Stable Diffusion Inpainting, to avoid exhaustive training of the entire diffusion model. This method enables a text-to-image model to work with the image condition by concatenating the reference image and the target image side by side. Consequently, it allows mutual attention across the reference and target images via self-attention inside the diffusion U-Net.

By adopting diffusion models conditioned on reference images, these approaches circumvent the need for complex pipelines. However, merely training a diffusion model to generate with the reference condition is insufficient to capture the correct reference-target correlation, potentially leading to inconsistent results with the reference image. On the contrary, our CorrFill achieves faithful inpainting by guiding the generation process through the correspondences between reference and target images.

Diffusion Models Reliable Generation. A series of work [1, 3, 4, 6, 7, 13, 18, 20] enhances the controllability of pre-trained text-guided diffusion models by text prompts, semantic maps, layouts, or other factors, without model retraining or fine-tuning. For instance, Balaji *et al.* [1] control the locations of generated objects by introducing cross-attention masks based on the semantic masks and corresponding text tokens, thereby encouraging semantic attributes appearing at specified image patches. Manukyan *et al.* [13] propose an inpainting approach that ensures faithfulness to the text prompt by reducing the self-attention scores of image tokens unrelated to the text prompt. Furthermore, they design an objective function to enhance the effect of the text prompt on the damaged regions by optimizing the latent input of the model using the gradients

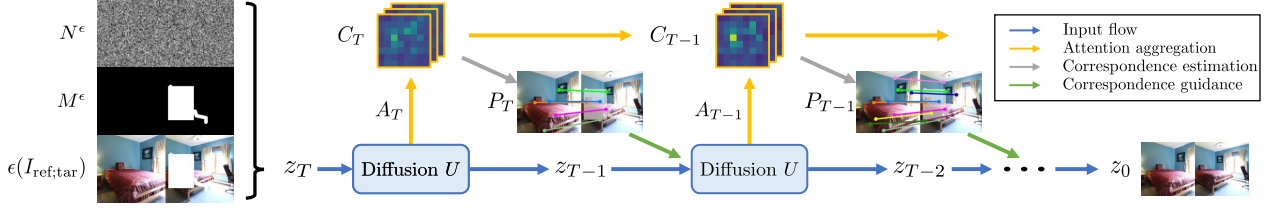


Figure 2. **Approach overview.** CorrFill jointly guides the inpainting and refines the estimated correspondences at each denoising step. The noise tensor N^ϵ , the downscaled mask M^ϵ and the encoded stitched image $\epsilon(I_{ref;tar})$ are concatenated into input latent tensor z_T . For each denoising step, the self-attention scores from the diffusion model are aggregated into a matching map C_t , and the correspondence P_t are computed from C_t , where P_t are used to guide the subsequent denoising step. For visual clarity, we use the real images to picture $\epsilon(I_{ref;tar})$ and z_0 .

of the objective function. Without additional learnable tokens or modules, our proposed CorrFill incorporates correspondence into the inpainting process using attention masks and latent tensor optimization, enabling training-free fine-grained control of diffusion models.

Diffusion Models and Correspondence. The powerful image priors of pre-trained diffusion models make them foundational models for numerous applications beyond image generation, such as image correspondence. Luo *et al.* [12] calculate semantic correspondence by aggregating Stable Diffusion features across different network layers and diffusion timesteps using a lightweight network. Zhang *et al.* [28] demonstrates the capability of zero-shot semantic correspondence inherent in Stable Diffusion features. They combine features from two vision foundation models, DINOv2 [16] and Stable Diffusion, and perform nearest neighbor searches based on the features to establish semantic correspondence. While these approaches estimate correspondence using features of intact images, we fully explore the potential of correspondence estimation with damaged inputs using pre-trained inpainting diffusion models.

3. Proposed Method

This section presents the proposed method, CorrFill, a correspondence-guided module for reference-based inpainting. Firstly, a method overview is provided. Then, we elaborate correspondence construction and refinement based on the reference-target attention scores. Finally, we present a joint process where correspondence estimation and guided inpainting are alternately performed to facilitate each other.

3.1. Overview

Reference-based image inpainting involves a reference image $I_{ref} \in \mathbb{R}^{h \times w \times 3}$ and a target image $I_{tar} \in \mathbb{R}^{h \times w \times 3}$ with damaged regions indicated by a binary mask $M \in \{0, 1\}^{h \times w}$. As depicted in Figure 2, our proposed method aims to restore the damaged regions of I_{tar} by referring to I_{ref} .

For ease of cross-image attention, we follow the practice used in [2] and horizontally stitch the reference and

target images to yield $I_{ref;tar} \in \mathbb{R}^{h \times 2w \times 3}$. CorrFill is developed based on pre-trained latent diffusion models [19]. To work in the latent space, the stitched image is encoded into $\epsilon(I_{ref;tar}) \in \mathbb{R}^{h' \times 2w' \times d}$, where $\epsilon(\cdot)$ is a variational autoencoder [10] and d is the dimension of the latent space. The image latent $\epsilon(I_{ref;tar})$ is then concatenated with the noise latent $N^\epsilon \in \mathbb{R}^{h' \times 2w' \times d}$ and the resized input mask $M^\epsilon \in \{0, 1\}^{h' \times 2w'}$, forming the input latent tensor $z_T \in \mathbb{R}^{h' \times 2w' \times (2d+1)}$ to the diffusion module.

For each denoising step t , it is carried out by a diffusion U-Net U , which takes the latent tensor z_t and correspondence $P_{t+1} \in [0, 1]^{h' \times w' \times 2}$ computed in the previous step as input and produces z_{t-1} via noise estimation. To compute correspondence, we utilize the self-attention maps produced in the denoising process. During denoising, the self-attention map $A_t \in \mathbb{R}^{(h' \times 2w') \times (h' \times 2w')}$ is computed and represents the patch-wise similarity in the stitched image $I_{ref;tar}$ at step t . We compile a matching map $C_t \in \mathbb{R}^{h' \times w' \times h' \times w'}$ to record the consensus on patch-wise similarities across the reference and target images of all attention maps. Namely, $C_t(i, j, \hat{i}, \hat{j})$ denotes the matching degree between patch (i, j) in the target and patch (\hat{i}, \hat{j}) in the reference. To aggregate information through the denoising process and stabilize the matching maps, C_t is estimated by considering both C_{t+1} and A_t . We further apply the geometric constraints to C_t to construct correspondence $P_t \in [0, 1]^{h' \times w' \times 2}$, where $P_t(i, j)$ is the corresponding normalized coordinate in the reference of patch (i, j) in the target. The correspondence P_t serves as the input and can facilitate denoising and inpainting in the next step $t - 1$.

3.2. Attention-Consensus Correspondence

With correspondence guidance, inpainting models can identify the most relevant parts to fill damaged regions, while avoiding interference from irrelevant parts. However, existing correspondence estimation approaches cannot find correspondences inside the damaged region. Inspired by semantic correspondence estimation using pre-trained diffusion models [12, 27, 28], we explore the capability of generalizing an inpainting diffusion model to joint corre-

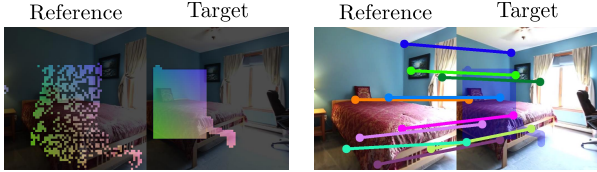


Figure 3. **Correspondences in the early stage.** The image on the left highlights the masked regions of the target and their most attended positions in the reference, indicated by colors, at the very first denoising step. The image on the right depicts a few correspondences computed at the first denoising step.

spondence estimation and image inpainting. Unlike methods [27, 28] using nearest neighbor search with diffusion features, we take self-attention scores as similarity matrices so that these scores can serve as the common domain for both correspondence estimation and image inpainting.

As shown in Figure 3, the self-attention scores present the correlation between references and targets even in the early generation stages. However, the attention map from a single attention layer is often less informative. To address this, we aggregate attention maps through accumulation across different layers. Specifically, we rescale averaged attention maps at different layers to a common size of $(h' \times 2w' \times h' \times 2w')$ and sum them up, resulting in aggregated attention map A_t . Since correspondence is established across the reference and target images, we consider only the parts of self-attention scores where queries are from the target and key-value pairs are from the reference. Therefore, the target-to-reference attention map $A_t^{\text{tar2ref}} \in \mathbb{R}^{h' \times w' \times h' \times w'}$, a submatrix of A_t , is extracted accordingly.

To calculate correspondence, we compute the matching map C_t by merging all aggregated attention maps until the current timestep, *i.e.*, $C_t = \sum_{\tau=t}^T A_\tau^{\text{tar2ref}}$. The reason we choose to calculate correspondences using consensus of the aggregated attention scores from multiple layers and timesteps is to eliminate the individual biases in certain layers and timesteps.

With the matching map C_t , the correspondence $P_t(i, j)$ for target token (i, j) is presented as the corresponding reference token and is determined via

$$P_t(i, j) = \underset{(\hat{i}, \hat{j})}{\operatorname{argmax}} C_t(i, j, \hat{i}, \hat{j}), \quad (1)$$

where (i, j) and (\hat{i}, \hat{j}) are the coordinates of the target and reference tokens, respectively.

3.3. Correspondence Refinement

As the self-attention mechanism is essential to propagating reference content to the damaged regions in the target, target query tokens attending to irrelevant reference tokens typically lead to incorrect inpainting results. Since the preliminary correspondences P_t are established by referring to

merely individual reference-target token pairs, they are not stable. Guiding the inpainting process solely on these correspondences fails to prevent the target tokens from attending to irrelevant tokens. To this end, we propose a correspondence refining strategy, including *filtering* and *smoothing*, to eliminate the inaccurate correspondence in P_t .

Correspondence Filtering. Given that the effective correspondences only reside in the overlapping areas of the reference and target images, it is clear that not every target token has a corresponding reference token. We observe that the target tokens not located in the overlapping regions tend to exhibit strong attention to certain reference tokens. We define these strongly attended but irrelevant reference tokens as *dominant tokens*. They need to be removed from correspondence constraints to avoid wrong feature propagation. Dominant tokens are identified by the presence of strong attention from diverse target tokens in P_t . In practice, we consider reference tokens with more than a certain number of corresponding target tokens as dominant, and their associated correspondences are probably outliers and, therefore, are excluded from P_t . We empirically set the threshold to four tokens and observe that this parameter is insensitive to the experimental results. Additionally, we notice that some target tokens within the overlapping regions are also affected by the dominant tokens, resulting in incorrect inpainting results. Hence, we save these excluded outlier correspondences as P_t^o , which are used to mitigate the adverse effects they caused through guidance.

Correspondence Smoothing. We introduce a smoothing mechanism based on the observation that when an incorrect inpainting result is present, a portion of target tokens at the center of the masked area (*i.e.*, the damaged region) exhibit incorrect correspondences. Conversely, their surrounding tokens, located around the edges of the mask, give more accurate correspondences and demonstrate attention scores consistent across different attention layers and timesteps. Therefore, we employ neighborhood weighted averages for smoothing correspondence, which corrects misleading correspondence, aiming to alleviate the presence of unwanted objects and incorrect geometry.

To calculate neighborhood weighted averages on the correspondence, we create a displacement matrix $D_t \in \mathbb{R}^{h' \times w'}$ indicating the differences between each target token and its corresponding reference tokens in coordinate, *i.e.*, $D_t(i, j) = P_t(i, j) - (\hat{i}, \hat{j})$. Next, we construct the consensus matrix $W_t \in \mathbb{R}^{h' \times w'}$ by assigning the matching score $C_t(i, j, P_t(i, j))$ to $W_t(i, j)$ for target token (i, j) , whose corresponding reference token is $P_t(i, j)$. For outlier correspondences P_t^o , we set their consensus value to zero, and therefore they are ignored during the smoothing process. The neighborhood weighted average of D_t is then

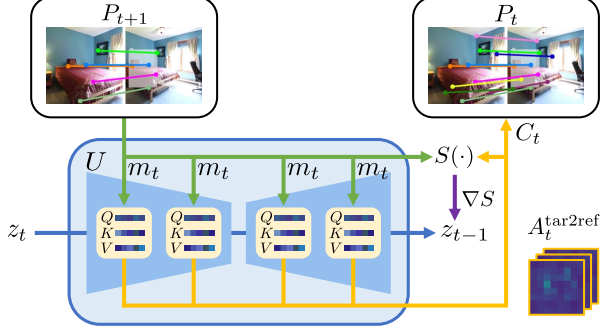


Figure 4. **Correspondence guidance in the diffusion U-Net.** At each denoising step t , the denoising process is guided by the correspondences estimated in the previous step, P_{t+1} , through attention masking with m_t and optimizing z_t using the objective function $S(\cdot)$. The generated attention maps A_t^{tar2ref} are then employed to further refine the estimated correspondences P_t by updating the matching map C_t .

calculated using W_t as weights as follows:

$$D_t^*(i, j) = \frac{1}{|W_t(i, j)|} \sum_{(\hat{i}, \hat{j}) \in \mathcal{N}(i, j)} D_t(\hat{i}, \hat{j}) \cdot W_t(\hat{i}, \hat{j}), \quad (2)$$

where $\mathcal{N}(i, j)$ is the set of neighborhood tokens of token (i, j) , and $|W_t(i, j)| = \sum_{(\hat{i}, \hat{j}) \in \mathcal{N}(i, j)} W_t(\hat{i}, \hat{j})$. In this formulation, we can propagate more accurate correspondences with higher degrees of consensus to those tokens of incorrect correspondences in the form of displacements, and the smoothed displacements are converted back to correspondences through $P_t^*(i, j) = D_t^*(i, j) + (i, j)$. The value of the smoothed correspondence P_t^* is then assigned back to the original correspondence: $P_t^* \rightarrow P_t$.

3.4. Cyclic Enhancement

By applying correspondence constraints to the denoising process, CorrFill establishes a cyclic enhancement that jointly improves the correspondence and inpainting processes at each iteration, progressively guiding the generation toward a faithful result. Figure 4 illustrates one cycle of the cyclic enhancement during a denoising step. Given the estimated correspondence P_{t+1} from the previous step, we guide the denoising process of the diffusion model by employing attention masks m_t across all self-attention layers and further enhancing the input latent z_t with an objective function S . The produced attention map A_t^{tar2ref} is then used to enhance the estimated correspondence P_{t+1} to P_t for the next step through updating the matching map C_t .

Attention Masking. To integrate correspondence constraints into the diffusion model, we employ attention masks within each self-attention layer. These attention masks are incorporated into the affinity matrix to modulate the influence of different value tokens.

The attention mechanism evaluates the contribution of value tokens through the affinity matrix, expressed as $QK^\top / \sqrt{d_a} \in \mathbb{R}^{(h' \times 2w') \times (h' \times 2w')}$, where Q and K are query and key token vectors, respectively, and d_a is the embedding dimension. For ease of discussion, we focus on operations conducted at a scale of $1/8$, while these operations are consistent across all attention layers, regardless of scale. The attention mask $m_t \in \mathbb{R}^{(h' \times 2w') \times (h' \times 2w')}$ adjusts the contribution of value tokens by adding either negative or positive values to the affinity matrix, resulting in the modified attentions: $(QK^\top + m_t) / \sqrt{d_a} \in \mathbb{R}^{(h' \times 2w') \times (h' \times 2w')}$.

We represent the attention mask in the shape of $h' \times 2w' \times h' \times 2w'$, which preserves the spatial context for both the queries and keys. We define a slice of the attention mask for a token (i, j) as $m_t^{ij} \in \mathbb{R}^{h' \times 2w'}$, denoting the part where the dot product between the query (i, j) and all keys occurs. The attention masks are composed according to the estimated correspondence P_{t+1} from the previous denoising step. For a target token (i, j) whose correspondence is not an outlier, the element in the slice m_t^{ij} is defined by

$$m_t^{ij}(\hat{i}, \hat{j}) = \begin{cases} v, & \text{if } (\hat{i}, \hat{j}) \in \mathcal{N}(P_{t+1}(i, j)), \\ -\infty, & \text{if } (\hat{i}, \hat{j}) \in \mathcal{R} - \mathcal{N}(P_{t+1}(i, j)), \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

where v represents a small positive number, $\mathcal{N}(P_{t+1}(i, j))$ denotes the neighboring tokens of $P_{t+1}(i, j)$, and \mathcal{R} refers to the set of all reference tokens. When the attention mask is applied to a self-attention layer, this slice of the mask boosts the attention values of the corresponding areas, thereby promoting attention for the relevant tokens. Conversely, it diminishes the attention values for other reference tokens, preventing them from being attended to.

For outlier tokens in P_{t+1}^o , the values assigned to their slices are defined as follows:

$$m_t^{ij}(\hat{i}, \hat{j}) = \begin{cases} -\infty, & \text{if } (\hat{i}, \hat{j}) \in \mathcal{R} \cap \mathcal{N}(P_{t+1}^o(i, j)), \\ 0, & \text{otherwise,} \end{cases} \quad (4)$$

This slice of the attention mask prevents the token (i, j) from attending to the irrelevant area, which is identified by the outlier correspondences. The remaining elements of the attention mask are assigned to 0, thereby preserving the original attention values for those tokens.

Latent Tensor Optimization. Similar to the observations made in recent studies on reliable generation within diffusion models [4, 13], we notice that solely employing attention masking is insufficient for steering inpainting towards the desired outcomes. To address this issue, we adopt a similar strategy, utilizing the produced constraints for further guidance by optimizing the latent tensor z_t with an objective function S . The core concept is to optimize z_t in

Method	RealEstate10K			MegaDepth			
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	
Paint-by-Example	Baseline	20.03	0.8528	0.1379	20.48	0.8274	0.1138
	+CorrFill	21.57(+1.54)	0.8724(+0.0196)	0.1097(-0.0282)	21.02(+0.54)	0.8343(+0.0069)	0.1014(-0.0124)
IP-Adapter-Plus	Baseline	21.26	0.8704	0.1127	21.33	0.8394	0.0989
	+CorrFill	25.10(+3.84)	0.8990(+0.0286)	0.0642(-0.0485)	22.14(+0.81)	0.8473(+0.0079)	0.0838(-0.0151)
Side-by-side	Baseline	23.32	0.8941	0.0856	22.89	0.8538	0.0850
	+CorrFill	25.81(+2.49)	0.9092(+0.0151)	0.0552(-0.0304)	23.24(+0.35)	0.8571(+0.0033)	0.0777(-0.0073)
LeftRefill	Baseline	26.71	0.9163	0.0443	23.60	0.8649	0.0653
	+CorrFill	26.97(+0.26)	0.9175(+0.0012)	0.0427(-0.0016)	23.60	0.8649	0.0651(-0.0002)

Table 1. **Quantitative results.** We demonstrate the evaluations of 4 different baselines: Paint-by-Example [25], IP-Adapter-Plus [26], Side-by-side [2] and LeftRefill [2], with and without the application of our CorrFill on the dataset RealEstate10K and MegaDepth.

a direction that aligns with the desired outcomes, specifically by ensuring that the attention of a token adheres to the pattern prescribed by P_{t+1} .

As depicted in Figure 4, we collect attention maps from all self-attention layers within U . Similar to the process producing A_t^{tar2ref} , the attention maps are reshaped, resized, and used to extract the target-to-reference submatrix, resulting in $(A_l)_t^{\text{tar2ref}}$, where l denotes the layer it is collected from. Instead of aggregating them, we calculate their gradients of the objective function S separately and update the input latent z_t by gradient descent. The objective function S is defined as follows:

$$S((A_l)_t^{\text{tar2ref}}) = \text{BCE}(\text{Sigmoid}(\text{Norm}((A_l)_t^{\text{tar2ref}})), E(P_{t+1})), \quad (5)$$

where function $\text{Norm}(\cdot)$ normalize matrix $(A_l)_t^{\text{tar2ref}}$, and $\text{BCE}(\cdot)$ is the weighted binary cross-entropy to $[0, 1]$ $E(\cdot)$ turns P_{t+1} into a one-hot tensor of the same shape as $(A_l)_t^{\text{tar2ref}}$. In this formulation, the input latent z_t is optimized toward a direction where its attention maps are encouraged to adhere to the correspondence constraint.

Implementation Details. Our proposed CorrFill is developed based on Stable Diffusion v2 Inpainting¹ [19], designed to work as a plug-in compatible with various diffusion models. The image resolution for both target and reference images is set to 512×512 , while the size of the latent representation for the stitched image $\epsilon(I_{\text{ref;tar}})$ is 64×128 . DDIM [21] sampling is used for efficient generation, with the number of sample steps set to 50. Further details, including the choices of parameters, are provided in the supplementary.

4. Experiments

4.1. Experimental Settings

Datasets. Since there is no publicly available benchmark for this task, we follow previous works [2, 30] to prepare the dataset for evaluation. We conducted experiments on

¹<https://huggingface.co/stabilityai/stable-diffusion-2-inpainting>

subsets sampled from two datasets: RealEstate10K [31] and MegaDepth [11]. We sampled 500 pairs of references and targets from each dataset and generated the inpainting masks based on the content of image pairs.

Before the sampling process, we cropped the images into squares and resized them to a resolution of 512×512 . To identify suitable images for experimentation with reference-based inpainting, we utilized the mid-level vision similarity metric DreamSim [8]. Images from the same scene are considered to form an image pair if the DreamSim distance is below 0.2. This threshold indicates that the image pair exhibits a certain degree of similarity in both appearance and semantics, making it suitable for evaluating reference-based inpainting. To prevent images in a pair from being overly similar, the pairs with distance below 0.1 are discarded for RealEstate10K, for its images are likely to be overly similar compared to those in MegaDepth.

To consistently reflect the capability of reference-based inpainting, we generate inpainting masks based on content in the references and targets. Specifically, we employ feature matching [23] to identify corresponding keypoints on the intact image pairs, subsequently creating masks that cover portions of these keypoints on the target images. This approach ensures that the recovery process for damaged regions relies on the content of the reference images. To simulate real-world applications, the masks are randomly generated, combining the shape of a rectangle with several strokes. The generated dataset will be publicly available.

Evaluation Metrics. We follow the common practice [2, 30, 32] and employ three evaluation metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM) and Learned Perceptual Image Patch Similarity (LPIPS). These metrics evaluate the similarities between restored targets and their ground truths.

4.2. Experimental Results

Our method serves as a plug-in module. We perform a thorough comparison of our method with four existing baselines in Table 1 and Figure 5.

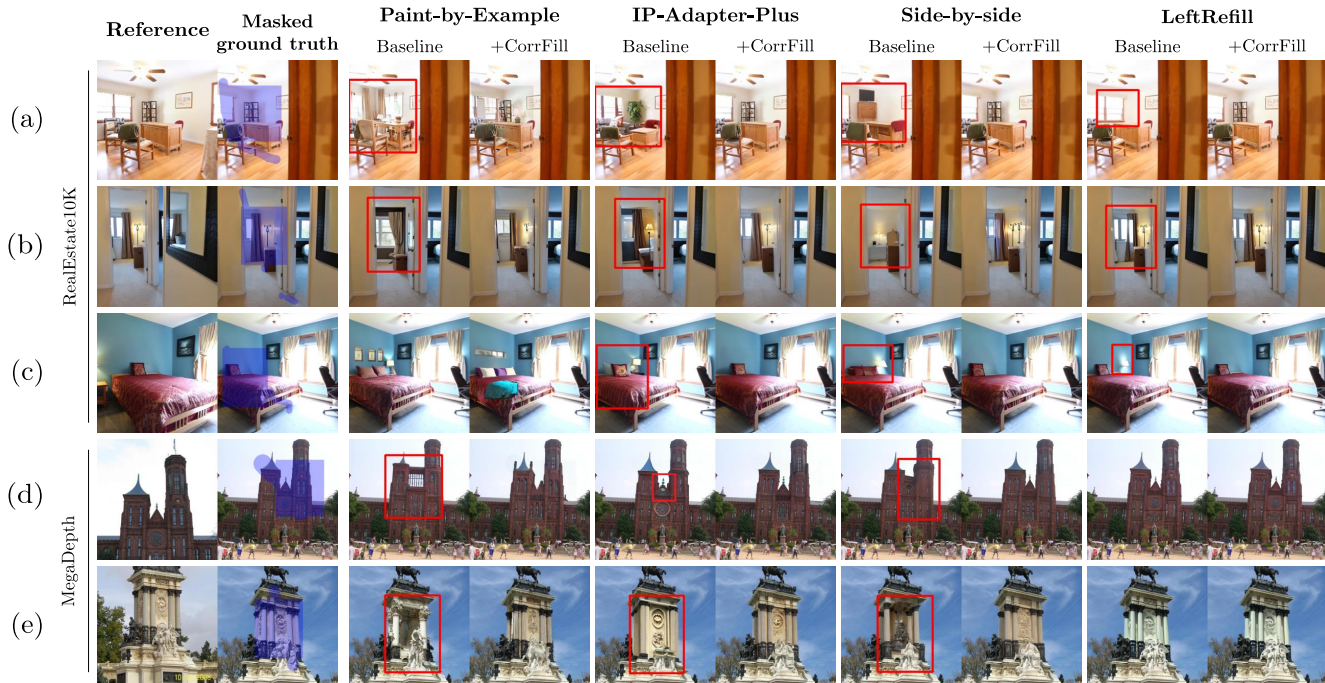


Figure 5. **Qualitative results.** We present the qualitative results with four different baselines and their counterparts integrated with our method on two datasets. We highlight the problematic regions in the results of the baseline methods that our approach can effectively address by enclosing them in red boxes. The inpainting masks are generated based on the content in image pairs.

Baselines. The comparison baselines achieve reference-based inpainting through various techniques. LeftRefill [2] stitches the reference and target images side by side, filling damaged regions with Stable Diffusion Inpainting [19], and incorporating prompt-tuning technique to learn an optimized prompt embedding specifically for the reference-based inpainting task. Side-by-side inpainting [2] is a variant of LeftRefill, which does not employ the prompt-tuning technique. In practice, we implement Side-by-side by providing an empty text prompt, which means that the model is not given any explicit instructions regarding our task.

Although Paint-by-Example [25] is designed for reference-based inpainting, its goal does not completely align with those of other baselines. Rather than restoring damaged regions by adhering to the fine-grained details of the reference image, this method focuses on generating plausible results based solely on the semantic attributes of the reference images. While it is conditioned on the reference images using their CLIP embeddings, it introduces an information bottleneck by considering only the global attributes, leading to a loss of fine-grained details. IP-Adapter-Plus [26] presents another approach that conditions the diffusion model on the CLIP embeddings of the reference images. In contrast to Paint-by-Example, IP-Adapter-Plus retains fine-grained details by incorporating all spatial tokens of CLIP embeddings.

To compare against these baseline methods, we integrate

CorrFill into them. For LeftRefill and Side-by-side, since they already employ the same formulation of stitched image inputs, we directly apply CorrFill to their diffusion models. For Paint-by-Example and IP-Adapter-Plus, we modify their inputs to match the stitched reference formulation and then apply our CorrFill.

Quantitative Results. We evaluate four baseline methods and their counterparts involving CorrFill module. As shown in Table 1, our method consistently shows improvement across all baselines on RealEstate10K.

With CorrFill, the performance of Side-by-side is elevated to 25.81dB in PSNR, increased by 2.49dB. Since the model of Side-by-side is not specifically designed for reference-based inpainting, these improvements highlight the effectiveness of CorrFill in enhancing the model’s awareness of reference-target correlations. While IP-Adapter-Plus preserves the appearance details in the reference images, it struggles to capture correct spatial correlations. CorrFill significantly improves its performance by 3.84dB in PSNR through incorporating correspondence constraints. Although the enhanced results of Paint-by-Example are inferior to those of other approaches due to the re-training of the model on a slightly different task, CorrFill still enhances their performance by incorporating additional reference content. CorrFill further improve the performance of LeftRefill, which is the state-of-the-art approach to the best of our knowledge.

Baselines	Module Components	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Side-by-side	Baseline	23.32	0.8941	0.0856
	+ Attention Masking	23.70	0.8995	0.0736
	+ Outlier Filtering	24.36	0.9031	0.0692
	+ Correspondence Smoothing	24.37	0.9030	0.0694
	+ Latent z_t Optimization	25.81	0.9092	0.0552
LeftRefill	Baseline	26.71	0.9163	0.0443
	+ Attention Masking	26.45	0.9153	0.0458
	+ Outlier Filtering	26.78	0.9165	0.0438
	+ Correspondence Smoothing	26.79	0.9166	0.0438
	+ Latent z_t Optimization	26.97	0.9175	0.0427

Table 2. **Ablation study on key components of CorrFill.** The ablation study on key components of CorrFill is conducted with Side-by-side and LeftRefill [2] baselines, on RealEstate10K dataset. It is important to note that outlier filtering and correspondence smoothing can only be implemented when attention masking is enabled.

The improvements on the challenging MegaDepth dataset are more limited due to the significant changes in viewpoints, which can lead to failures in correspondence estimation. Additionally, the nature of the dynamic scenes may diminish the benefits gained from strictly adhering to the reference images. Despite these challenges, CorrFill still demonstrates clear improvements for Paint-by-Example, IP-Adapter-Plus, and Side-by-side across all evaluation metrics.

Qualitative Results. Figure 5 presents a quality comparison of four baseline methods and their variants integrated with CorrFill. In the comparison involving IP-Adapter-Plus, Side-by-side, and LeftRefill on the RealEstate10K dataset, our method effectively addresses the issues present in the baseline results, including the removal of unwanted objects in (c) and the correction of incorrect scene layouts in (a),(b), and (c). While our method does not completely resolve all issues in Paint-by-Example, it still achieves a higher level of faithfulness in the results. In the comparisons on MegaDepth, results exhibiting greater faithfulness can be observed when compared to most of the baselines.

4.3. Ablation Study

We perform an ablation study by incrementally activating different key components of CorrFill to assess the impact of each component in Table 2. As attention masking serves as a prerequisite for outlier filtering and correspondence smoothing, our analysis reveals that the correspondence refinement strategies effectively enhance performance when applied in conjunction with attention masking. Although the improvement brought by correspondence smoothing may seem subtle, it can be the crucial component for correcting incorrect content in certain cases, and an example is provided in the supplementary. The optimization of the latent tensor z_t further boosts performance, while also benefiting from the correspondence refinements.

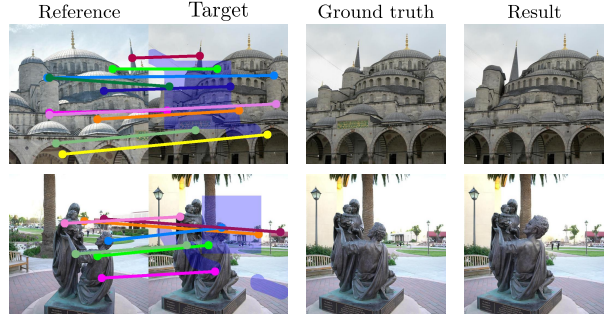


Figure 6. **Failure cases.** The image on the left depicts the estimated correspondences. The image on the right shows the inpainting result of CorrFill integrated on LeftRefill [2]. In the first row, the repetitive structures and complex geometry introduce ambiguity in correspondence estimation. In the second row, the incorrect orientation of the statue’s head demonstrates that correspondence constraints in 2D space are inadequate when faced with significant changes in viewpoint due to a lack of 3D awareness.

5. Conclusion

We propose CorrFill, a training-free module that incorporates correspondence constraints into reference-based image inpainting diffusion models. CorrFill achieve higher degrees of faithfulness to the reference images in the inpainting results by guiding the inpainting process with correspondence between the reference and target images. To perform this guidance, we exploit the capability of diffusion models to estimate correspondence during the inpainting process, and we utilize this correspondence to constrain the inpainting through self-attention masking and input latent optimization. Experimental results demonstrate the effectiveness of CorrFill in enhancing reference-target correlations in the inpainting results for multiple baseline methods. CorrFill improves performance across these baselines on RealEstate10K and MegaDepth datasets, pushing the limits for state-of-the-art reference-based inpainting methods.

Limitations. As shown in the first row of Figure 6, CorrFill fails to faithfully restore damaged regions when encountered with scenes featuring complex geometry or repetitive structures, which introduce significant ambiguity when estimating correspondences. In the second row, although CorrFill successfully captures the reference-target correlation, the incorrect orientation of the statue’s head suggests that our correspondence constraints in 2D space are susceptible to the significant geometric variations of objects, which require advanced 3D prior.

Acknowledgement. This work was supported in part by the National Science and Technology Council (NSTC) under grants 112-2221-E-A49-090-MY3, 111-2628-E-A49-025-MY3, and 112-2634-F-002-005. This work was funded in part by NVIDIA.

References

- [1] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. **2**
- [2] Chenjie Cao, Yunuo Cai, Qiaole Dong, Yikai Wang, and Yanwei Fu. Leftrefill: Filling right canvas based on left reference through generalized text-to-image diffusion model. In *CVPR*, 2024. **1, 2, 3, 6, 7, 8**
- [3] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *TOG*, 2023. **2**
- [4] Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. In *WACV*, 2024. **2, 5**
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. **1**
- [6] Dave Epstein, Allan Jabri, Ben Poole, Alexei Efros, and Aleksander Holynski. Diffusion self-guidance for controllable image generation. In *NeurIPS*, 2023. **2**
- [7] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. In *ICLR*, 2023. **2**
- [8] Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data. In *NeurIPS*, 2023. **6**
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. **1, 2**
- [10] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2022. **3**
- [11] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *CVPR*, 2018. **6**
- [12] Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion hyperfeatures: Searching through time and space for semantic correspondence. In *NeurIPS*, 2023. **3**
- [13] Hayk Manukyan, Andranik Sargsyan, Barsegh Atanyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Hd-painter: High-resolution and prompt-faithful text-guided image inpainting with diffusion models. *arXiv preprint arXiv:2312.14091*, 2023. **2, 5**
- [14] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *AAAI*, 2024. **2**
- [15] Seoung Wug Oh, Sungho Lee, Joon-Young Lee, and Seon Joo Kim. Onion-peel networks for deep video completion. In *ICCV*, 2019. **1**
- [16] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *TMLR*, 2024. **3**
- [17] 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. Learning transferable visual models from natural language supervision. In *ICML*, 2021. **2**
- [18] Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map alignment. In *NeurIPS*, 2023. **2**
- [19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. **1, 2, 3, 6, 7**
- [20] Jaskirat Singh, Stephen Gould, and Liang Zheng. High-fidelity guided image synthesis with latent diffusion models. In *CVPR*, 2023. **2**
- [21] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. **6**
- [22] Luming Tang, Nataniel Ruiz, Qinghao Chu, Yuanzhen Li, Aleksander Holynski, David E. Jacobs, Bharath Hariharan, Yael Pritch, Neal Wadhwa, Kir Aberman, and Michael Rubinstein. Realfill: Reference-driven generation for authentic image completion. *TOG*, 2024. **1, 2**
- [23] Shitao Tang, Jiahui Zhang, Siyu Zhu, and Ping Tan. Quadtree attention for vision transformers. In *ICLR*, 2022. **6**
- [24] DeJia Xu, Xingqian Xu, Wenyan Cong, Humphrey Shi, and Zhangyang Wang. Reference-based painterly inpainting via diffusion: Crossing the wild reference domain gap. *arXiv preprint arxiv:2307.10584*, 2023. **1**
- [25] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *CVPR*, 2023. **1, 2, 6, 7**
- [26] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arxiv:2308.06721*, 2023. **1, 2, 6, 7**
- [27] Junyi Zhang, Charles Herrmann, Junhwa Hur, Eric Chen, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. Telling left from right: Identifying geometry-aware semantic correspondence. In *CVPR*, 2024. **3, 4**
- [28] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. In *NeurIPS*, 2023. **3, 4**
- [29] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. **2**

- [30] Yunhan Zhao, Connelly Barnes, Yuqian Zhou, Eli Shechtman, Sohrab Amirghodsi, and Charless Fowlkes. Geofill: Reference-based image inpainting with better geometric understanding. In *WACV*, 2023. [1](#), [2](#), [6](#)
- [31] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. [6](#)
- [32] Yuqian Zhou, Connelly Barnes, Eli Shechtman, and Sohrab Amirghodsi. Transfill: Reference-guided image inpainting by merging multiple color and spatial transformations. In *CVPR*, 2021. [1](#), [2](#), [6](#)

CorrFill: Enhancing Faithfulness in Reference-based Inpainting with Correspondence Guidance in Diffusion Models – Supplementary Materials

Kuan-Hung Liu¹ Cheng-Kun Yang^{2*} Min-Hung Chen³ Yu-Lun Liu¹ Yen-Yu Lin¹

¹National Yang Ming Chiao Tung University ²National Taiwan University ³NVIDIA

<https://corrfill.github.io/>

In this supplementary material, we provide implementation details, showcase examples to support our proposed approaches, and present advanced analyses of the proposed method.

A. Implementation Details

A.1. Implementation Details of Baselines

For comparisons, we implement all baseline methods using the Python library Diffusers [3]. For Paint-by-Example [4], we utilize their publicly released model weights. Side-by-side [1] is essentially an inpainting base model, and we directly utilize Stable Diffusion v2 Inpainting model¹ as its implementation. In the case of LeftRefill [1], we use the same model of Side-by-side and incorporate LeftRefill’s learned prompt embedding. We integrate IP-Adapter-Plus [5] module into a Stable Diffusion Inpainting model using their released pre-trained weights.

A.2. Implementation Details of CorrFill

CorrFill modify baseline models by substituting the attention processing function across all self-attention layers. Correspondence estimation and attention masking are then carried out in the substituted function. We also collect the attention maps used to optimize input latent tensor z_t in the attention processing function, and the gradients are computed in the denoising main loop of the diffusion models. Since optimizing z_t requires additional memory, a gradient accumulation strategy can be employed to trade off inference time for lower memory requirements. We conduct the experiments using an NVIDIA RTX A5000 GPU with 24GB of memory.

A.3. Details of Dataset Sampling

RealEstate10K is a video dataset comprising approximately 80,000 clips sourced from YouTube. Given that the clips are recorded by cameras with stable trajectories, adjacent frames tend to exhibit high similarity. Therefore, when

selecting image pairs from RealEstate10K, we specifically consider frames that are separated by 30 frames during the sampling process.

A.4. Choices of Parameters

The parameters used in the comparisons presented in the main papers are reported in Table 1. Step_a and Step_o represent the number of steps guided by attention masking and latent tensor optimization, respectively, out of a total of 50 sampling steps. Win_a is the radius that determines the neighborhood of a token used in the creation of attention masks, and Win_s is the radius that determines the neighborhood for the weighted average used in attention smoothing. Str_a and Str_o indicate the value v added to the attention mask and the weight for controlling the guidance strength of latent tensor optimization, respectively.

We selected the parameters by evaluating the subsets of our datasets. During this evaluation, we tested various parameter settings and observed their responses in the results of different baseline methods and datasets. The general strategy is to increase the influence of guidance for the combinations that can significantly benefit from enhanced faithfulness.

B. Effectiveness of Proposed Components

B.1. Attention Smoothing

In the quantitative ablation study presented in the main paper, the performance gains from attention smoothing are not particularly significant. However, we provide one example demonstrating how attention smoothing serves as a crucial component in achieving accurate inpainting results in Figure 1.

B.2. Correspondence Update Policies

In this section, we demonstrate the effectiveness of two policies including cyclic enhancement and accumulation of attention maps over timesteps. We conduct a comparison of the correctness of the estimated correspondences against two counterparts excluding the two policies on

*Now at MediaTek Inc., Taiwan.

¹<https://huggingface.co/stabilityai/stable-diffusion-2-inpainting>

Parameter	Paint-by-Example		IP-Adapter-Plus		Side-by-side		LeftRefill	
	RealEstate10K	MegaDepth	RealEstate10K	MegaDepth	RealEstate10K	MegaDepth	RealEstate10K	MegaDepth
Step _a	50	25	25	25	50	25	5	5
Step _o	50	50	50	50	50	50	50	5
Win _a	4(t)	4(t)	0.3(i)	5(t)	2(t)	3(t)	0.3(i)	2(t)
Win _s	0.4	0.4	0.2	0.2	0.2	0.2	0.05	0.2
Str _a	1	1	1	1	1	1	1	0
Str _o	2	0.5	2	0.5	2	0.5	0.5	0.5

Table 1. **List of parameters.** The comparisons presented in the main papers are conducted using these parameters. (t) indicates that the value refers to the number of tokens, and (i) denotes that the value is the ratio to the size of encoded images, *i.e.*, h' . For Win_s, all the values are the ratios to the size of encoded images.

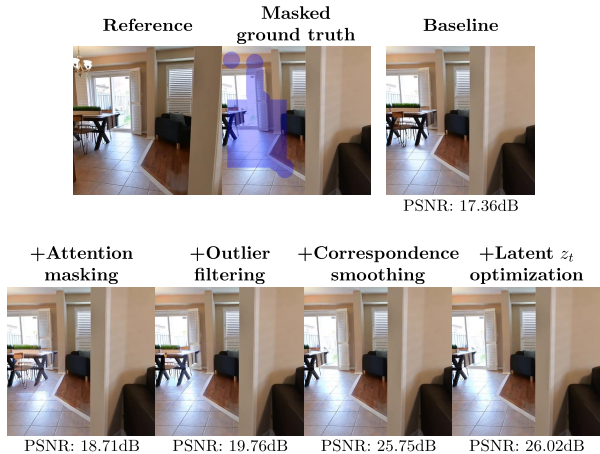


Figure 1. **Importance of Smoothing.** An example where correspondence smoothing is the pivotal component for correcting the incorrect geometry in the result of the baseline [1].

RealEstate10K. To estimate the correctness of the correspondences, we generate pseudo-ground truth correspondences using an image matching method [2]. We define a correspondence with an error within the size of one token as a correct correspondence. The counterpart that does not accumulate attention scores over time utilizes the most recently produced correspondences for guidance. The counterpart without cyclic enhancement is guided by the correspondences computed in the first step. The average numbers of correct correspondences during different stages of the inpainting process are illustrated in Figure 2. The counterpart “No acc” fails to achieve stability, while “No cyc.” relies on the correspondence produced in the first step for guidance, resulting in inferior results. The PSNR performance results for “Ours”, “No acc.”, and “No cyc.” are 27.39dB, 27.34dB, and 27.25dB, respectively.

C. Further Analysis

C.1. Time Efficiency

We analyze the average execution time for the inpainting of a single input with different key components enabled, following the experimental settings described earlier. The

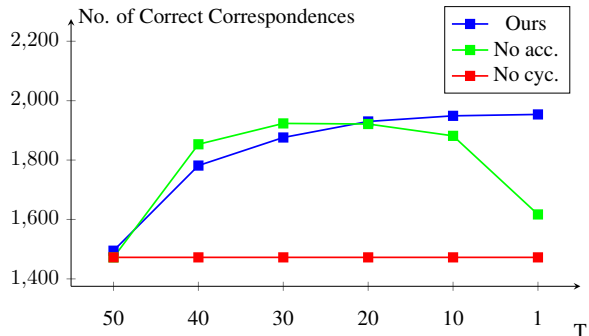


Figure 2. **Numbers of correct correspondences.** The graph illustrates the numbers of correct correspondences for three versions of CorrFill. “T” denotes the timesteps of the reverse process, where inpainting progresses from $T = 50$ to 0. “Ours” represents our proposed method, which utilizes cyclic enhancement and estimates correspondences using aggregated attention scores across different timesteps. “No acc.” and “No cyc.” are the counterparts that exclude the accumulation of attention maps and cyclic enhancement, respectively.

Method	Execution Time(s)	Change(s)
Baseline	6.69	-
+ Attention Masking	13.77	+7.08
+ Outlier Filtering	14.97	+1.20
+ Correspondence Smoothing	15.76	+0.79
+ Latent z_t Optimization	66.52	+50.76

Table 2. **Time analysis of key components of CorrFill.** The execution times for the inpainting of an input were measured while incrementally enabling the key components. The baseline used in the analysis is LeftRefill.

average execution times with LeftRefill as the baseline are reported in Table 2, which indicates that the latent input optimization contributes the most additional execution time within the proposed method. The increase in execution time is primarily attributed to the necessity of gradient calculation during each denoising iteration.

C.2. Extreme Case

Since CorrFill is an improvement method designed to enhance faithfulness, it encounters certain extreme cases



Figure 3. **Results with large masks.** The inpainting and outpainting results for the baseline method and CorrFill are presented. The first two rows depict the inpainting results, while the last row illustrates the outpainting results. All masks cover 50% of the target images. CorrFill cannot consistently enhance the results due to the significant degradation in the inpainting performance of the baseline method.

that challenge its performance, particularly when baseline models struggle to address them. While we previously discussed the issue of significant geometric variation in the main paper, another notable challenge for the baseline models involves large masks. The ratios of masked pixels for our generated pairs of inputs typically range from 10% to 40%. We find that when faced with larger masks, the inpainting results produced by LeftRefill tend to degrade to a point where CorrFill is unable to enhance faithfulness effectively. Figure 3 illustrates this limitation of CorrFill that it relies on the robustness of the baseline model. While CorrFill successfully improves the results for the first row, it does not yield similar improvements for the other cases.

References

- [1] Chenjie Cao, Yunuo Cai, Qiaole Dong, Yikai Wang, and Yanwei Fu. Leftrefill: Filling right canvas based on left reference through generalized text-to-image diffusion model. In *CVPR*, 2024. 1, 2
- [2] Xuelun Shen, zhipeng cai, Wei Yin, Matthias Müller, Zijun Li, Kaixuan Wang, Xiaozhi Chen, and Cheng Wang. GIM: Learning generalizable image matcher from internet videos. In *ICLR*, 2024. 2
- [3] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models, 2022. 1
- [4] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by ex-

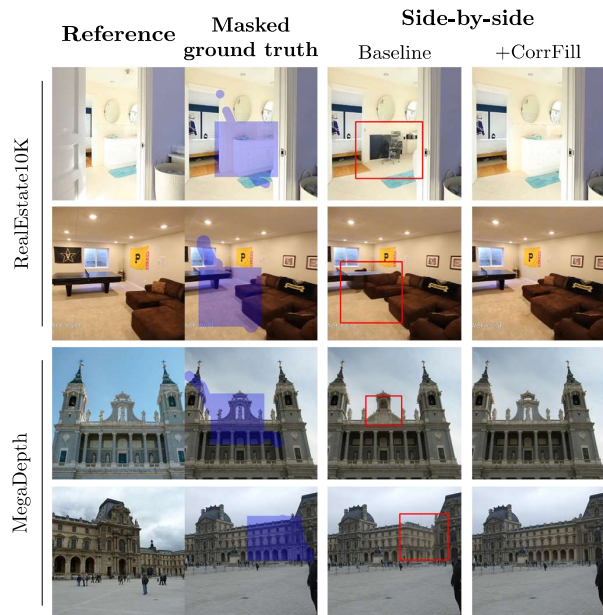


Figure 4. **Additional results with Side-by-side.** Problematic regions addressed by CorrFill are highlighted within the red boxes.

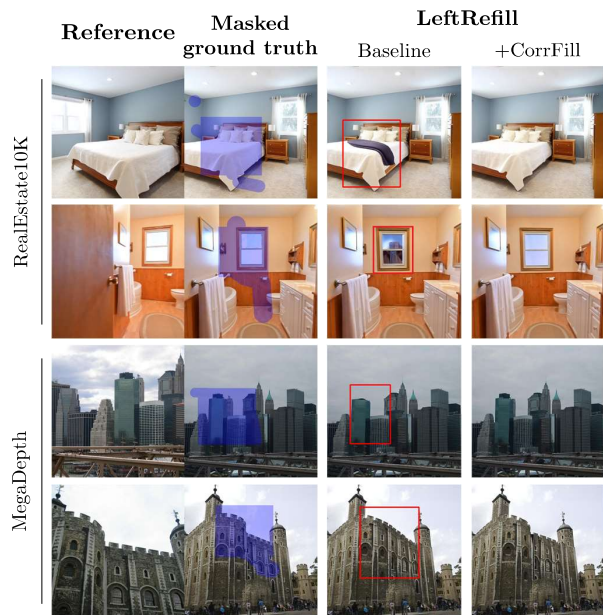


Figure 5. **Additional results with LeftRefill.** Problematic regions addressed by CorrFill are highlighted within the red boxes.

ample: Exemplar-based image editing with diffusion models. In *CVPR*, 2023. 1

- [5] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arxiv:2308.06721*, 2023. 1