

Completing Visual Objects via Bridging Generation and Segmentation

Xiang Li¹ Yinpeng Chen² Chung-Ching Lin² Hao Chen¹ Kai Hu¹ Rita Singh¹ Bhiksha Raj^{1,3}
Lijuan Wang² Zicheng Liu²

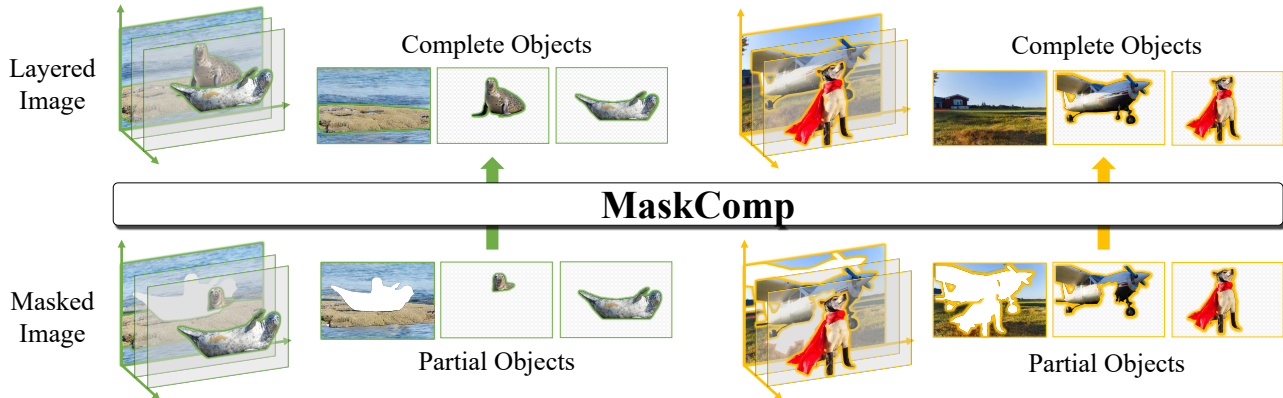


Figure 1. Given an image with mutual-occluded objects, MaskComp effectively completes occluded objects achieving a layered image.

Abstract

This paper presents a novel approach to object completion, with the primary goal of reconstructing a complete object from its partially visible components. Our method, named **MaskComp**, delineates the completion process through iterative stages of generation and segmentation. In each iteration, the object mask is provided as an additional condition to boost image generation, and, in return, the generated images can lead to a more accurate mask by fusing the segmentation of images. We demonstrate that the combination of one generation and one segmentation stage effectively functions as a mask denoiser. Through alternation between the generation and segmentation stages, the partial object mask is progressively refined, providing precise shape guidance and yielding superior object completion results. Our experiments demonstrate the superiority of MaskComp over existing approaches, e.g., ControlNet and Stable Diffusion, establishing it as an effective solution for object completion.

¹Carnegie Mellon University ²Microsoft ³MBZUAI. Correspondence to: Xiang Li <xl6@andrew.cmu.edu>.

1. Introduction

In recent years, creative image editing has attracted substantial attention and seen significant advancements. Recent breakthroughs in image generation techniques have delivered impressive results across various image editing tasks, including image inpainting (Xie et al., 2023), composition (Yang et al., 2023a) and colorization (Chang et al., 2023). However, another intriguing challenge lies in the domain of object completion (Fig. 1). This task involves the restoration of partially occluded objects within an image, representing the image as a layered stack of objects and background, which can potentially enable a number of more complicated editing tasks such as object layer switching. Unlike other conditional generation tasks, e.g., image inpainting, which only generates and integrates complete objects into images, object completion requires seamless alignment between the generated content and the given partial object, which imposes more challenges to recover realistic and comprehensive object shapes.

To guide the generative model in producing images according to a specific shape, additional conditions can be incorporated (Koley et al., 2023; Yang et al., 2023b). Image segmentation has been shown to be a critical technique for enhancing the realism and stability of generative models by providing pixel-level guidance during the synthesis process. Recent research, as exemplified in the latest study by Zhang et al. (Zhang et al., 2023), showcases that, by supplying object segmentations as additional high-quality masks for

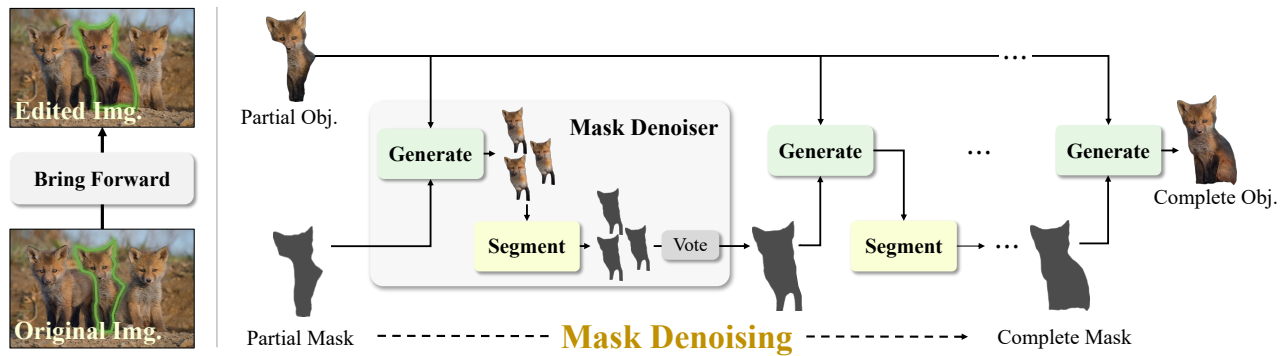


Figure 2. **Illustration of iterative mask denoising (IMD).** Starting from an initial partial object and its corresponding mask, IMD utilizes alternating generation and segmentation stages to progressively refine the partial mask until it converges to the complete mask. With the complete mask as the condition, the final complete object can be seamlessly generated.

shaping the objects, it becomes possible to generate complex images of remarkable fidelity.

In this paper, we present MaskComp, a novel approach that bridges image generation and segmentation for effective object completion. MaskComp is rooted in a fundamental observation: the quality of the resulting image in the mask-conditioned generation is directly influenced by the quality of the conditioned mask (Zhang et al., 2023). That says the more detailed the conditioned mask, the more realistic the generated image. Based on this observation, unlike prior object completion methods that solely rely on partially visible objects for generating complete objects, MaskComp introduces an additional mask condition combined with an iterative mask denoising (IMD) process, progressively refining the incomplete mask to provide comprehensive shape guidance to object completion.

Our approach formulates the partial mask as a noisy form of the complete mask and the IMD process is designed to iteratively denoise this noisy partial mask, eventually leading to the attainment of the complete mask. As illustrated in Fig. 2, each IMD step comprises two crucial stages: generation and segmentation. The generation stage’s objective is to produce complete object images conditioning on the visible portion of the target object and an object mask. Meanwhile, the segmentation stage is geared towards segmenting the object mask within the generated images and aggregating these segmented masks to obtain a superior mask that serves as the condition for the subsequent IMD step. By seamlessly integrating the generation and segmentation stages, we demonstrate that each IMD step effectively operates as a mask-denoising mechanism, taking a partially observed mask as input and yielding a progressively more complete mask as output. Consequently, through this iterative mask denoising process, the originally incomplete mask evolves into a satisfactory complete object mask, enabling the generation of complete objects guided by this refined mask.

The effectiveness of MaskComp is demonstrated by its capacity to address scenarios involving heavily occluded objects and its ability to generate realistic object representations through the utilization of mask guidance. In contrast to recent progress in the field of image generation research, our contributions can be succinctly outlined as follows:

- We explore and unveil the benefits of incorporating object masks into the object completion task. A novel approach, MaskComp, is proposed to seamlessly bridge the generation and segmentation.
- We formulate the partial mask as a form of noisy complete mask and introduce an iterative mask denoising (IMD) process, consisting of alternating generation and segmentation stages, to refine the object mask and thus improve the object completion.
- We conduct extensive experiments for analysis and comparison, the results of which indicate the strength and robustness of MaskComp against previous methods, e.g., Stable Diffusion.

2. Related Works

Conditional image generation. Conditional image generation (Lee et al., 2022; Gafni et al., 2022; Li et al., 2023f; Ye et al., 2023; Yu et al., 2021; Yan et al., 2019; Guo et al., 2021; Wan et al., 2021) involves the process of creating images based on specific conditions. These conditions can take various forms, such as layout (Li et al., 2020; Sun & Wu, 2019; Zhao et al., 2019), sketch (Koley et al., 2023), or semantic masks (Gu et al., 2019). For instance, Cascaded Diffusion Models (Ho et al., 2022) utilize ImageNet class labels as conditions, employing a two-stage pipeline of multiple diffusion models to generate high-resolution images. Meanwhile, in the work by (Schwag et al., 2022), diffusion models are guided to produce novel images from

low-density regions within the data manifold. Another noteworthy approach is CLIP (Radford et al., 2021), which has gained widespread adoption in guiding image generation in GANs using text prompts (Galatolo et al., 2021; Gal et al., 2022; Zhou et al., 2021b). In the realm of diffusion models, Semantic Diffusion Guidance (Liu et al., 2023) explores a unified framework for diffusion-based image generation with language, image, or multi-modal conditions. Dhariwal et al. (Dhariwal & Nichol, 2021) employ an ablated diffusion model that utilizes the gradients of a classifier to guide the diffusion process, balancing diversity and fidelity. Furthermore, Ho et al. (Ho & Salimans, 2022) introduce classifier-free guidance in conditional diffusion models, incorporating score estimates from both a conditional diffusion model and a jointly trained unconditional diffusion model.

Image segmentation. In the realm of image segmentation, traditional approaches have traditionally leaned on domain-specific network architectures to tackle various segmentation tasks, including semantic, instance, and panoptic segmentation (Long et al., 2015; Chen et al., 2015; He et al., 2017; Neven et al., 2019; Newell et al., 2017; Wang et al., 2020b; Cheng et al., 2020; Wang et al., 2021; 2020a; Li et al., 2023b;e;d;c;a; 2022b). However, recent strides in transformer-based methodologies, have highlighted the effectiveness of treating these tasks as mask classification challenges (Cheng et al., 2021; Zhang et al., 2021; Cheng et al., 2022; Carion et al., 2020). MaskFormer (Cheng et al., 2021) and its enhanced variant (Cheng et al., 2022) have introduced transformer-based architectures, coupling each mask prediction with a learnable query. Unlike prior techniques that learn semantic labels at the pixel level, they directly link semantic labels with mask predictions through query-based prediction. Notably, the Segment Anything Model (SAM) (Kirillov et al., 2023) represents a cutting-edge segmentation model that accommodates diverse visual and textual cues for zero-shot object segmentation. Similarly, SEEM (Zou et al., 2023) is another universal segmentation model that extends its capabilities to include object referencing through audio and scribble inputs. By leveraging those foundation segmentation models, e.g., SAM and SEEM, a number of downstream tasks can be boosted (Ma & Wang, 2023; Cen et al., 2023; Yu et al., 2023).

3. MaskComp

3.1. Problem Definition and Key Insight

We address the object completion task, wherein the objective is to predict the image of a complete object $I_c \in \mathbb{R}^{3 \times H \times W}$, based on its visible (non-occluded) part $I_p \in \mathbb{R}^{3 \times H \times W}$.

We first discuss the high-level idea of the proposed **Iterative Mask Denoising (IMD)** and then illustrate the module de-

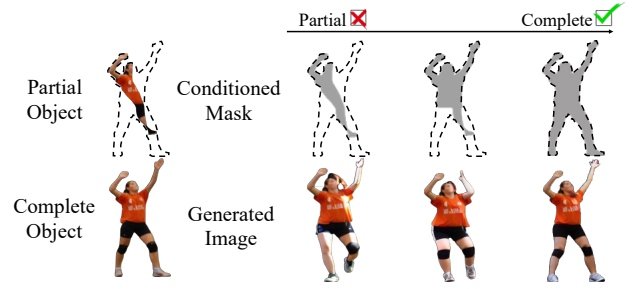


Figure 3. Object completion with different mask conditions.

tails in Section 3.3 and Section 3.4. The core of IMD is based on an essential observation: In the mask-conditioned generation, the quality of the generated object is intricately tied to the quality of the conditioned mask. As shown in Fig. 3, we visualize the completion result of the same partial object but with different conditioning masks. We notice a more complete object mask condition will result in a more complete and realistic object image. Based on this observation, high-quality occluded object completion can be achieved by providing a complete object mask as the condition.

3.2. Iterative Mask Denoising

However, in real-world scenarios, the complete object mask is not available. To address this problem, we propose the IMD process which leverages intertwined generation and segmentation processes to approach the partial mask to the complete mask gradually. Given a partially visible object I_p and its corresponding partial mask M_p , the conventional object completion task aims to find a generative model \mathcal{G} such that $I_c \leftarrow \mathcal{G}(I_p)$, where I_c is the complete object. Here, we additionally add the partial mask M_p to the condition $I_c \leftarrow \mathcal{G}(I_p, M_p)$, where M_p can be assumed as an addition of the complete mask and a noise $M_p = M_c + \Delta$. By introducing a segmentation model \mathcal{S} , we can find a mask denoiser $\mathcal{S} \circ \mathcal{G}$ from the object completion model:

$$M_c \leftarrow \mathcal{S} \circ \mathcal{G}(I_p, M_c + \Delta) \quad (1)$$

where $M_c = \mathcal{S}(I_c)$. Starting from the visible mask $M_0 = M_p$, as shown in Fig. 2, we repeatedly apply the mask denoiser $\mathcal{S} \circ \mathcal{G}$ to gradually approach the visible mask M_p to complete mask M_c . In each step, the input mask is denoised with a stack of generation and segmentation stages. Specifically, as the $\mathcal{S} \circ \mathcal{G}(\cdot)$ includes a generative process, we can obtain a set of estimations of denoised mask $\{M_t^{(k)}\}$. Here, we utilize a function $\mathcal{V}(\cdot)$ to find a more complete and reasonable mask from the N sampled masks and leverage it as the input mask for the next iteration to further denoise. The updating rule can be written as:

$$M_t^{(k)} = \mathcal{S} \circ \mathcal{G}(I_p, \hat{M}_{t-1}), \quad \hat{M}_t = \mathcal{V}(M_t^{(1)}, \dots, M_t^{(N)}) \quad (2)$$

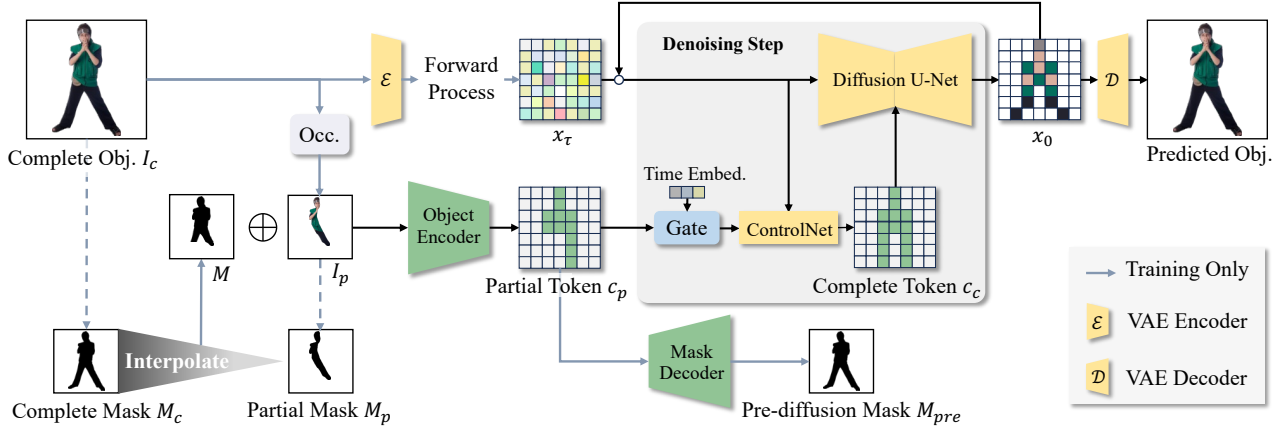


Figure 4. Illustration of CompNet (generation stage of MaskComp). The CompNet aims to recover the complete object I_c from the partial object I_p and a mask M . An object encoder is utilized to extract partial token c_p which is gated and fed to the ControlNet to form the complete token c_c . The complete token c_c serves as the condition to the diffusion U-Net to guide the conditional denoising process. In addition, a pre-diffusion mask is predicted from the partial token to encourage the object encoder to capture shape information.

Method	Objective	Objective with Segm.	Object Comp.
ControlNet	$I_p \leftarrow \mathcal{G}(I_p, M_p)$	$M_p \leftarrow \mathcal{S} \circ \mathcal{G}(I_p, M_p)$	✗
CompNet	$I_c \leftarrow \mathcal{G}(I_p, M_p)$	$M_c \leftarrow \mathcal{S} \circ \mathcal{G}(I_p, M_p)$	✓

Table 1. Objective difference with ControlNet.

where N is the number of sampled images in each iteration. With a satisfactory complete mask \hat{M}_T after T iterations, the object completion can be achieved accordingly by $\mathcal{G}(I_p, \hat{M}_T)$. The mathematical explanation of the process will be discussed in Section 3.5.

3.3. Generation Stage

We introduce **CompNet** as the generative model \mathcal{G} which aims to recover complete objects based on partial conditions. We build CompNet based on popular ControlNet (Zhang et al., 2023) while making fundamental modifications to enable object completion. As shown in Table 1, the target of ControlNet is to generate images strictly based on the given conditions, i.e., $I_p \leftarrow \mathcal{G}(I_p, M_p)$, making it unable to complete object. Differently, CompNet is designed to recover the object. With a segmentation network, it can act as a mask denoiser to refine the conditioned mask, i.e., $M_c \leftarrow \mathcal{S} \circ \mathcal{G}(I_p, M_p)$.

Mask condition. As illustrated on the left side of Fig. 4, we begin with a complete object I_c and its corresponding mask M_c . Our approach commences by occluding the complete object, retaining only the partially visible portion as I_p . Recall that the mask-denoising procedure initiates with the partial mask M_p and culminates with the complete mask M_c . To facilitate this iterative denoising, the model must effectively handle any mask that falls within the interpolation between the initial partial mask and the target complete mask. Consequently, we introduce a mask M with an occlu-

sion rate positioned between the partial and complete masks as a conditioning factor for the generative model. During training, we conduct the random occlusion process (detailed in Appendix C) twice for each complete mask M_c . The partial mask M_p is achieved by considering the occluded areas in both occlusion processes. The interpolated mask M is generated by using one of the occlusions.

Diffusion model. Diffusion models have achieved notable progress in synthesizing unprecedented image quality and have been successfully applied to many text-based image generation works (Rombach et al., 2022; Zhang et al., 2023). For our object completion task, the complete object can be generated by leveraging the diffusion process.

Specifically, the diffusion model generates image latent x by gradually reversing a Markov forward process. As shown in Figure 4, starting from $x_0 = \mathcal{E}(I_c)$, the forward process yields a sequence of increasing noisy tokens $\{x_\tau | \tau \in [1, T_G]\}$, where $x_\tau = \sqrt{\bar{\alpha}_\tau}y_0 + \sqrt{1 - \bar{\alpha}_\tau}\epsilon$, ϵ is the Gaussian noise, and α_τ decreases with the timestep τ . For the denoising process, the diffusion model progressively denoises a noisy token from the last step given the conditions $c = (I_p, M, E)$ by minimizing the following loss function: $\mathcal{L} = \mathbb{E}_{\tau, x_0, \epsilon} \|\epsilon_\theta(x_\tau, c, \tau) - \epsilon\|_2^2$. I_p , M , and E are the partial object, conditioned mask, and text prompt respectively.

CompNet architecture. Previous work (Zhang et al., 2023) has demonstrated an effective way to add additional control to generative diffusion models. We follow this architecture and make necessary modifications to adapt the architecture to object completion. As shown in Fig. 4, given the visible object I_p and the conditioning mask M , we first concatenate them and extract the partial token c_p with an object encoder.

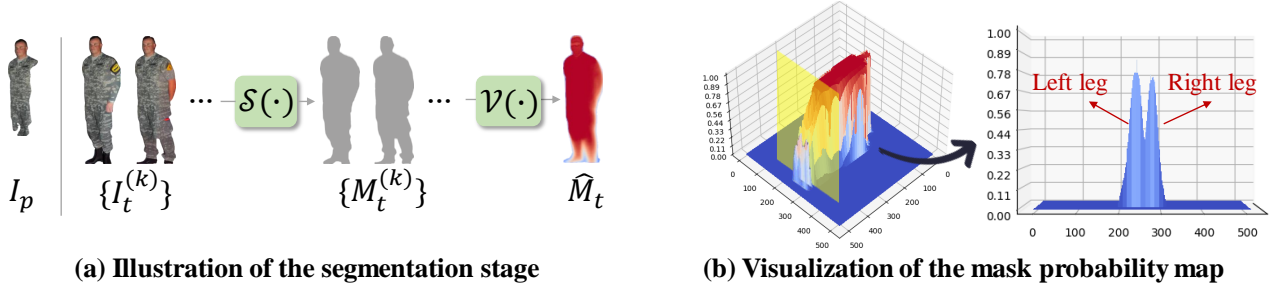


Figure 5. We calculate the mask probability map by averaging and normalizing the masks of sampled images. We show a cross-section of the lower leg to better visualize (shown as yellow).

Different from ControlNet (Zhang et al., 2023) assuming the condition is accurate, the object completion task relies on incomplete conditions. Specifically, in the early diffusion steps, the condition information is vital to complete the object. Nevertheless, in the later steps, inaccurate information in the condition can degrade the generated object. To tackle this problem, we introduce a time-variant gating operation to adjust the importance of conditions in the diffusion steps. We learn a linear transform $f : \mathbb{R}^C \rightarrow \mathbb{R}^1$ upon the time embedding $e_t \in \mathbb{R}^C$ and then apply it to the partial token as $f(e_t) \cdot c_p$ before feeding it to the ControlNet. In this way, the importance of visible features can be adjusted as the diffusion steps forward. The time embedding used for the gating operation is shared with the time embedding for encoding the diffusion step in the stable diffusion.

To encourage the object encoder to capture shape information, we introduce an auxiliary path to predict the complete object mask from the partial token c_p . Specifically, a feature pyramid network (Lin et al., 2017) is leveraged as the mask decoder which takes c_p and the multi-scale features from the object encoder as input and outputs a pre-diffusion mask M_{pre} . We encourage mask completion with supervision as

$$\mathcal{L}_{mask} = \mathcal{L}_{dice}(M_c, M_{pre}) + \lambda_{ce} \mathcal{L}_{ce}(M_c, M_{pre}) \quad (3)$$

where \mathcal{L}_{dice} and \mathcal{L}_{ce} are Dice loss (Li et al., 2019) and BCE loss respectively. λ_{ce} is a constant.

3.4. Segmentation Stage

In the segmentation stage, illustrated in Fig. 5 (a), our approach initiates by sampling N images denoted as $\{I_t^{(k)}\}_{k=1}^N$ from the generative model, where t is the IMD step. Subsequently, we employ an off-the-shelf object segmentation model denoted as $\mathcal{S}(\cdot)$ to obtain the shapes (object masks) $\{M_t^{(k)}\}$ from these sampled images.

To derive an improved mask for the subsequent IMD step, we seek a function $\mathcal{V}(\cdot)$ that can produce a high-quality mask prediction from the set of N generated masks. Interestingly, though the distribution of sampled images is complex, we notice the distribution of masks has good properties. In Fig. 5 (b), we provide a visualization of the probability map

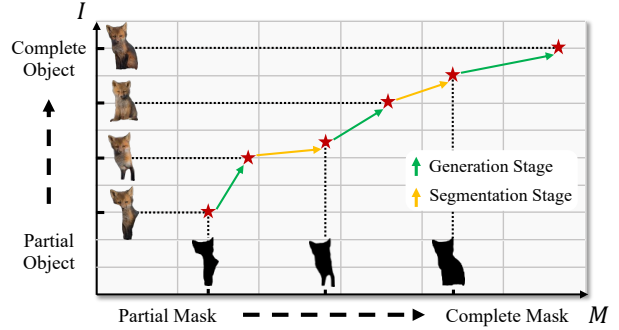


Figure 6. Mutual-beneficial sampling.

associated with a set of object masks with the same conditions, which is computed by taking the normalized average of the masks. To enhance the visualization of this probability distribution, we focus on a specific cross-section of the fully occluded portion in image I_p (the lower leg, represented as a yellow section) and visualize the probability as a function of the horizontal coordinate which demonstrates an obvious unimodal and symmetric property. Leveraging this observation, we can find an improved mask by taking the high-probability region. The updating can be achieved by conducting a voting process across the N estimated masks, as defined by the following equation:

$$\hat{M}_t[i, j] = \begin{cases} 1, & \text{if } \frac{\sum_{k=1}^N M_t^{(k)}[i, j]}{N} \geq \tau \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where $[i, j]$ denotes the coordinate, and τ is the threshold employed for the mask voting process.

3.5. Discussion

In this section, we will omit the conditioned partial image I_p for simplicity.

Joint modeling of mask and object. In practical scenarios where the complete object mask M_c is unavailable, modeling object completion through a marginal probability $p(I_c | M_c)$ becomes infeasible. Instead, it necessitates the more challenging joint modeling of objects and masks, denoted as $p(I, M)$, where the images and masks can range from partial to complete. Let us understand the joint distribu-

Method	AHP (Zhou et al., 2021a)				DYCE (Ehsani et al., 2018)			
	FID-G ↓	FID-S ↓	Rank ↓	Best ↑	FID-G ↓	FID-S ↓	Rank ↓	Best ↑
ControlNet	40.2	45.4	3.4	0.10	42.4	49.4	3.4	0.08
Kandinsky 2.1	43.9	39.2	3.2	0.11	44.3	47.7	3.4	0.06
Stable Diffusion 1.5	35.7	41.4	3.2	0.12	31.2	43.4	3.4	0.11
Stable Diffusion 2.1	30.8	39.9	3.1	0.14	30.0	41.1	3.0	0.12
MaskComp (Ours)	16.9	21.3	2.1	0.53	20.0	25.4	1.9	0.63

Table 2. **Quantitative evaluation on object completion task.** The computing of FID-G and FID-S only considers the object areas within ground truth and foreground regions segmented by SAM, respectively, to eliminate the influence of the generated background. The Rank denotes the average ranking in the user study. The Best denotes the percentage of samples that are ranked as the best. ↓ and ↑ denote the smaller the better and the larger the better respectively.

tion by exploring its marginals. Since the relation between mask and image is one-to-many (each object image only has one mask while the same mask can be segmented from multiple images), the $p(M|I)$ is actually a Dirac delta distribution δ and only the $p(I|M)$ is a real distribution. This way, the joint distribution of mask and image is discrete and complex, making the modeling difficult. To address this issue, we introduce a slack condition to the joint distribution $p(I, M)$ that *the mask and image can follow a many-to-many relation*, which makes its marginal $p(M|I)$ a real distribution and permits $p(I|M)$ to predict image I that has a different shape as the conditioned M and vice versa.

Mutual-beneficial sampling. After discussing the joint distribution that we are targeting, we introduce the mathematical explanation of MaskComp. MaskComp introduces the alternating modeling of two marginal distributions $p(I|M)$ (generation stage) and $p(M|I)$ (segmentation stage), which is actually a Markov Chain Monte Carlo-like (MCMC-like) process and more specifically Gibbs sampling-like. It samples the joint distribution $p(I, M)$ by iterative sampling from the marginal distributions. Two core insights are incorporated in MaskComp: (1) providing a mask as a condition can effectively enhance object generation and (2) fusing the mask of generated object images can result in a more accurate and complete object mask. Based on these insights, we train CompNet to maximize $p(I|M)$ and leverage mask voting to maximize the $p(M|I)$. As shown in Fig. 6, MaskComp develops a mutual-beneficial sampling process from the joint distribution $p(I, M)$, where the object mask is provided to boost the image generation and, in return, the generated images can lead to a more accurate mask by fusing the segmentation of images. Through alternating sampling from the marginal distributions, we can effectively address the object completion task.

4. Experiment

4.1. Experimental Settings

Dataset. We evaluate MaskComp on two popular datasets: AHP (Zhou et al., 2021a) and DYCE (Ehsani et al., 2018). AHP is an amodal human perception dataset with 56,302

images with annotations of integrated humans. DYCE is a synthetic dataset with photo-realistic images and the natural configuration of objects in indoor scenes. For both datasets, the non-occluded object and its corresponding mask for each object are available. We train MaskComp on the AHP and a filtered subset of OpenImage v6 (Kuznetsova et al., 2020). OpenImage is a large-scale dataset offering heterogeneous annotations. We select a subset of OpenImage that contains 429,358 objects as a training set of MaskComp.

Evaluation metrics. In accordance with previous methods (Zhou et al., 2021a), we evaluate image generation quality Fréchet Inception Distance (FID). The background is removed with object masks before evaluation. As the FID score cannot reflect the object completeness, we further conduct a user study, leveraging human assessment to compare the quality and completeness of images. During the assessment, given a partial object, the participants are required to rank the generated object from different methods based on their completeness and quality. We calculate the averaged ranking and the percentage of the image being ranked as the first place as the metrics (details available in the Appendix).

Implementation details. For the generation stage, we train the CompNet with frozen Stable Diffusion (Rombach et al., 2022) on the AHP dataset for 50 epochs. The learning rate is set for $1e-5$. We adopt batchsize = 8 and an Adam (Loshchilov & Hutter, 2017) optimizer. The image is resized to 512×512 for both training and inference. The object is cropped and resized to have the longest side 360 before sticking on the image. For a more generalized setting, we train the CompNet on a subset of the OpenImage (Kuznetsova et al., 2020) dataset for 36 epochs. We generate text prompts using BLIP (Li et al., 2022a) for all experiments (prompts are necessary to train ControlNet). For the segmentation stage, we leverage SAM (Kirillov et al., 2023) as $\mathcal{S}(\cdot)$. We vote mask with a threshold of $\tau = 0.5$. During inference, if no other specification, we conduct the IMD process for 5 steps with $N = 5$ images for each step. We give the class label as the text prompt to facilitate the CompNet to effectively generate objects. All baseline methods are given the same text prompts during the experiments. More implementation details are available in the appendix. The

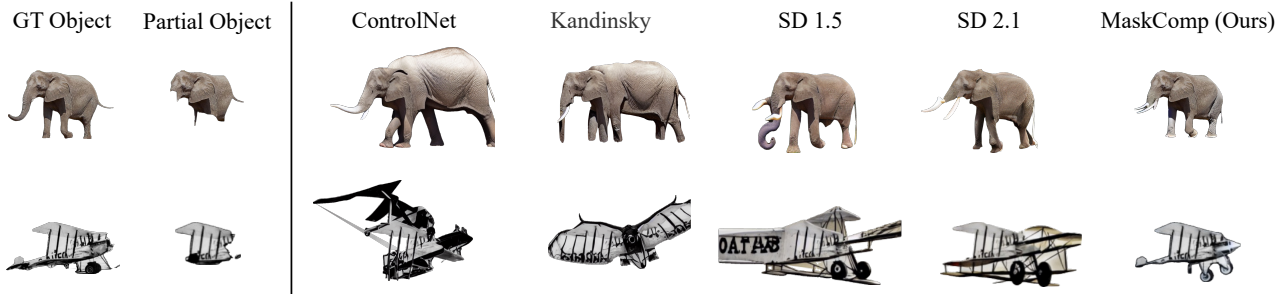


Figure 7. **Qualitative comparison against ControlNet, Kandinsky and Stable Diffusion.** The background is filtered out for better visualization. More results are available in the Appendix.

Mask	Partial	Intermed.	Complete	Occ.	20%	40%	60%	80%	Comp.	Gen.	Segm.	Total	Model	Baseline	MaskComp
FID	16.9	15.3	12.7	FID	13.4	15.7	17.2	29.9	Second	14.3	1.2	15.5	FID	29.4	16.9
(a) Conditioned mask.			(b) Occlusion rate.				(c) Inference time.			(d) Amodal baseline.					

Table 3. **Ablation of MaskComp on AHP dataset.** We ablate (a) the different conditioning masks during inference, (b) the occlusion rate during inference, (c) the inference time of each component in an IMD step, and (d) the performance compared with the amodal baseline.

code will be made publicly available.

4.2. Main Results

Quantitative results. We compare the MaskComp with state-of-the-art methods, ControlNet (Zhang et al., 2023), Kandinsky 2.1 (Shakhmatov et al., 2023), Stable Diffusion 1.5 (Rombach et al., 2022) and Stable Diffusion 2.1 (Rombach et al., 2022) on AHP (Zhou et al., 2021a) and DYCE (Ehsani et al., 2018) dataset. The results in Table 2 indicate that MaskComp consistently outperforms other methods, as evidenced by its notably lower FID scores, signifying the superior quality of its generated content. We conducted a user study to evaluate object completeness in which participants ranked images generated by different approaches. MaskComp achieved an impressive average ranking of 2.1 and 1.9 on the AHP and DYCE datasets respectively. Furthermore, MaskComp also generates the highest number of images ranked as the most complete and realistic compared to previous methods. We consider the introduced mask condition and the proposed IMD process benefits the performance of MaskComp, where the additional conditioned mask provides robust shape guidance to the generation process and the proposed iterative mask denoising process refines the initial conditioned mask to a more complete shape, further enhancing the generated image quality.

Qualitative results. We present visual comparisons among ControlNet, Kandinsky 2.1, Stable Diffusion 1.5, and Stable Diffusion 2.1, illustrated in Fig. 7. Our visualizations showcase MaskComp’s ability to produce realistic and complete object images given partial images as the condition, whereas previous approaches exhibit noticeable artifacts and struggle to achieve realistic object completion. In addition, without mask guidance, it is common for previous methods

to generate images that fail to align with the partial object.

4.3. Analysis

In this section, we provide an experimental analysis of MaskComp. All the results are evaluated with GT masks to filter out the background, i.e., FID-G.

Performance with different mask conditions. We evaluated the quality of generated images when conditioned on the same partial images along with three distinct types of masks: (1) partial mask (mask of the partial image), (2) intermediate mask (less occlusion than partial), and (3) complete mask. As shown in Table 3a, the model achieves its highest performance when it is conditioned with complete object masks, whereas relying solely on partial masks yields less optimal results. These results provide strong evidence that the quality of the conditioned mask significantly influences the quality of the generated images.

Performance with different occlusion rates. We perform ablation studies to assess the resilience of MaskComp under varying occlusion levels. As presented in Table 3b, we evaluate MaskComp at different occlusion levels (proportion of the obscured area relative to the complete object) ranging from 20% to 80%, and the results indicate that its performance does not degrade significantly up to 60% occlusion.

Inference time. Table 3c reports the inference time of each component in IMD (with a single NVIDIA V100 GPU). Although MaskComp’s throughput is reduced due to the inclusion of multiple diffusion processes in each IMD step, it is capable of attaining a higher degree of accuracy in visual object completion. Based on our empirical experiments, reducing the number of diffusion steps during the first few IMD steps can increase model speed without sacri-

Model	CLIPSeg	SEEM	SAM	T	1	3	5	7	N	4	5	6	Gating	✓	✗
FID	19.9	18.1	16.9	FID	24.7	19.4	16.9	16.1	FID	17.4	16.9	16.8	FID	16.9	18.2
(a) Segmentation model.				(b) IMD step number.					(c) # of sampled images.				(d) Condition gating.		

Table 4. **Design choices for IMD on AHP dataset.** We ablate (a) the impact of different segmentation networks, (b) IMD step number, (c) the number of sampled images in the segmentation stage, and (d) the gating operation in the CompNet.

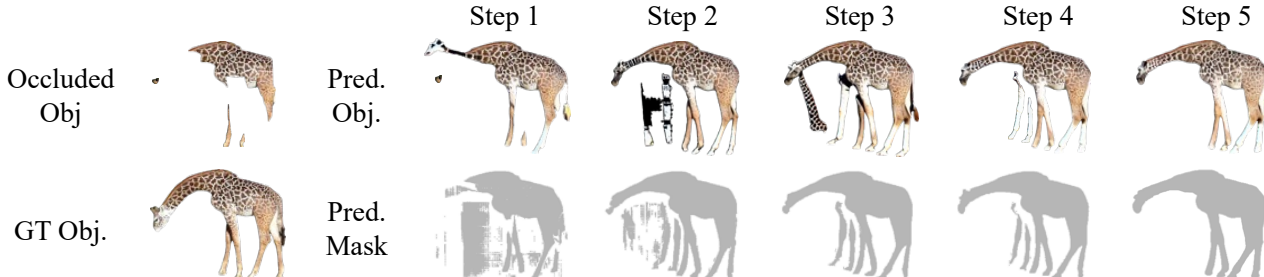


Figure 8. **Visualization of the IMD process.** For each step, we randomly demonstrate one generated image and the averaged mask for all generated images. We omit the input mask which has the same shape as the input occluded object.

Noise degree	Iter. 1	Iter. 3	Iter. 5	Iter. 7	Iter. 9
15% area	28.4	22.7	18.9	17.2	16.5
10% area	26.4	21.4	18.1	17.0	16.4
5% area	24.9	19.6	17.0	16.2	16.0
No noise	24.7	19.4	16.9	16.1	15.9

Table 5. Performance against segmentation errors on AHP dataset.

ficing much performance. With this idea incorporated into MaskComp, the average running time could be reduced to 2/3 of the original time with FID slightly increasing by 0.50. While beyond the scope of this study, we expect more advanced techniques could be explored to optimize the tradeoff between model speed and performance.

Comparison to amodal segmentation baseline. Amodal segmentation has a similar objective to the proposed IMD process. To demonstrate the effectiveness of MaskComp, we construct an amodal baseline that generates amodal masks from the SOTA amodal segmentation method (Tran et al., 2022) and then utilize ControlNet to generate images based on the amodal masks. As shown in Table 3d, we notice that our method outperforms the amodal baseline by a considerable margin, which could be attributed to the strong mask completion capability of the proposed IMD process.

Impact of different segmentation networks. We adopt SAM to obtain object masks at the segmentation stage. To study the impacts of different segmenters, we replace SAM with two smaller segmentation networks, CLIPSeg (Lüddecke & Ecker, 2022) and SEEM (Zou et al., 2023). Table 4a shows that the FID score with CLIPSeg (19.9) is slightly higher than with SAM (16.9), but remains competitive against other state-of-the-art methods, e.g., Stable Diffusion 2.1 (30.8 reported in Table 2). MaskComp is an iterative mask denoising (IMD) process that progressively

refines a partial object mask to boost image generation. The results support our hypothesis that the impact of the segmenter is modest.

Design choices in IMD. We conduct experiments to ablate the design choices in IMD and their impacts on the completion performance. We first study the effect of IMD step number. With a larger step number, IMD can better advance the partial mask to the complete mask. As shown in Table 4b, we notice that the image quality keeps increasing and slows down at a step number of 5. In this way, we choose 5 as our IMD step number. After that, we ablate the number of sampled images in the segmentation stage in Table 4c. We notice more sampled images generally leading to a better performance. We leverage an image number of 5 with the efficiency consideration. As we leverage the diffusion-based method for image generation, we ablate the iterations for the diffusion process. As shown in Table 4d, we notice the gating operation improves the generation quality by 1.3 FID, indicating the necessity of conditional gating.

Robustness to segmentation errors. We conduct experiments to manually add random errors to masks. As shown in Table 5, we ablate on the number of iterations and the degree of segmentation error. We observe that segmentation errors will increase the convergence iteration number while not affecting the final performance significantly. As IMD is a reciprocal process intended to provide effective control for later-generated masks to be refined based on adaptive feedback, mask errors are mitigated and not propagated.

Visualization of iterative mask denoising. To provide a clearer depiction of the IMD process, as depicted in Fig. 8, we present visualizations of the generated image and the averaged mask for each step. In the initial step, we observe the emergence of artifacts alongside the object. As we

progress through the steps, both the image and mask quality exhibit continuous improvement.

5. Conclusion

In this paper, we introduce MaskComp, a novel approach for object completion. MaskComp addresses the object completion task by seamlessly integrating conditional generation and segmentation, capitalizing on the crucial observation that the quality of generated objects is intricately tied to the quality of the conditioned masks. We augment the object completion process with an additional mask condition and propose an iterative mask denoising (IMD) process. This iterative approach gradually refines the partial object mask, ultimately leading to the generation of satisfactory objects by leveraging the complete mask as a guiding condition. Our extensive experiments demonstrate the robustness and effectiveness of MaskComp, particularly in challenging scenarios involving heavily occluded objects.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In *ECCV*, 2020. 3
- Cen, J., Zhou, Z., Fang, J., Shen, W., Xie, L., Zhang, X., and Tian, Q. Segment anything in 3d with nerfs. *arXiv preprint arXiv:2304.12308*, 2023. 3
- Chang, Z., Weng, S., Zhang, P., Li, Y., Li, S., and Shi, B. L-coins: Language-based colorization with instance awareness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 19221–19230, June 2023. 1
- Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 3
- Cheng, B., Collins, M. D., Zhu, Y., Liu, T., Huang, T. S., Adam, H., and Chen, L.-C. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 3
- Cheng, B., Schwing, A. G., and Kirillov, A. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021. 3
- Cheng, B., Misra, I., Schwing, A. G., Kirillov, A., and Girdhar, R. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 3
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009. 15
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 3
- Ehsani, K., Mottaghi, R., and Farhadi, A. Segan: Segmenting and generating the invisible. In *CVPR*, 2018. 6, 7
- Gafni, O., Polyak, A., Ashual, O., Sheynin, S., Parikh, D., and Taigman, Y. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pp. 89–106. Springer, 2022. 2
- Gal, R., Patashnik, O., Maron, H., Bermano, A. H., Chechik, G., and Cohen-Or, D. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 3
- Galatolo, F. A., Cimino, M. G., and Vaglini, G. Generating images from caption and vice versa via clip-guided generative latent space search. *arXiv preprint arXiv:2102.01645*, 2021. 3
- Gu, S., Bao, J., Yang, H., Chen, D., Wen, F., and Yuan, L. Mask-guided portrait editing with conditional gans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3436–3445, 2019. 2
- Guo, D., Zhao, H., Cheng, Y., Zheng, H., Gu, Z., and Zheng, B. Painting from part. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14779–14788, 2021. 2
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. Mask r-cnn. In *ICCV*, 2017. 3
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3
- Ho, J., Saharia, C., Chan, W., Fleet, D. J., Norouzi, M., and Salimans, T. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23(47):1–33, 2022. 2
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C., Lo, W.-Y., et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 3, 6

- Koley, S., Bhunia, A. K., Sain, A., Chowdhury, P. N., Xiang, T., and Song, Y.-Z. Picture that sketch: Photorealistic image generation from abstract sketches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6850–6861, June 2023. 1, 2
- Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont-Tuset, J., Kamali, S., Popov, S., Mallocci, M., Kolesnikov, A., et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7):1956–1981, 2020. 6
- Lee, D., Kim, C., Kim, S., Cho, M., and Han, W.-S. Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11523–11532, 2022. 2
- Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022a. 6
- Li, X., Sun, X., Meng, Y., Liang, J., Wu, F., and Li, J. Dice loss for data-imbalanced nlp tasks. *arXiv preprint arXiv:1911.02855*, 2019. 5
- Li, X., Wang, J., Li, X., and Lu, Y. Video instance segmentation by instance flow assembly. *IEEE Transactions on Multimedia*, 2022b. 3
- Li, X., Cao, H., Zhao, S., Li, J., Zhang, L., and Raj, B. Panoramic video salient object detection with ambisonic audio guidance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 1424–1432, 2023a. 3
- Li, X., Lin, C.-C., Chen, Y., Liu, Z., Wang, J., and Raj, B. Paintseg: Training-free segmentation via painting. *arXiv preprint arXiv:2305.19406*, 2023b. 3
- Li, X., Wang, J., Xu, X., Li, X., Raj, B., and Lu, Y. Robust referring video object segmentation with cyclic structural consensus. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22236–22245, 2023c. 3
- Li, X., Wang, J., Xu, X., Yang, M., Yang, F., Zhao, Y., Singh, R., and Raj, B. Towards noise-tolerant speech-referring video object segmentation: Bridging speech and text. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2283–2296, 2023d. 3
- Li, X., Wen, Y., Yang, M., Wang, J., Singh, R., and Raj, B. Rethinking voice-face correlation: A geometry view. *arXiv preprint arXiv:2307.13948*, 2023e. 3
- Li, Y., Cheng, Y., Gan, Z., Yu, L., Wang, L., and Liu, J. Bachgan: High-resolution image synthesis from salient object layout. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8365–8374, 2020. 2
- Li, Y., Liu, H., Wu, Q., Mu, F., Yang, J., Gao, J., Li, C., and Lee, Y. J. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22511–22521, 2023f. 2
- Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B., and Belongie, S. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2117–2125, 2017. 5
- Liu, X., Park, D. H., Azadi, S., Zhang, G., Chopikyan, A., Hu, Y., Shi, H., Rohrbach, A., and Darrell, T. More control for free! image synthesis with semantic diffusion guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 289–299, 2023. 3
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021. 15
- Long, J., Shelhamer, E., and Darrell, T. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 3
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- Lüddecke, T. and Ecker, A. Image segmentation using text and image prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7086–7096, 2022. 8
- Ma, J. and Wang, B. Segment anything in medical images. *arXiv preprint arXiv:2304.12306*, 2023. 3
- Neven, D., De Brabandere, B., Proesmans, M., and Van Gool, L. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *CVPR*, 2019. 3
- Newell, A., Huang, Z., and Deng, J. Associative embedding: End-to-end learning for joint detection and grouping. In *NeurIPS*, 2017. 3
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 3

- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022. 4, 6, 7, 14, 15, 16
- Sehwag, V., Hazirbas, C., Gordo, A., Ozgenel, F., and Canton, C. Generating high fidelity data from low-density regions using diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11492–11501, 2022. 2
- Shakhmatov, A., Razzhigaev, A., Nikolich, A., Arkhipkin, V., Pavlov, I., Kuznetsov, A., and Dimitrov, D. kandinsky 2.1, 2023. 7
- Sun, W. and Wu, T. Image synthesis from reconfigurable layout and style. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10531–10540, 2019. 2
- Tran, M., Vo, K., Yamazaki, K., Fernandes, A., Kidd, M., and Le, N. Aisformer: Amodal instance segmentation with transformer. *arXiv preprint arXiv:2210.06323*, 2022. 8
- Wan, Z., Zhang, J., Chen, D., and Liao, J. High-fidelity pluralistic image completion with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4692–4701, 2021. 2
- Wang, H., Zhu, Y., Green, B., Adam, H., Yuille, A., and Chen, L.-C. Axial-DeepLab: Stand-alone axial-attention for panoptic segmentation. In *ECCV*, 2020a. 3
- Wang, H., Zhu, Y., Adam, H., Yuille, A., and Chen, L.-C. MaX-DeepLab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021. 3
- Wang, X., Kong, T., Shen, C., Jiang, Y., and Li, L. SOLO: Segmenting objects by locations. In *ECCV*, 2020b. 3
- Xie, S., Zhang, Z., Lin, Z., Hinz, T., and Zhang, K. Smart-brush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22428–22437, June 2023. 1
- Yan, X., Wang, F., Liu, W., Yu, Y., He, S., and Pan, J. Visualizing the invisible: Occluded vehicle segmentation and recovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7618–7627, 2019. 2
- Yang, B., Gu, S., Zhang, B., Zhang, T., Chen, X., Sun, X., Chen, D., and Wen, F. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18381–18391, 2023a. 1
- Yang, Z., Wang, J., Gan, Z., Li, L., Lin, K., Wu, C., Duan, N., Liu, Z., Liu, C., Zeng, M., and Wang, L. Reco: Region-controlled text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14246–14255, June 2023b. 1
- Ye, H., Kuen, J., Liu, Q., Lin, Z., Price, B., and Xu, D. Seggen: Supercharging segmentation models with text2mask and mask2img synthesis. *arXiv preprint arXiv:2311.03355*, 2023. 2
- Yu, T., Feng, R., Feng, R., Liu, J., Jin, X., Zeng, W., and Chen, Z. Inpaint anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023. 3
- Yu, Y., Zhan, F., Lu, S., Pan, J., Ma, F., Xie, X., and Miao, C. Wavefill: A wavelet-based generation network for image inpainting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 14114–14123, 2021. 2
- Zhang, L., Rao, A., and Agrawala, M. Adding conditional control to text-to-image diffusion models, 2023. 1, 2, 4, 5, 7
- Zhang, W., Pang, J., Chen, K., and Loy, C. C. K-Net: Towards unified image segmentation. In *NeurIPS*, 2021. 3
- Zhao, B., Meng, L., Yin, W., and Sigal, L. Image generation from layout. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8584–8593, 2019. 2
- Zhou, Q., Wang, S., Wang, Y., Huang, Z., and Wang, X. Human de-occlusion: Invisible perception and recovery for humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3691–3701, 2021a. 6, 7
- Zhou, Y., Zhang, R., Chen, C., Li, C., Tensmeyer, C., Yu, T., Gu, J., Xu, J., and Sun, T. Lafite: Towards language-free training for text-to-image generation. *arXiv preprint arXiv:2111.13792*, 2021b. 3
- Zou, X., Yang, J., Zhang, H., Li, F., Li, L., Gao, J., and Lee, Y. J. Segment everything everywhere all at once. *arXiv preprint arXiv:2304.06718*, 2023. 3, 8

A. More Experiments

In this section, we provide more ablation experiments and analysis of MaskComp. We conducted ablation experiments to determine the design choice in the segmentation stage.

Iter	20	40	50	Occ.	Rectangle	Oval	Object	Comp.	✓	✗
FID	16.9	15.7	15.1	FID	15.3	15.1	16.9	FID	16.9	19.4
(a) Iteration for diffusion.			(b) Occlusion type.			(c) Availability of complete object.				
Strategy	Logits (V)	Logits (M)	Mask (V)	Mask (M)	\mathcal{L}_{mask}	✓	✗			
FID	16.9	17.2	17.6	17.0	FID	16.9	17.7			
(d) Voting strategies.					(e) Mask loss.					

Table 6. **More ablation of MaskComp.** We report the performance with the AHP dataset. (a) We ablate the iteration number of the diffusion model. (b) We report the performance with different types of occlusion. (c) We report the performance of MaskComp trained with or without the complete objects. (d) We ablate voting strategies. V: voting. M: Mean. (e) We ablate the effectiveness of adding intermediate supervision to predict the complete mask.

Iteration for diffusion model. Since the number of diffusion steps has a large impact on the inference speed, we conduct the ablation studies about iteration steps for the diffusion process in Table 6a. We notice a larger diffusion step will lead to a better performance. After the number of diffusion steps is larger than 40, the performance improvement becomes slow.

Occlusion type To understand the influence of occlusion type, we conduct an ablation study as shown in Table 6b. Three types of occlusion types are employed. Specifically, the occlusions are randomly generated by controlling the size and location. For object occlusion, we occlude the object by the mask of itself. We notice that the occlusion with a more complicated object shape will impose more challenges on the proposed model.

Availability of complete objects during training Occlusion is a prevalent occurrence in images, posing a significant challenge for object completion, primarily due to the unavailability of ground-truth complete objects. This motivates us to investigate the performance of MaskComp without ground-truth complete objects during training. As MaskComp relies on a mask denoising process that does not necessarily request the complete objects during training, it is possible to adapt MaskComp to the scenarios without complete objects available during training. As shown in Table 6c, we report the performance on AHP dataset with the models trained on AHP (with complete objects) and OpenImage (without complete objects) respectively. We notice that the performance of MaskComp trained on OpenImage is just slightly lower than that trained with AHP dataset, indicating that MaskComp has the potential to be adapted to the scenarios without ground-truth complete objects. More qualitative comparisons are available in the Section B.

Voting strategies. To investigate the impact of different voting strategies in the segmentation stage, we conduct experiments to evaluate different voting approaches as shown in Table 6d. We notice voting with logits achieves the best performance. The current design choice of using SAM and voting with logits is based on the ablation results.

Mask loss. We leverage auxiliary mask prediction from the partial token c_p to encourage the object encoder to capture object shape information. As shown in Table 6e, we ablate the effectiveness of adding additional mask supervision. The results indicate that the incorporation of mask prediction can benefit the final object completion performance.

B. More Discussion

Image diffusion v.s. Mask denoising. During the training of the image diffusion model, Gaussian noise is introduced to the original image. A denoising U-Net is then trained to predict this noise and subsequently recover the image to its clean state during inference.

Type	Noise	Network	Target
Image diffusion	Gaussian	UNet	Predict added noise
Mask denoising	Occlusion	Mask denoiser $\mathcal{S} \circ \mathcal{G}(\cdot)$	Predict denoised mask

Table 7. Analogy between image diffusion and mask denoising.

Similarly, in the context of the proposed iterative mask denoising (IMD) process, we manually occlude the complete object (which can be assumed as adding noise) and train a generative model to recover the complete object. During inference, as shown in Eq. (1), we employ an iterative approach that combines the segmentation and generation model $\mathcal{S} \circ \mathcal{G}(\cdot)$ functioning as a denoiser. This denoiser progressively denoises the partial mask to achieve a complete mask, following a similar principle to the denoising diffusion process. By drawing parallels between image diffusion and mask denoising, we establish an analogy, as depicted in Table 7. We can notice that the mask-denoising process shares the spirits of the image diffusion process and the only difference is that mask denoising does not explicitly calculate the added noise but directly predicts the denoised mask. In this way, MaskComp can be assumed as a double-loop denoising process with an inner loop for image denoising and an outer loop for mask denoising.

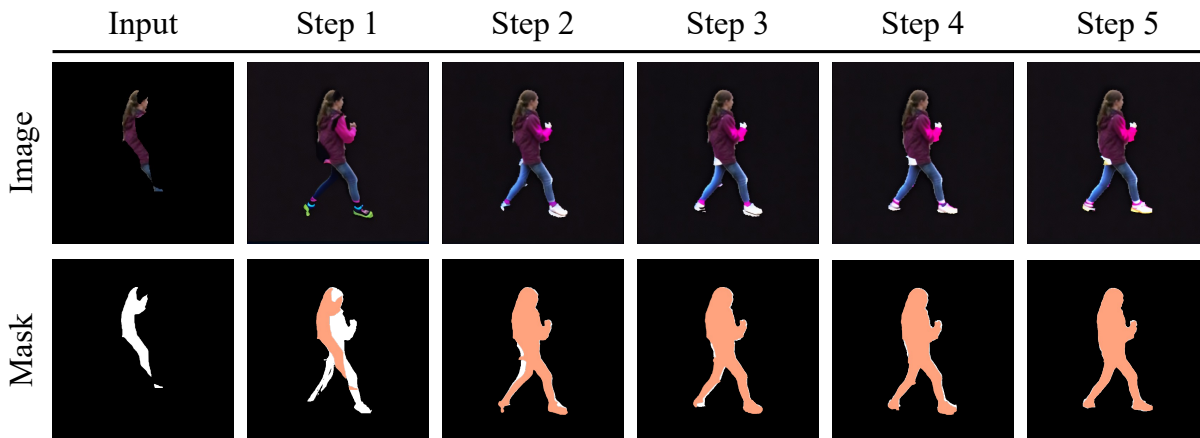


Figure 9. Visualization of IMD process with model trained without complete objects. To better visualize the iterative mask denoising process, we denote the overlapping masked area from the last iteration as orange. We can notice that the object shape is gradually refined and converged to a complete shape.

Training without complete object. In the context of image diffusion, though multiple forward steps are involved to add noise to the image, the network only learns to predict the noise added in a single step during training. Therefore, if we possess a set of noisy images generated through forward steps, the original image is not required during the training. This motivates us to explore the feasibility of training MaskComp without relying on the complete mask. Similar to image diffusion, given a partial mask, we can further occlude it and learn to predict the partial mask before further occlusion. In this way, MaskComp can be leveraged in a more generic scenario without the strict demand for complete objects. We have discussed the quantitative results in Section 4.3. Here, we visualize the IMD process with a model trained without complete objects (on OpenImage). To better visualize the object shape updating, we denote the overlapping masked area from the last step as orange. We can notice that the object shape gradually refines and converges to the complete shape as the IMD process forwards. Interestingly, the IMD process can learn to complete the object even if only a small portion of the complete object was available in the dataset during the training. We consider this property to make it possible to further generalize MaskComp to the scenarios in which a complete object is not available.

What will the marginal distribution $p(I|M)$ and $p(M|I)$ be like without the slack condition? The relation between mask and object image is one-to-many. The $p(I|M)$ models a filling color operation that paints the color within the given mask area. And as each object image only corresponds to one mask, the $p(M|I)$ is a deterministic process that can be modeled by a delta function δ . Previous methods generally leverage the unslacked setting. For example, the ControlNet assumes the given mask condition can accurately reflect the object shape and therefore, it can learn to fill colors to the masked regions.

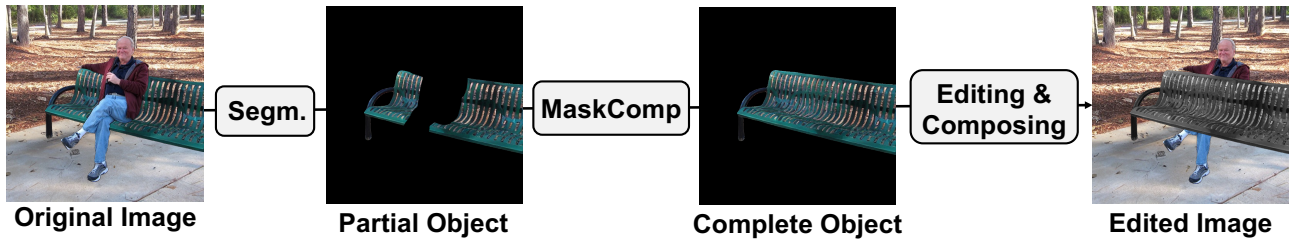


Figure 12. Illustration of potential application.

Analysis of the aggregation among masks. Given a partial object I_p , the generation model samples from a distribution that contains both realistic and unrealistic images. As shown in Fig. 10, the image distribution is typically complex and the expectation cannot represent a realistic image. However, when we consider the shape of the generated images, we are excited to find that the expectation leads to a more realistic shape. We consider this observation interesting as, for most of the other generation tasks (non-conditioned), averaging the object shape will just yield an unrealistic random shape. This observation serves as one of our core observations to build the IMD process. Here, SAM serves as the tool to extract the object shape and voting is a way to binarize the expectation of object shapes to a binary mask.

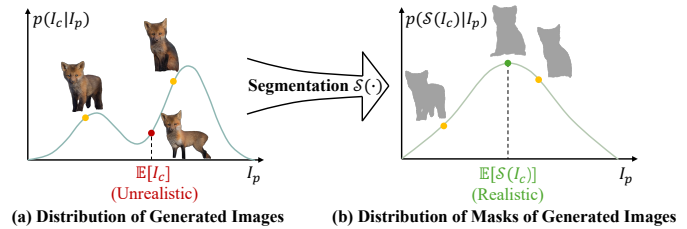


Figure 10. Distribution of image and mask.

Potential applications. Object completion is a fundamental technique that can boost a number of applications. For example, a straightforward application is the image editing. With the object completion, we can modify the layer of the objects in an image as we modify the components in the PowerPoint. It is possible to bring forward and edit objects as shown in Fig. 12. In addition, object completion is also an important technique for data augmentation. We hope MaskComp can shed light on more applications leveraging object completion.

Background objects in the generated images. The training of CompNet aims to learn an intra-object correlation. We leverage a black background to eliminate the influence of background objects. However, we notice that even if we train the network with the black background as ground truth, it is still possible to generate irrelevant objects in the background. As shown in Fig. 11, we visualize an image that generates a leather bag near the women. We consider the generated background object can result from the learned inter-object correlation from the frozen Stable Diffusion model (Rombach et al., 2022). As the generated background object typically will not be segmented in the segmentation stage, it will not influence the performance of MaskComp.

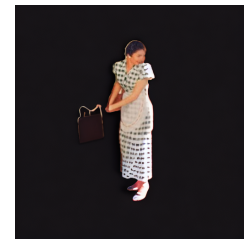


Figure 11. BG objects.

C. More Experiments

Failure case analysis. We present a failure case in Fig. 14, where MaskComp exhibits a misunderstanding of the pose of a person bending over, resulting in the generation of a hat at the waist. We attribute this generation of an unrealistic image to the uncommon pose of the partial human. Given that the majority of individuals in the AHP training set have their heads up and feet down, MaskComp may have a tendency to generate images in this typical position. We consider that with a more diverse dataset, including images of individuals in unusual poses, MaskComp could potentially yield superior results in handling similar cases.

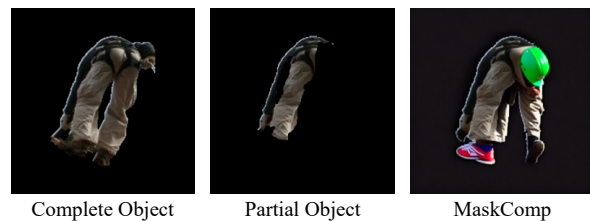


Figure 14. Failure case.

Impact of segmentation errors in intermediate stages.

Despite the robust capabilities of the CompNet and SAM models, they can still generate low-quality images and inaccurate segmentation results. In Fig. 15, we show a case where the intermediate stage of IMD produces a human with an extra right arm. To address this, we implement three key strategies: (1) **Error Mitigation during Segmentation with SAM:** As shown in Fig. 15, SAM effectively filters out incorrectly predicted components, such as a misidentified right arm, resulting in a more coherent shape for subsequent iterations.

SAM’s robust instance understanding capability extends to not only accurately segmenting objects with regular shapes but also filtering out irrelevant parts when additional objects/parts are generated. (2) **Error Suppression through Mask Voting:** In cases where only a few generated images exhibit errors, the impact of these errors can be mitigated through mask voting. The generated images are converted to masks, and if only a minority displays errors, their influence is diminished through the voting operation. (3) **Error Tolerance in IMD Iteration:** We train the CompNet to handle a wide range of occluded masks. Consequently, if the conditioned mask undergoes minimal improvement or degradation due to the noises in a given iteration, it can still be improved in the subsequent iteration. While this may slightly extend the convergence time, it is not anticipated to have a significant impact on the ultimate image quality.

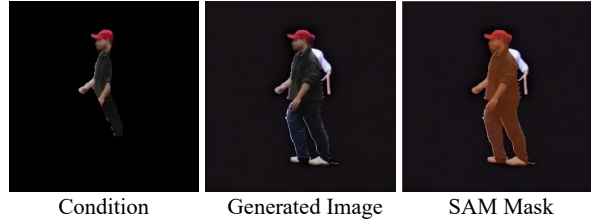


Figure 15. Visualization of the intermediate stage in IMD.

More implementation details. We leverage two types of occlusion strategies during the training of CompNet. First, we randomly sample a point on the object region, and then randomly occlude a rectangle area with the sampled point as the centroid. The width and height of the rectangle are determined by the width and height of the bounding box of the ground truth object. We uniformly sample a ratio within $[0.2, 0.9]$ and apply it to the ground truth width and height to occlude the object. Second, we randomly occlude the object by shifting its mask. Specifically, we randomly shift its mask by a range of $[0.17, 0.25]$ and occluded the region within the shifted mask. We equally leverage these two occlusion strategies during training. For the object encoder to extract partial token c_p in the CompNet, we utilize a Swin-Transformer (Liu et al., 2021) pre-trained on ImageNet (Deng et al., 2009) with an additional convolution layer to accept the concatenation of mask and image as input. We initialize the CompNet with the pre-trained weight of ControlNet with additional mask conditions. To segment objects in the segmentation stage, we give a mix of box and point prompts to the Segment Anything Model (SAM). Specifically, we uniformly sample three points from the partial object as the point prompts and we leverage an extended bounding box of the partial object as the box prompts. We also add negative point prompts at the corners of the box to further improve the segmentation quality.

More visualization. As shown in Fig. 13, we provide more qualitative comparisons with Stable Diffusion (Rombach et al., 2022). We notice that Stable Diffusion tends to complete irrelevant objects to the complete parts and thus leads to an unrealism of objects. Instead, MaskComp is guided by a mask shape and successfully captures the correct object shape thus achieving superior results.

Details of user study. There are 16 participants in the user study. All participants have relevant knowledge to understand the task. During the assessment, each participant is provided with instructions and an example to understand the task. We list the instructions as follows.

Task: Given the partial object (lower left), generate the complete object (upper left).

Instruction:

- Ranking images 1-5, put the best on the left and the worst on the right.
- Please focus on the foreground object and ignore the difference presented in the background.
- Original image is provided as a good example.
- The criteria for ranking are founded on object quality, encompassing aspects such as completeness, realism, sharpness, and more.
- It must be strictly ordered (no tie).
- Please rank the image in the following form: 1;2;3;4;5 or 5;4;3;2;1 (Use a colon to separate, no space at the beginning)

Completing Visual Objects via Bridging Generation and Segmentation

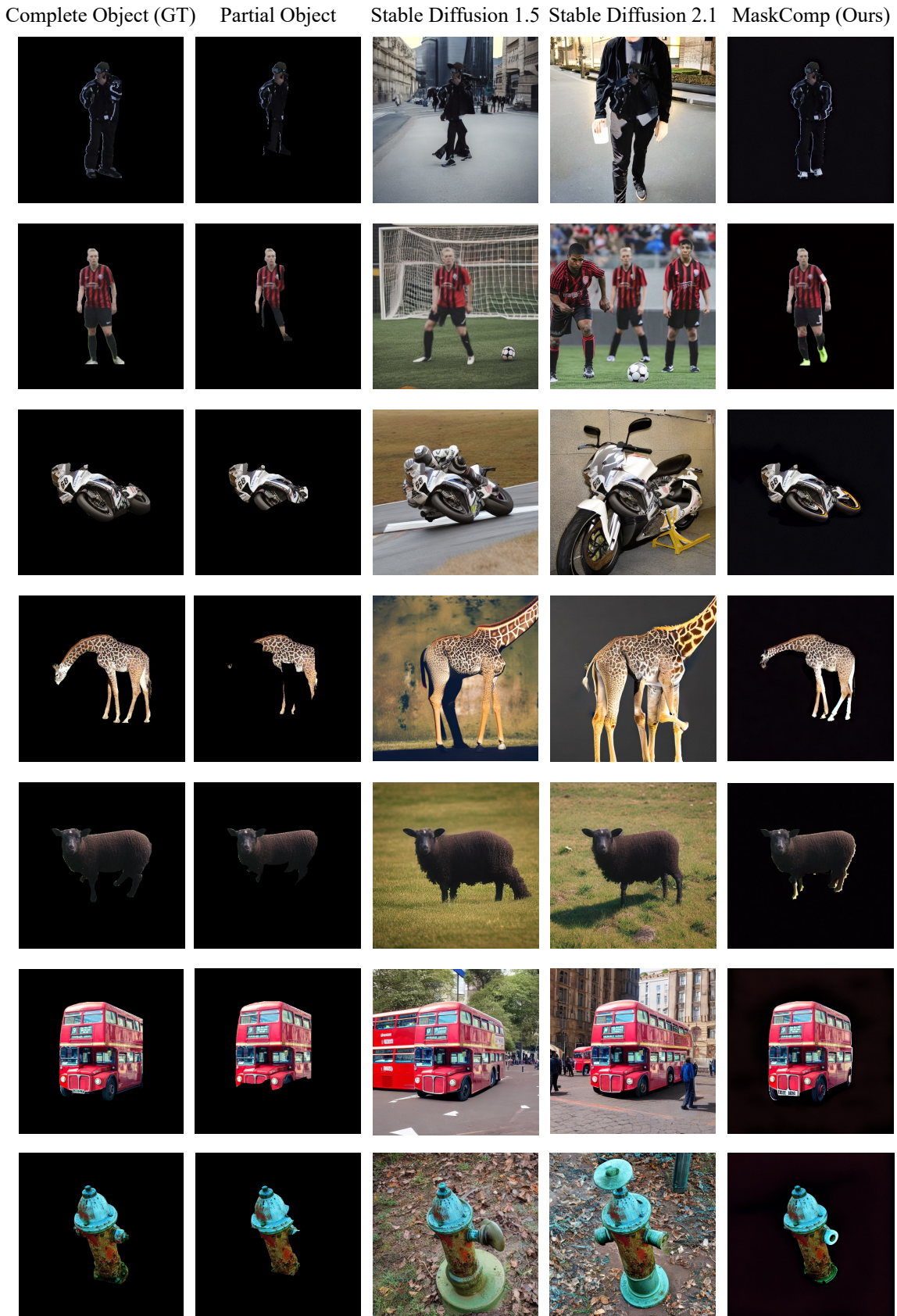


Figure 13. More qualitative comparison with Stable Diffusion (Rombach et al., 2022).