

SPATIALCOMPOSER: 3D SPATIAL OBJECT INSERTION VIA IMAGE GAUSSIAN COMPOSITION

Anonymous authors

Paper under double-blind review

ABSTRACT

With the rapid advancement of open-world image generation models in recent years, a series of image editing tasks have achieved excellent performance. However, considering object insertion as a representative example, this task still presents three primary challenges. First, the inserted object should maintain identity consistency with the reference object while preserving the original scene in non-edited regions. Second, the spatial position and scale of the inserted object should be reasonable and align with user expectations. Third, the inserted object should harmonize with other image components, typically involving object style and surface illumination harmonization. To address these challenges, we propose SpatialComposer, which leverages depth-aware image Gaussians to construct a spatially-structured scene representation from a single scene image and models object insertion as Gaussian composition, thereby achieving effective preservation of scene and object identity while enabling precise control over the scale and 3D spatial position of the inserted object. Subsequently, based on pre-trained diffusion generative models, we introduce a simple yet effective refinement method for the object harmonization process. By designating only the Gaussian components corresponding to the inserted object as trainable parameters, SpatialComposer avoids unintended modifications to other regions while simultaneously addressing both object-scene integration and scene detail preservation. Furthermore, recognizing that current object insertion benchmarks lack consideration for depth-aware position control, we construct a specialized benchmark featuring high-resolution scene images with substantial depth complexity. Comprehensive evaluations demonstrate that SpatialComposer achieves comparable or superior performance over state-of-the-art object insertion approaches across all three aforementioned challenges.

1 INTRODUCTION

Object insertion is a fundamental image editing task involving seamless integration of specified objects into target scene locations. Recent diffusion-based models (e.g., Stable Diffusion (Rombach et al., 2022b), FLUX (Labs, 2024)) have enabled significant evolution beyond previous methods (Tripathi et al., 2019; Zhang et al., 2020; Liu et al., 2021; Niu et al., 2022; Cong et al., 2020; Ling et al., 2021; Sofiiuk et al., 2021; Cong et al., 2022).

Current methods fall into two categories: *training-based methods* (Song et al., 2025; Wang et al., 2025; Chen et al., 2024b) that fine-tune pre-trained inpainting models conditioned on reference objects, and *training-free methods* (Chen et al., 2024c; Wang et al., 2024) that manipulate intermediate features and attention during inference. These methods enable intuitive personalized customization of real or synthesized images using arbitrary references, providing accessible automated tools for creating and modifying visual content.

Despite promising results, object insertion faces three main challenges: 1). *Object and scene consistency*. The forward diffusion process randomness and encoder-decoder downsampling information loss make maintaining detail consistency challenging for current methods. 2). *3D Spatial and scale controllability*. Existing approaches control insertion through object masks or bounding boxes with optional text prompts, specifying only two-dimensional positioning and failing to accurately control depth-related spatial position and scale. Our experiments show that even with depth text prompts,

054 previous methods often fail to meet expectations. 3). *Style and illumination harmony*. Current meth-
 055 ods struggle with style consistency, natural lighting, and color harmonization needed for seamless
 056 visual integration.

057 In this paper, we propose a novel depth-aware object insertion method based on Gaussian Kerbl
 058 et al. (2023) representation, which we term SpatialComposer. In order to tackle the above three
 059 challenges, our approach contains three main steps: Gaussian fitting, Gaussian composition, and
 060 Object refinement. In the first step, by leveraging a pre-trained monocular depth estimation network,
 061 combined with a back-projection Gaussian initialization strategy, we efficiently construct consistent
 062 object and scene Gaussian representations with meaningful spatial structure. In the second step,
 063 object insertion is implemented through the composition of object Gaussians with scene Gaussians.
 064 Through user-specified scaling and translation operations applied to the object Gaussians, our ap-
 065 proach achieves precise control over both the depth and scale of the inserted object. In the third step,
 066 following the composition of object and scene representations, we introduce a simple yet effective
 067 refinement method for inserted objects based on pre-trained inpainting diffusion models and illumi-
 068 nation harmonization models. With separable scene Gaussians and object Gaussians, we designate
 069 only the object Gaussian components as trainable parameters, thereby preventing unintended mod-
 070 ifications to other scene regions. This design circumvents the limitations of pre-trained diffusion
 071 generative models regarding resolution constraints and reconstruction fidelity. It also preserves the
 072 integrity of unmodified scene areas while harnessing the supervisory signals provided by pre-trained
 073 diffusion models for harmonizing inserted objects with the scene.

074 The existing TF-ICON benchmark (Lu et al., 2023b) for image-guided object insertion consists
 075 of generated 512×512 scene images with simple spatial structures, where insertions lack spatial
 076 relationships with scene components. This benchmark is insufficient for our evaluation. To fill this
 077 research gap, we collect a novel benchmark dataset comprising over 200 high-resolution scenes with
 078 complex structures and over 200 objects spanning 9+ categories including animals, vegetables, food,
 079 people, and vehicles. Scenes include over 100 real photographs and over 100 artistic images across
 080 multiple styles: oil painting, cartoon, anime, watercolor, and pixel art. We constructed over 200
 081 semantically coherent insertion cases. We name the dataset Depth-Aware Object Insertion (DAOI)
 082 Dataset and will release it soon.

083 Ablation experiments demonstrate the effectiveness of our image Gaussian representation and ini-
 084 tialization method. Using a single image and relative depth estimates, we rapidly construct scene
 085 Gaussians with meaningful depth while maintaining high-quality reconstruction. Leveraging Gaus-
 086 sian compositionality, scaling and translating object Gaussians enables precise control of scale and
 087 placement at any location. The proposed refinement module exhibits strong generalization and
 088 excellent performance without additional fine-tuning. Unlike existing methods, our approach han-
 089 dles ultra-high-resolution images while preserving fine details and demonstrates robust performance
 090 across diverse artistic styles, confirming practical applicability.

091 In summary, our main contributions include: (1) We propose a depth-aware image Gaussian rep-
 092 resentation and a back-projection initialize strategy to enhance the performance of Gaussian fitting
 093 while providing meaningful depth information. (2) To the best of our knowledge, we are the first to
 094 consider the precise control of 3D spatial position and scale for inserted objects in object insertion
 095 tasks. Based on our proposed representation, we provide a solution to this problem. This rep-
 096 resentation also achieves stronger preservation of non-edited scene regions and object identity. (3)
 097 We collect a high-resolution object insertion benchmark featuring complex scene spatial structures,
 098 providing a high-quality evaluation platform for advancing research in this domain. (4) Experimen-
 099 tal results demonstrate the superior performance and robustness of SpatialComposer across diverse
 100 stylistic and real scenes.

101 2 RELATED WORK

103 2.1 TRAINING-BASED OBJECT INSERTION

104 Built upon pre-trained text-to-image or inpainting models, training-based object insertion methods
 105 employ techniques such as LoRA (Hu et al., 2022) and Adapters (Houlsby et al., 2019) to fine-tune
 106 the models using paired data, i.e., scene images with and without specific objects. These approaches
 107 typically introduce new network structures to accept the scene image, reference object, and spa-

108 tial mask as conditions. Several works (Yang et al., 2023; Song et al., 2022; Chen et al., 2024a;
 109 Canet Tarrés et al., 2024; Yuan et al., 2024; Kulal et al., 2023; Zhang et al., 2023a; Song et al., 2024;
 110 He et al., 2024) extract reference features via visual encoders and trainable modules, influencing
 111 UNet through cross-attention, summation, or custom fusion. Some methods (Zhang et al., 2023b;
 112 Chen et al., 2024b) use composite images with pasted objects as conditions. Insert Anything and
 113 UniCombine (Song et al., 2025; Wang et al., 2025) adopt Diffusion Transformer (DiT) for generation
 114 control. For environmental harmony, ZeroComp (Zhang et al., 2025) conditions on depth, surface
 115 normal, albedo, and shading. Objectmate (Winter et al., 2024b) uses 2×2 grids as input to produce
 116 coherent insertion results. Anydoor (Chen et al., 2024b) combines ID extractors with high-pass
 117 filters for dual feature control. These methods require large paired datasets and have insufficient
 118 generalization ability. To avoid large dataset dependency, DreamCom and DreamEdit(Lu et al.,
 119 2023a; Li et al., 2023) fine-tune embeddings using DreamBooth(Ruiz et al., 2023) with few im-
 120 ages. DreamEdit(Li et al., 2023) uses DDIM Inversion(Mokady et al., 2023), while DreamCom(Lu
 121 et al., 2023a) employs masked attention control. OmniPaint(Yu et al., 2025) trains separate inser-
 122 tion/removal models with cycle consistency loss. To reduce paired data, ObjectDrop (Winter et al.,
 123 2024a) trains removal models on small datasets, then collects large synthetic datasets for insertion
 124 training. Overall, when there exists a significant gap between training and inference data, the preser-
 125 vation of non-edited scene regions and object identity, as well as the harmony between objects and
 126 scenes, exhibit substantial degradation. These methods are also unable to control the depth-related
 127 spatial positioning of object insertion. In contrast, SpatialComposer is training-free, enables pre-
 128 cise control over 3D spatial positioning, and has been validated to deliver stable performance across
 129 real-world scenes and a variety of stylistic scenes.

129 2.2 TRAINING-FREE OBJECT INSERTION

130
 131 Instead of parameter updating, training-free object insertion methods manipulate the intermediate
 132 features or attention maps during the inference stage. TF-ICON (Lu et al., 2023b) uses three-branch
 133 inference with DDIM-inverted noise from scene and reference images, along with composite noise
 134 placing resized object noise into target regions. Scene and object features guide cross-attention in
 135 the composite branch. Similarly,(Li et al., 2024) employs three-branch inference with self-attention
 136 fusion for object preservation and text-guided modification. FreeCompose(Chen et al., 2024c) iter-
 137 atively optimizes editing through key-value replacement and DDS-loss (Hertz et al., 2023). Prime-
 138 Composer (Wang et al., 2024) introduces a correlation diffuser computing cross-attention between
 139 target and reference features in UNet self-attention layers, with attention maps injected during de-
 140 noising for appearance preservation alongside region-constrained cross-attention. Constrained by
 141 the reconstruction and generation capabilities of the underlying foundation models, these methods
 142 exhibit limitations in preserving non-edited regions and maintaining object identity. Likewise, they
 143 do not provide effective control over the spatial placement of inserted objects.

144 3 METHOD

145
 146 The pipeline of SpatialComposer is illustrated in Fig. 1. It comprises three main components: Gaus-
 147 sian fitting, Gaussian composition and object refinement. We begin by introducing the founda-
 148 tional concepts of 3D Gaussians in Sec. 3.1, and we also introduce the latent diffusion model in
 149 Appendix A. Subsequently, we provide detailed descriptions of our proposed Depth-Aware Image
 150 Gaussian (DA-ImgGS) representation, its initialization, fitting and composition in Sec. 3.2 and the
 151 Diffusion-Based Object Refinement method in Sec. 3.3.

152 3.1 3D GAUSSIAN SPLATTING (3DGS)

153
 154 Our Gaussian representation builds upon standard 3D Gaussians with task-specific adaptations.
 155 Hence, we first introduce the foundational concepts of 3DGS in this section. 3DGS is an explicit
 156 scene representation using 3D Gaussians Kerbl et al. (2023), with each Gaussian \mathcal{G} defined by mean
 157 $\boldsymbol{\mu} \in \mathbb{R}^3$ and covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$:

$$158 \mathcal{G}(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}. \quad (1)$$

159
 160 Additionally, each Gaussian includes spherical harmonics (SH) $\mathbf{c} \in \mathbb{R}^k$ for color and opacity $\alpha \in \mathbb{R}$.
 161 The covariance matrix decomposes as $\Sigma = R S S^T R^T$, where $R \in \mathbb{R}^{3 \times 3}$ is the rotation matrix and

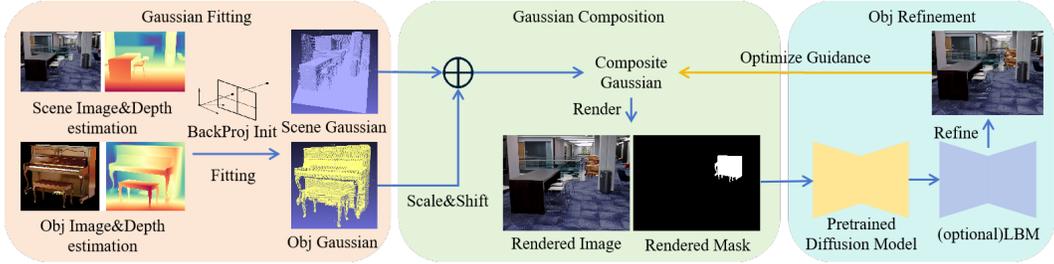


Figure 1: We first initialize and fit both scene and object images using our proposed depth-aware image Gaussian representation. The object is then scaled and positioned at the desired location through scaling and translation operations to generate the composed Gaussians. Subsequently, we employ a refinement method based on pre-trained diffusion models to optimize the object Gaussians, thereby achieving harmonization between the object and the scene.

$S = \text{diag}([s_x, s_y, s_z])$ is the diagonal scale matrix. The rotation matrix R is constructed from quaternion $\mathbf{v} = [r_w, r_x, r_y, r_z]$. For rendering, 3D Gaussians are projected onto pixel coordinates given camera pose W , with pixel-space covariance matrix defined as:

$$\Sigma_{pix} = JW\Sigma W^T J^T, \tag{2}$$

where J is the Jacobian matrix of the affine approximation of the projection transformation. The color at each pixel can then be obtained through alpha-blending of N overlapping Gaussians at that pixel in depth order:

$$\mathbf{c}_{pix} = \sum_i^N \mathbf{c}_i \alpha_i \prod_j^{i-1} (1 - \alpha_j), \tag{3}$$

where \mathbf{c}_i and α_i represent the color and density of each Gaussian at that pixel, weighted by the covariance matrix Σ . The differentiable rendering process enables end-to-end optimization of all Gaussian parameters based on image reconstruction loss.

3.2 DEPTH-AWARE IMAGE GAUSSIAN (DA-IMGGS)

We modify the parameter configuration, initialization, and optimization of standard 3D Gaussian to develop depth-aware image Gaussian. This adaptation accommodates our application requirements, enhances fitting performance on single images, while endowing the scene representation with usable spatial structure.

Standard 3DGS defines Gaussian means in world coordinates as 3D positions. During rendering, these world coordinates are transformed to camera coordinates through the camera’s extrinsic and intrinsic matrices before being projected onto the image plane, enabling image rendering from arbitrary viewpoints. However, for applications requiring rendering only from a fixed camera pose, we define our depth-aware image Gaussian directly in Normalized Device Coordinates (NDC). NDC represents a standardized coordinate system in computer graphics where coordinates are normalized to the range $[-1, 1]$. This coordinate system enables graphics rendering to adapt seamlessly across display devices with different resolutions and aspect ratios.

The covariance matrix in standard 3DGS represents an ellipsoid in 3D space, which projects to ellipses on different image planes during rendering. Since image Gaussians only require rendering onto a fixed image plane, we reduce the covariance matrix to 2D, where $\Sigma_{2d} \in \mathbb{R}^{2 \times 2}$ represents an elliptical Gaussian distribution parallel to the fixed image plane. In standard 3DGS, object colors observed from different viewpoints depend on multiple factors including viewing angle, lighting conditions, and material properties. Higher-order spherical harmonics decompose this lighting and material information into coefficient sets, providing efficient and smooth representations of complex lighting effects such as shadows, specular highlights, and viewpoint-dependent color variations. However, since image Gaussian rendering operates from a fixed viewpoint without 3D scene lighting variations, simple RGB values sufficiently capture the required surface colors. This reduction in Gaussian representation parameters improves scene fitting quality. Finally, the parameters of our

depth-aware image Gaussian include: mean $\mathbf{u} \in \mathbb{R}^3$, where the 2D covariance matrix is decomposed via Cholesky decomposition as $\Sigma_{2d} = L \cdot L^T$, requiring storage of only three elements from the lower triangular matrix L , denoted as $L = (l_{11}, l_{21}, l_{22})^T \in \mathbb{R}^3$, opacity $\alpha \in \mathbb{R}$, and RGB color representation $\mathbf{c} \in \mathbb{R}^3$.

When supervised with only a single image and its corresponding monocular depth estimation, random initialization of Gaussians followed by optimization produces floating Gaussian components in mid-air. These components disrupt the correct spatial structure of the scene and lead to poor fitting performance. Therefore, we adopt a back-projection initialization method based on the input image and depth estimation. The mean of each initial Gaussian component in NDC is determined using the normalized pixel coordinates and depth estimation of the corresponding pixel. Color values are initialized directly from the pixel RGB values. The initial axis lengths of the ellipse corresponding to the covariance matrix are computed based on the scene image resolution, focal length, and pixel depth values by analyzing the distribution range of initial Gaussian components in NDC space. This back-projection initialization approach enables rapid fitting of image Gaussians while preventing the emergence of floating Gaussian components. The supervision during optimization combines multiple loss terms: L1 and SSIM losses between the rendered and source scene images, regularization terms for occupancy and covariance, and L1 loss with respect to the depth estimation values:

$$\mathcal{L}_{gs} = (1 - \lambda_{ssim})\mathcal{L}_1 + \lambda_{ssim}\mathcal{L}_{ssim} + \mathcal{L}_{1_depth} + \lambda_{reg-\alpha}\mathcal{L}_{reg-\alpha} + \lambda_{reg-\Sigma}\mathcal{L}_{reg-\Sigma}, \quad (4)$$

we set λ_{ssim} to 0.2, $\lambda_{reg-\alpha}$ and $\lambda_{reg-\Sigma}$ to 0.01. After fitting Gaussians to the scene and object images, the object can be placed at any spatial location within the scene Gaussians at arbitrary scale by scaling and translating the object Gaussians. This process can be interactively accomplished by users in real-time within the visualization interface. Details are provided in Appendix D.

3.3 DIFFUSION-BASED OBJECT REFINEMENT

After achieving depth-aware object insertion through Gaussian combination of scene and object, the inserted object requires further refinement to harmonize its style and surface illumination with the surrounding environment. We propose a simple yet effective object refinement approach that leverages a pretrained diffusion model to harmonize the rendered image. Diffusion models pretrained on large-scale datasets inherently capture diverse output distributions encompassing various visual styles and lighting conditions. For cross-domain object insertion tasks (real object insert to stylized scene), we directly exploit the knowledge embedded in the pretrained FLUX.1-Fill-dev model to refine the inserted object, thereby achieving style harmonization between the object and scene.

We denote the combined Gaussians of scene and object as \mathcal{G}_{com} , where R represents the differentiable rendering process. Based on the combined Gaussians, we can simultaneously render the composite image and the opacity of object Gaussians to obtain the object mask: $\mathbf{I}_{com}, \mathbf{M}_{obj} = R(\mathcal{G}_{com})$. We introduce an coefficient $s \in [0, 1]$ to control the refinement strength. When $s = 0$, the method outputs the original image without modification, while $s = 1$ performs complete inpainting of the masked region. Let T denote the total number of inference steps in the diffusion model, corresponding to T different noise levels. Given refinement strength s , the corresponding noise level is computed as $T' = \lceil sT \rceil$. Let \mathcal{E} denote the encoder of the diffusion model. The latent variable corresponding to the composite image is $\mathbf{z}_{com} = \mathcal{E}(\mathbf{I}_{com})$. The initial latent variable for the inference process is then defined as:

$$\mathbf{z}_{T'} = \sigma(T')\epsilon + (1 - \sigma(T'))\mathbf{z}_{com}, \epsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (5)$$

DiT models such as FLUX process inputs by dividing them into patches and converting the values within each patch into tokens. Let the size of each patch be $k = h \times w$, with a total of N patches. We apply similar processing to \mathbf{M}_{obj} , transforming it into a patch-wise binary mask $\mathbf{M}_{patch} \in \{0, 1\}^N$:

$$m_n = \left[\sum_{i=1}^k \mathbf{M}_{obj.n,i} > 0 \right] \in \{0, 1\}. \quad (6)$$

We denote the diffusion model used for refinement as G_{refine} , which follows the flow matching framework. At each time step t during the refinement process, we first perform a single denoising step on the noisy latent variables.

$$\mathbf{z}_{t-1} = \mathbf{z}_t - \Delta t \cdot G_{refine}(\mathbf{z}_t, t|c), \quad (7)$$

where c represents the textual condition for the model. In our implementation, we employ a simple template format "XXX style of a XXX" that provides basic descriptions of object categories and styles. Subsequently, we perform an overwrite operation on the latent variables based on the patch-wise mask.

$$\mathbf{z}_{t-1} = \mathbf{M}_{patch} \odot \mathbf{z}_{t-1} + (1 - \mathbf{M}_{patch}) \odot \mathbf{z}_{ref}(t-1), \quad (8)$$

with

$$\mathbf{z}_{ref}(t) = \begin{cases} \sigma(t)\epsilon + (1 - \sigma(t))\mathbf{z}_{com} & t > 0, \\ \mathbf{z}_{com} & t = 0. \end{cases} \quad (9)$$

This ensures precise reconstruction of the original scene image in regions outside the object mask.

In real-scene object insertion, object refinement focuses primarily on achieving consistency between object surface lighting and the surrounding environment. Since widely-used open-source text-to-image diffusion models and inpainting models such as Stable Diffusion or FLUX possess limited knowledge in lighting harmonization, we additionally leverage the pre-trained illumination harmonization model LBM (Chadebec et al., 2025), denoted as G_{light} . For object refinement in real-scene, we first apply low-strength refinement to harmonize the inserted object, then perform additional lighting harmonization on the refined image using G_{light} . The refinement process for real-scene domain object insertion can be expressed as:

$$\mathbf{I}_{harmony} = G_{light}(\text{Refinement}(\mathbf{I}_{com}, s)). \quad (10)$$

After completing the refinement process on the rendered image following object insertion, the harmonized image from the diffusion model can directly provide supervision for Gaussian parameter updates due to the differentiable nature of the rendering process. However, during the refinement process, diffusion models often introduce unwanted modifications or detail loss in non-target regions due to constraints in their generative capabilities and supported resolutions. Therefore, we set only the object Gaussians as learnable to avoid affecting the scene Gaussians.

4 EXPERIMENTS

4.1 EXPERIMENT SETUP

Data Preparation Due to the lack of high-quality object insertion datasets with complex spatial relationships and cross-domain scenarios, we construct a high-resolution dataset containing diverse real photographs and artistic scenes with complex spatial structures. For scenes, we select high-quality photographs from the DIODE dataset (Vasiljevic et al., 2019), high-resolution artworks from WikiArt (Saleh & Elgammal, 2024), frames from anime and cartoon videos, and generate additional scenes using GPT-Image-1 (OpenAI, 2024) and FLUX-pro-1.1-ultra. Stylized scenes encompass anime, cartoon, comic, pixel art, watercolor, and over ten oil painting styles from WikiArt. Real-world scenes include diverse indoor/outdoor environments, various lighting conditions (strong, weak, nighttime), and different light sources (natural, artificial, colored). This yields 225 high-quality scenes across multiple domains. For objects, we collect web resources and use SAM (Kirillov et al., 2023) to segment objects from COCO (Lin et al., 2014), obtaining 255 objects spanning over 9 categories including animals, vegetables, food, people, and vehicles. We provide over 200 semantically coherent scene-object combinations for evaluation. All experiments are conducted on this dataset. We also report results on the TF-ICON benchmark in Sec 4.5.

Experiment Details For Gaussian fitting, we employ the Adam optimizer with learning rates of 0.01 for means, $2.5e - 3$ for opacities, 0.01 for RGB values, and 0.1 for covariances. We fix the global random seed to 42. Since our method involves only forward inference of the refinement model and Gaussian fitting, a single A6000 GPU is sufficient to support object refinement using FLUX.1-Fill-dev.

Baselines To validate our method’s effectiveness, we compare against both training-based and training-free object insertion approaches. The training-based baselines include InsertAnything, Uni-Combine, and AnyDoor, while the training-free baselines comprise Freecompose, Primecomposer.



351 Figure 2: Visual comparison results between our SpatialComposer and baselines. Zoom in to ob-
352 serve details.

355 4.2 QUANTITATIVE EVALUATION

357 Existing metrics measure non-edited region preservation by calculating distance between low-level
358 features in generated and original scenes. We compute PSNR between these regions and the original
359 scene. For object identity, prior works (Lu et al., 2023b; Wang et al., 2024; Chen et al., 2024c) mea-
360 sure distance between high-level semantic or low-level perceptual features from edited regions and
361 reference objects. Our method considers depth-related positioning during insertion. After handling
362 occlusion relationships, feature distance between edited regions and references actually increases.
363 When domain gaps exist between scenes and objects, color and texture modifications are required for
364 seamless integration, degrading quantitative metrics. We therefore exclude these measures. Follow
365 prior works, we adopt the cosine similarity between the CLIP image embedding of the edited region
366 in the result image and the CLIP text embedding of the corresponding text prompt for the edited
367 region (formatted as "XXX style / professional photograph of a XXX") to com-
368 prehensively measure both inserted object style and category. CLIP-based similarity metrics do
369 not account for alignment with human visual perception, often yielding results inconsistent with
370 human subjective evaluation. Therefore, we introduce VQAScore (Lin et al., 2024), a metric specifi-
371 cally designed for generated image assessment that demonstrates superior alignment with human
372 subjective judgment, to evaluate both the overall visual quality of inserted objects and their sur-
373 rounding regions, as well as their alignment with text prompts (also formatted as "XXX style /
374 professional photograph of a XXX"). However, spatial positioning and scale reason-
375 ableness remains difficult to assess automatically.

376 Given the lack of automatic metrics aligned with human perception for interactive generation tasks,
377 this field relies on user studies for evaluation. We recruited 40 participants to evaluate SpatialCom-
378 poser and baselines across four dimensions: overall quality, object-environment harmony, spatial
379 positioning and scale reasonableness, and consistency with reference objects and original scenes.

Table 1: Quantitative comparison of different methods on different benchmarks. **Bold** indicates the best result, underline denotes the second best, and [†] marks training-free methods.

Method	Our DAOI Dataset			TF-ICON Dataset		
	PSNR _{bg}	CLIPScore	VQAScore	PSNR _{bg}	CLIPScore	VQAScore
AnyDoor	31.52	0.287	0.520	23.74	0.279	0.493
UniCombine	32.03	0.255	0.447	26.61	<u>0.286</u>	0.567
InsertAnything	<u>36.82</u>	<u>0.296</u>	<u>0.535</u>	<u>31.91</u>	0.285	0.557
Freecompose [†]	23.56	0.281	0.529	21.33	0.283	0.546
Primecomposer [†]	16.82	0.268	0.445	12.32	0.272	<u>0.590</u>
SpatialComposer [†]	46.99	0.299	0.552	33.14	0.307	0.653

Table 2: User study on our dataset across different evaluation dimensions.

	AnyDoor	UniComb	InsertAnything	Freecomp [†]	Primecomp [†]	SpatialComp [†]
Qual	2.94	4.96	<u>8.92</u>	1.75	0.55	80.88
Harm	2.11	6.71	<u>7.17</u>	1.84	0.74	81.43
Reas	2.30	4.32	<u>7.08</u>	1.38	0.74	84.19
Cons	2.21	5.61	<u>9.38</u>	1.75	0.46	80.61

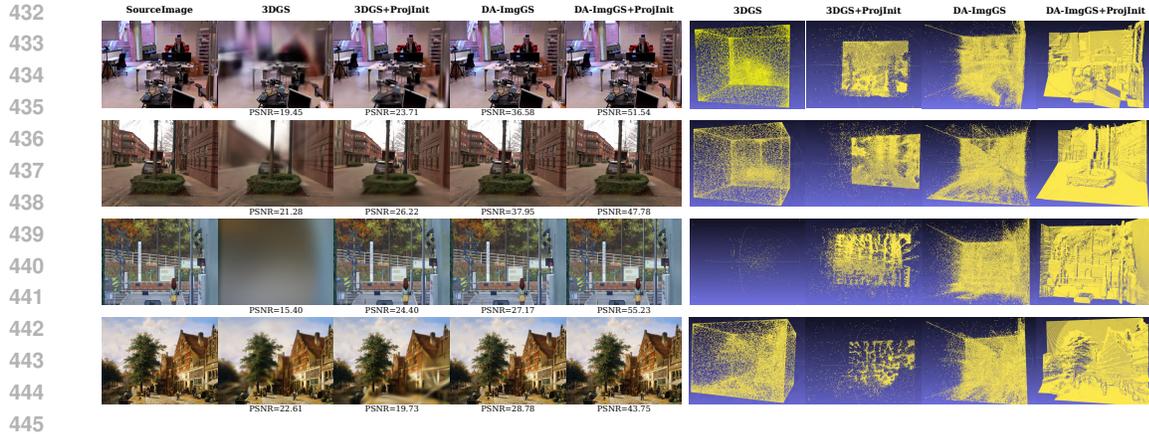
Participants identified the best method per dimension, and we computed average vote percentages across all dimensions. Table. 2 shows SpatialComposer achieves comparable or superior performance in automatic metrics. Our dataset comprises high-resolution real photographs and artworks with complex spatial structures and diverse styles at insertion locations. These characteristics challenge existing methods. SpatialComposer overcomes these through depth-aware Gaussian representations and effective use of pre-trained diffusion models, achieving significantly superior user study performance.

4.3 QUALITATIVE EVALUATION

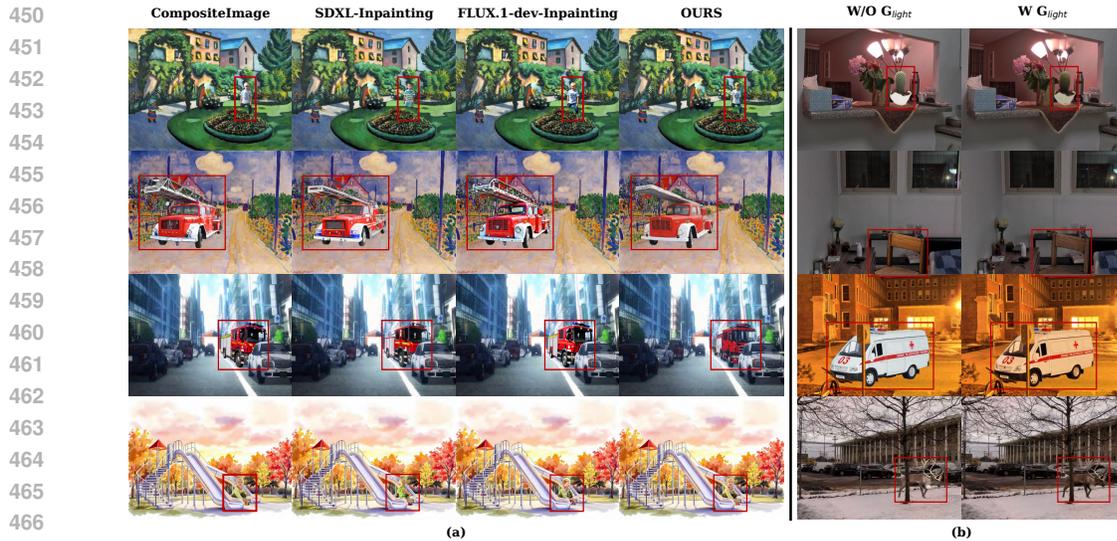
As shown in Fig. 2, InsertAnything demonstrates the best performance among the baselines, showing some understanding of object depth relationships and displaying reasonable depth relations in certain insertion results. However, its capabilities for style harmonization in cross-domain insertion and illumination processing in real scenes remain limited. UniCombine exhibits stronger style harmonization capabilities for cross-domain object insertion but frequently generates results with significant errors or unreasonable outcomes. AnyDoor shows limitations in object-scene harmony and the visual quality in edited regions. FreeCompose lacks understanding of depth information and demonstrates limited harmonization performance. Beyond the issues encountered by other methods, PrimeComposer also exhibits significant problems in preserving the original scene image. Overall, SpatialComposer achieves comparable or superior performance in terms of image quality, reasonableness of spatial positioning and scale, consistency of object identity with non-edited regions of the scene, and harmonization between the scene and objects. More results are provided in Appendix B. Due to file size constraints, we provide vector-format versions of the visualization comparison figures in the supplementary material.

4.4 ABLATION STUDY

Gaussian representation and initialization method We simplify standard 3D Gaussians to better adapt them to our task and employ a back-projection-based initialization method. This approach not only significantly enhances scene fitting quality but also achieves superior spatial structure representation. As shown in Fig. 3, our proposed Gaussian representation and initialization method substantially improve the fitting quality of scene reconstruction. We also visualize the Gaussian means under different configurations, demonstrating that our proposed Gaussian settings and initialization method yield the most reasonable spatial structure, which serves as the foundation for subsequent object insertion operations.



446 Figure 3: Comparison of fitting performance under different Gaussian representations and initializa-
 447 tion strategies. Both the proposed depth-aware image Gaussian and initialization strategy demon-
 448 strate improved fitting quality and spatial structure.
 449



468 Figure 4: (a) demonstrates the compatibility of our pipeline with different pre-trained diffusion
 469 models and the effectiveness of our refinement method; (b) shows that incorporating pre-trained
 470 illumination harmonization models in real-scene object insertion enhances the harmony between
 471 object surface lighting and the environment.
 472

473

474 **Refinement Foundation Model** In the refinement processing of inserted objects, our pipeline is
 475 compatible with different foundation models. Fig. 4 (a) compares the refinement results across
 476 different foundation models. Our proposed refinement method achieves the best overall performance
 477 in maintaining object identity while achieving style consistency between objects and scenes.

478 **Illumination Harmonization Model in Real Scenes** In Fig. 4 (b), we validate the limitations of pre-
 479 trained text-to-image and inpainting models when addressing object surface illumination harmoniza-
 480 tion in real scenes, as well as the necessity of incorporating pretrained illumination harmonization
 481 models. The illumination harmonization model better captures environmental lighting information,
 482 including intensity, color, and direction.

483 **Refinement Strength** Finally, we validate the impact of the strength coefficient s on the results
 484 during the refinement process. As shown in Fig. 5, higher refinement strength produces results that
 485 are more harmonious with the scene style but simultaneously weakens the preservation of object
 identity.



504 Figure 5: The impact of different object refinement intensities on the results.

505

506

507 the inserted object and the scene. In our experiments, we select $s = 0.6$ as the default setting for

508 stylized scene object insertion and $s = 0.2$ for real scene object insertion.

509 4.5 RESULTS ON TF-ICON BENCHMARK

510

511

512

513 Prior works most closely related to ours, TF-ICON and Primecomposer, primarily employ the TF-

514 ICON Benchmark. This dataset comprises 95 stylized scene object insertion cases and 237 real

515 scene object insertion cases, constructed from approximately 30 scene images and about 100 object

516 images. Due to its low scene resolution, synthetic scene generation, and simplistic object insertion

517 scenarios, this dataset cannot fully satisfy the requirements of our task. Although the object in-

518 scription locations in this dataset do not involve complex spatial relationships, SpatialComposer still

519 achieves performance comparable to or superior to existing methods in terms of overall quality. The

520 quantitative results presented in Table.1 further support our observations. We also present a visual

521 comparison of different methods on the TF-ICON benchmark in Appendix Fig.6.

522 5 CONCLUSION AND LIMITATION

523

524

525 In this paper, we introduce SpatialComposer, which effectively reconstructs high-fidelity scene rep-

526 resentations with meaningful depth structure, enabling precise control over object scale and 3D spa-

527 tial positioning during insertion. To address insertion disharmony, we employ a refinement method

528 based on pretrained diffusion models, achieving reasonable harmonization that aligns object and

529 lighting with the surrounding scene. Furthermore, we constructed the DAOI dataset, in which we

530 collected over 200 high-resolution scene images with complex spatial structures and diverse styles,

531 along with over 200 multi-category object images. Experiments demonstrate our method’s superi-

532 ority, achieving comparable or superior performance across multiple evaluation dimensions.

533

534 In our task setting, given only a single image of the insertion object, we cannot construct a complete

535 Gaussian representation of the object and can only model the visible surface shown in the single im-

536 age. Consequently, in the subsequent Gaussian composition operations, our method only supports

537 scaling and translation of the object Gaussians, but does not support rotation of the object Gaus-

538 sians. Reconstructing individual object Gaussians from single object images represents an actively

539 pursued research direction in the 3D Gaussian field. With the rapid advancement of techniques in

539 this domain, incorporating object 3D Gaussian reconstruction models based on a single object image

539 may provide an effective approach to address the limitations of our method.

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. Our research focuses on image editing tasks, where all data are sourced from open-source datasets, generated through model API calls, and publicly available internet data, with no involvement of personal privacy information. All data used in this study are publicly available and were used in accordance with their respective licenses and terms of use. We have properly cited all data sources and respected the original creators' rights. While our method demonstrates improvements in object insertion, we recognize that like other AI technologies, it could potentially be misused to generate unsafe visual content. We encourage responsible deployment and further research into safety mechanisms. We believe this work contributes positively to the research community and poses minimal ethical concerns when used responsibly.

REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our work. Our constructed dataset DAOI will be publicly available alongside our framework and experimental source code upon publication, while TF-ICON serves as a public benchmark. Our proposed method is detailed in Sec. 3, which includes the hyperparameters involved in the loss terms. Additional experimental details such as optimizers, learning rates, and random seeds are provided in Sec. 4. We commit to making the complete framework and experimental source code, as well as the constructed benchmark, publicly available as soon as possible after publication.

REFERENCES

- Gemma Canet Tarrés, Zhe Lin, Zhifei Zhang, Jianming Zhang, Yizhi Song, Dan Ruta, Andrew Gilbert, John Collomosse, and Soo Ye Kim. Thinking outside the bbox: Unconstrained generative object compositing. In *European Conference on Computer Vision*, pp. 476–495. Springer, 2024.
- Clément Chadebec, Onur Tasar, Sanjeev Sreetharan, and Benjamin Aubin. Lbm: Latent bridge matching for fast image-to-image translation. *arXiv preprint arXiv:2503.07535*, 2025.
- Jiaxuan Chen, Bo Zhang, Qingdong He, Jinlong Peng, and Li Niu. Mureobjectstitch: Multi-reference image composition. *arXiv preprint arXiv:2411.07462*, 2024a.
- Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor: Zero-shot object-level image customization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6593–6602, 2024b.
- Zhekai Chen, Wen Wang, Zhen Yang, Zeqing Yuan, Hao Chen, and Chunhua Shen. Freecompose: Generic zero-shot image composition with diffusion prior. In *European Conference on Computer Vision*, pp. 70–87. Springer, 2024c.
- Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixing Ling, Weiyuan Li, and Liqing Zhang. Dovenet: Deep image harmonization via domain verification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8394–8403, 2020.
- Wenyan Cong, Xinhao Tao, Li Niu, Jing Liang, Xuesong Gao, Qihao Sun, and Liqing Zhang. High-resolution image harmonization via collaborative dual transformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18470–18479, 2022.
- Jixuan He, Wanhua Li, Ye Liu, Junsik Kim, Donglai Wei, and Hanspeter Pfister. Affordance-aware object insertion via mask-aware dual diffusion. *arXiv preprint arXiv:2412.14462*, 2024.
- Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2328–2337, 2023.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.

- 594 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
595 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 596
- 597 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D Gaussian splat-
598 ting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42(4):1–14,
599 2023.
- 600 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
601 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-*
602 *ings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
- 603
- 604 Sumith Kulal, Tim Brooks, Alex Aiken, Jiajun Wu, Jimei Yang, Jingwan Lu, Alexei A Efros, and Kr-
605 ishna Kumar Singh. Putting people in their place: Affordance-aware human insertion into scenes.
606 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
607 17089–17099, 2023.
- 608 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 609
- 610 Pengzhi Li, Qiang Nie, Ying Chen, Xi Jiang, Kai Wu, Yuhuan Lin, Yong Liu, Jinlong Peng, Chengjie
611 Wang, and Feng Zheng. Tuning-free image customization with image and text guidance. In
612 *European Conference on Computer Vision*, pp. 233–250. Springer, 2024.
- 613 Tianle Li, Max Ku, Cong Wei, and Wenhui Chen. Dreamedit: Subject-driven image editing. *arXiv*
614 *preprint arXiv:2306.12624*, 2023.
- 615
- 616 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
617 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
618 *conference on computer vision*, pp. 740–755. Springer, 2014.
- 619 Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and
620 Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European*
621 *Conference on Computer Vision*, pp. 366–384. Springer, 2024.
- 622
- 623 Jun Ling, Han Xue, Li Song, Rong Xie, and Xiao Gu. Region-aware adaptive instance normalization
624 for image harmonization. In *Proceedings of the IEEE/CVF conference on computer vision and*
625 *pattern recognition*, pp. 9361–9370, 2021.
- 626 Liu Liu, Zhenchen Liu, Bo Zhang, Jiangtong Li, Li Niu, Qingyang Liu, and Liqing Zhang. Opa:
627 object placement assessment dataset. *arXiv preprint arXiv:2107.01889*, 2021.
- 628
- 629 Lingxiao Lu, Jiangtong Li, Bo Zhang, and Li Niu. Dreamcom: Finetuning text-guided inpainting
630 model for image composition. *arXiv preprint arXiv:2309.15508*, 2023a.
- 631 Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. Tf-icon: Diffusion-based training-free cross-
632 domain image composition. In *Proceedings of the IEEE/CVF International Conference on Com-*
633 *puter Vision*, pp. 2294–2305, 2023b.
- 634
- 635 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for
636 editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF conference*
637 *on computer vision and pattern recognition*, pp. 6038–6047, 2023.
- 638 Li Niu, Qingyang Liu, Zhenchen Liu, and Jiangtong Li. Fast object placement assessment. *arXiv*
639 *preprint arXiv:2205.14280*, 2022.
- 640
- 641 OpenAI. gpt-image-1. <https://platform.openai.com/docs/models#gpt-image>,
642 2024. Accessed: 2025-09-24.
- 643 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
644 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
645 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022a.
- 646
- 647 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
resolution image synthesis with latent diffusion models. *arXiv preprint arXiv:2112.10752*, 2022b.

- 648 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
649 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*
650 *ceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22500–
651 22510, 2023.
- 652 Babak Saleh and Ahmed Elgammal. Wikiart: A large-scale dataset of artworks,
653 dec 2024. URL [https://service.tib.eu/ldmservice/dataset/](https://service.tib.eu/ldmservice/dataset/wikiart--a-large-scale-dataset-of-artworks)
654 [wikiart--a-large-scale-dataset-of-artworks](https://service.tib.eu/ldmservice/dataset/wikiart--a-large-scale-dataset-of-artworks).
655
- 656 Konstantin Sofiiuk, Polina Popenova, and Anton Konushin. Foreground-aware semantic represen-
657 tations for image harmonization. In *Proceedings of the IEEE/CVF winter conference on applica-*
658 *tions of computer vision*, pp. 1620–1629, 2021.
- 659 Wensong Song, Hong Jiang, Zongxing Yang, Ruijie Quan, and Yi Yang. Insert anything: Image
660 insertion via in-context editing in dit. *arXiv preprint arXiv:2504.15009*, 2025.
661
- 662 Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim, and
663 Daniel Aliaga. Objectstitch: Generative object compositing. *arXiv preprint arXiv:2212.00932*,
664 2022.
- 665 Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim,
666 He Zhang, Wei Xiong, and Daniel Aliaga. Imprint: Generative object compositing by learning
667 identity-preserving representation. In *Proceedings of the IEEE/CVF Conference on Computer*
668 *Vision and Pattern Recognition*, pp. 8048–8058, 2024.
- 669 Shashank Tripathi, Siddhartha Chandra, Amit Agrawal, Amrith Tyagi, James M Rehg, and Visesh
670 Chari. Learning to generate synthetic data via compositing. In *Proceedings of the IEEE/CVF*
671 *Conference on Computer Vision and Pattern Recognition*, pp. 461–470, 2019.
672
- 673 Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z. Dai, An-
674 drea F. Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R. Walter, and Gregory
675 Shakhnarovich. DIODE: A Dense Indoor and Outdoor DEpth Dataset. *CoRR*, abs/1908.00463,
676 2019. URL <http://arxiv.org/abs/1908.00463>.
- 677 Haoxuan Wang, Jinlong Peng, Qingdong He, Hao Yang, Ying Jin, Jiafu Wu, Xiaobin Hu, Yanjie
678 Pan, Zhenye Gan, Mingmin Chi, et al. Unicombine: Unified multi-conditional combination with
679 diffusion transformer. *arXiv preprint arXiv:2503.09277*, 2025.
- 680 Yibin Wang, Weizhong Zhang, Jianwei Zheng, and Cheng Jin. Primecomposer: Faster progressively
681 combined diffusion for image composition with attention steering. In *Proceedings of the 32nd*
682 *ACM International Conference on Multimedia*, pp. 10824–10832, 2024.
- 683 Daniel Winter, Matan Cohen, Shlomi Fruchter, Yael Pritch, Alex Rav-Acha, and Yedid Hoshen.
684 Objectdrop: Bootstrapping counterfactuals for photorealistic object removal and insertion. In
685 *European Conference on Computer Vision*, pp. 112–129. Springer, 2024a.
- 686 Daniel Winter, Asaf Shul, Matan Cohen, Dana Berman, Yael Pritch, Alex Rav-Acha, and Yedid
687 Hoshen. Objectmate: A recurrence prior for object insertion and subject-driven generation. *arXiv*
688 *preprint arXiv:2412.08645*, 2024b.
- 689 Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and
690 Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *Proce-*
691 *edings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18381–18391,
692 2023.
- 693 Yongsheng Yu, Ziyun Zeng, Haitian Zheng, and Jiebo Luo. Omnipaint: Mastering object-oriented
694 editing via disentangled insertion-removal inpainting. *arXiv preprint arXiv:2503.08677*, 2025.
- 695 Ziyang Yuan, Mingdeng Cao, Xintao Wang, Zhongang Qi, Chun Yuan, and Ying Shan. Customnet:
696 Object customization with variable-viewpoints in text-to-image diffusion models. In *Proceed-*
697 *ings of the 32nd ACM International Conference on Multimedia*, pp. 10976–10984, 2024.
- 698 Bo Zhang, Yuxuan Duan, Jun Lan, Yan Hong, Huijia Zhu, Weiqiang Wang, and Li Niu. Controlcom:
699 Controllable image composition using diffusion model. *arXiv preprint arXiv:2308.10040*, 2023a.
700
701

Lingzhi Zhang, Tarmily Wen, Jie Min, Jiancong Wang, David Han, and Jianbo Shi. Learning object placement by inpainting for compositional data augmentation. In *European Conference on Computer Vision*, pp. 566–581. Springer, 2020.

Xin Zhang, Jiaxian Guo, Paul Yoo, Yutaka Matsuo, and Yusuke Iwasawa. Paste, inpaint and harmonize via denoising: Subject-driven image editing with pre-trained diffusion model. *arXiv preprint arXiv:2306.07596*, 2023b.

Zitian Zhang, Frédéric Fortier-Chouinard, Mathieu Garon, Anand Bhattad, and Jean-François Lalonde. Zerocomp: Zero-shot object compositing from image intrinsics via diffusion. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 483–494. IEEE, 2025.

A LATENT DIFFUSION MODEL

In this work, we leverage pre-trained text-conditioned latent diffusion models (Rombach et al., 2022a) to guide the object refinement process. These models operate in a learned latent space through an encoder-decoder architecture, where $\mathcal{E}(\cdot)$ and $\mathcal{D}(\cdot)$ represent the encoder and decoder components, respectively. The diffusion process involves forward noise addition followed by reverse denoising operations within this latent representation. Given an input image \mathbf{x}_0 , we first encode it into the latent space as $\mathbf{z}_0 = \mathcal{E}(\mathbf{x}_0)$. During training, this latent representation is progressively corrupted through the forward diffusion process, transforming \mathbf{z}_0 into \mathbf{z}_t :

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), \quad (11)$$

for $t \in [1, T]$, where $\bar{\alpha}_t = \prod_{s=1}^t 1 - \beta_s$, and β_s represents the variance schedule at timestep s . Subsequently, a denoising U-Net is trained to predict the added noise conditioned on c using the following objective function:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), t} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t, c)\|_2^2], \quad (12)$$

where $\boldsymbol{\epsilon}_\theta$ represents the denoising U-Net.

B ADDITIONAL VISUAL COMPARISON RESULTS

In Figs. 7 and 8, we present additional visualized comparisons of generation results.

C MULTI-OBJECT INSERTION EXAMPLES

SpatialComposer supports inserting multiple objects in a single scene. Since the Gaussians of these objects and the scene Gaussians are distinguishable, it avoids the accumulation of influence on other Gaussians when inserting multiple objects. In Fig. 9, we show some results of multiple object insertion in scenes.

D INTERACTIVE GAUSSIAN COMPOSITION BASED ON VISUALIZATION

In the scene-object composition phase, we leverage the point clouds of Gaussian means from both the scene and object, which are preserved during the Gaussian fitting process, and visualize these two point clouds in a unified coordinate system using the Open3D library. The system supports 360-degree rotation of the entire coordinate system via mouse interaction, and allows real-time adjustment of the object Gaussian’s scale, x-coordinate, y-coordinate, and z-coordinate through four pairs of keyboard keys, with continuous tracking and output of the scale and coordinate parameters during the adjustment process. Upon completion of the adjustment and closure of the visualization window, the scale and coordinate parameters are passed to the corresponding functions to perform operations on the object Gaussian, which is then composed with the scene Gaussian. As illustrated in Fig. 10, through this visualization-based Gaussian composition approach, we can precisely control the scale and position of the object Gaussian, enabling us to place the chair behind two pillars at different depths within the scene.

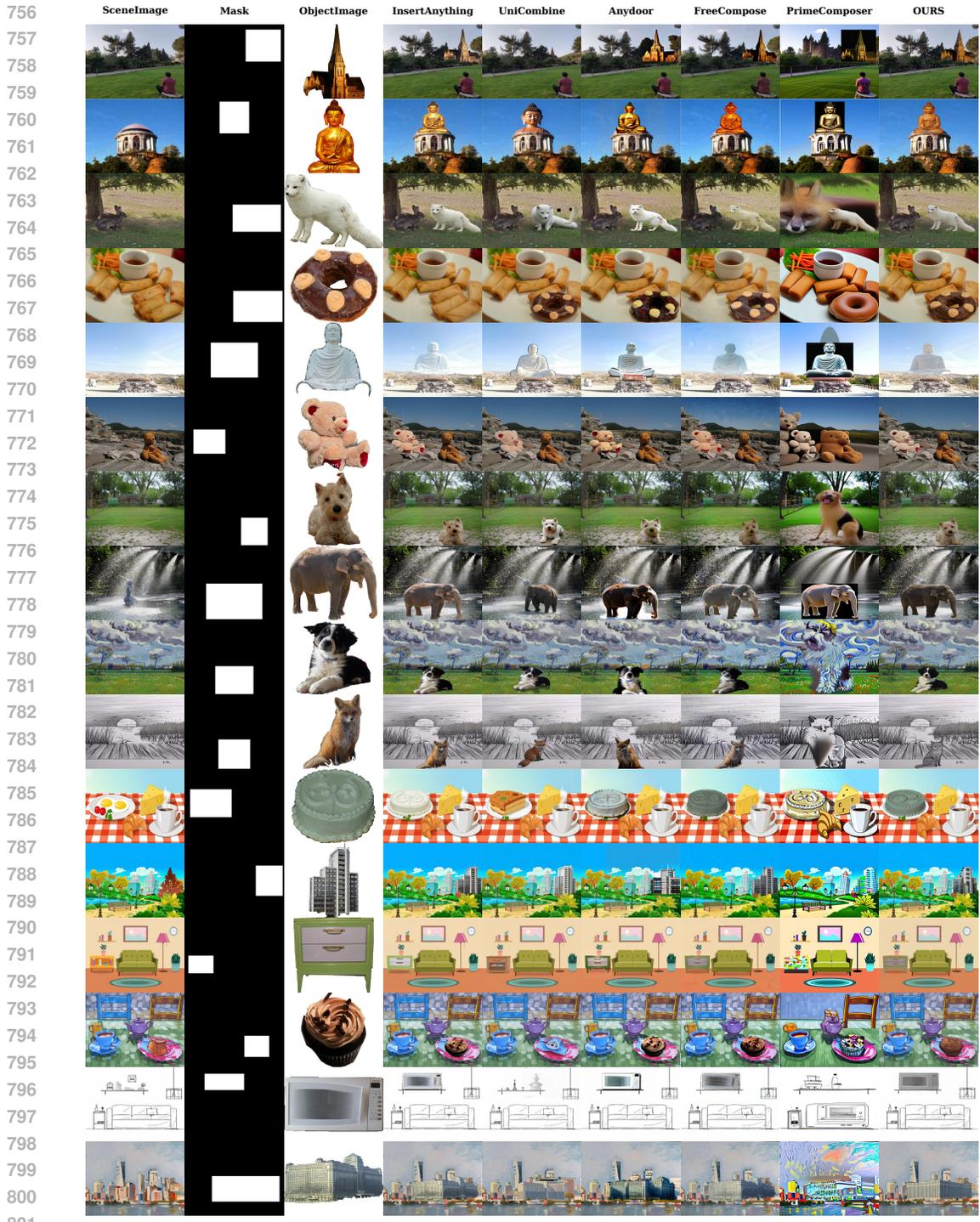
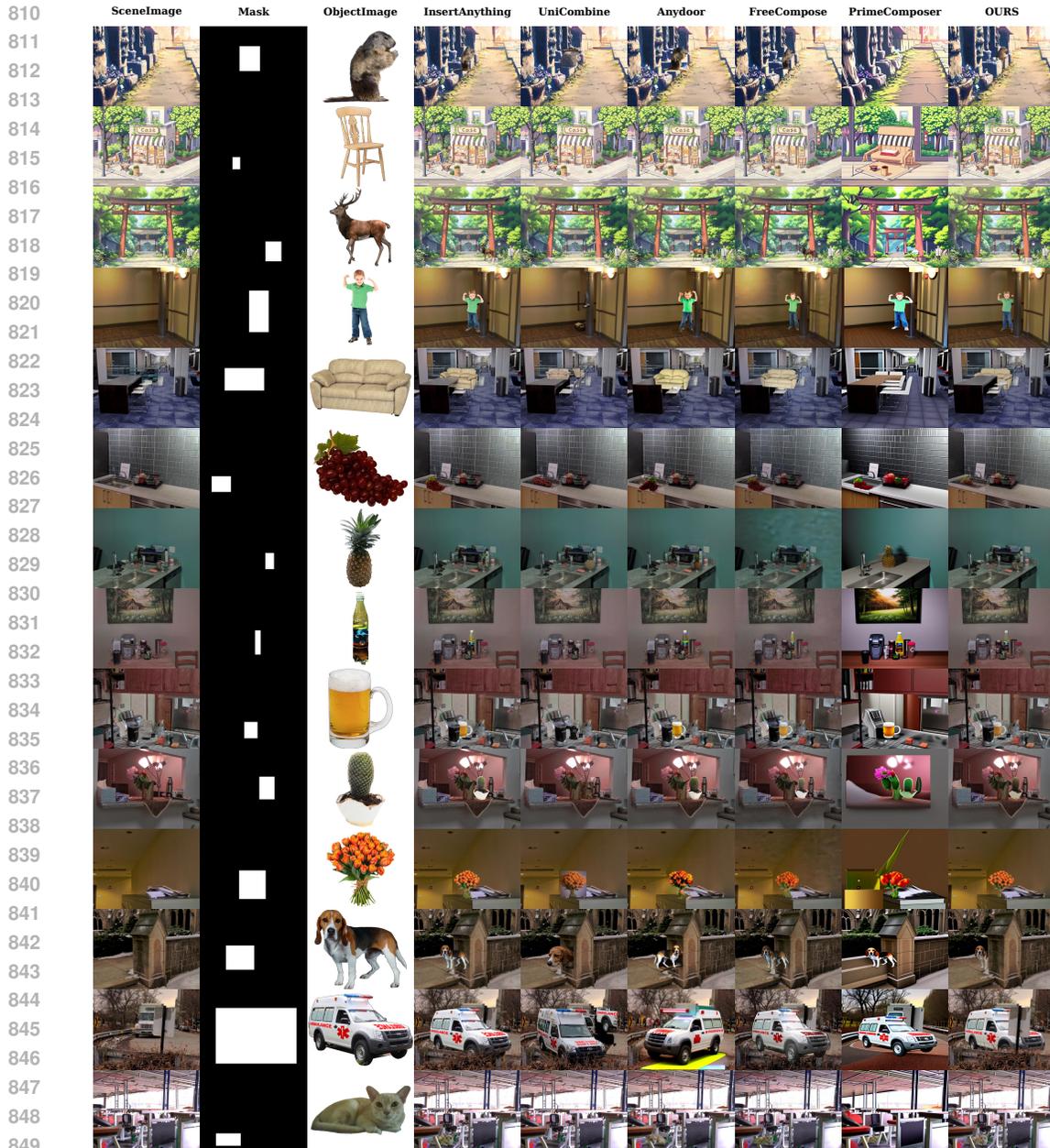


Figure 6: Visualization of comparative results on TF-ICON benchmark, zoomed in for detailed observation.

802
803
804
805
806
807
808
809



850
851 Figure 7: Visualization of comparative results, zoomed in for detailed observation.

852
853
854 **E THE USE OF LARGE LANGUAGE MODELS**

855
856 During the writing process of this paper, we utilized large language models to enhance the
857 manuscript quality, including employing large language models to correct grammatical errors and
858 modify wording and expressions to achieve a more formal and academic style.

859
860
861
862
863

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



Figure 8: Visualization of comparative results, zoomed in for detailed observation.

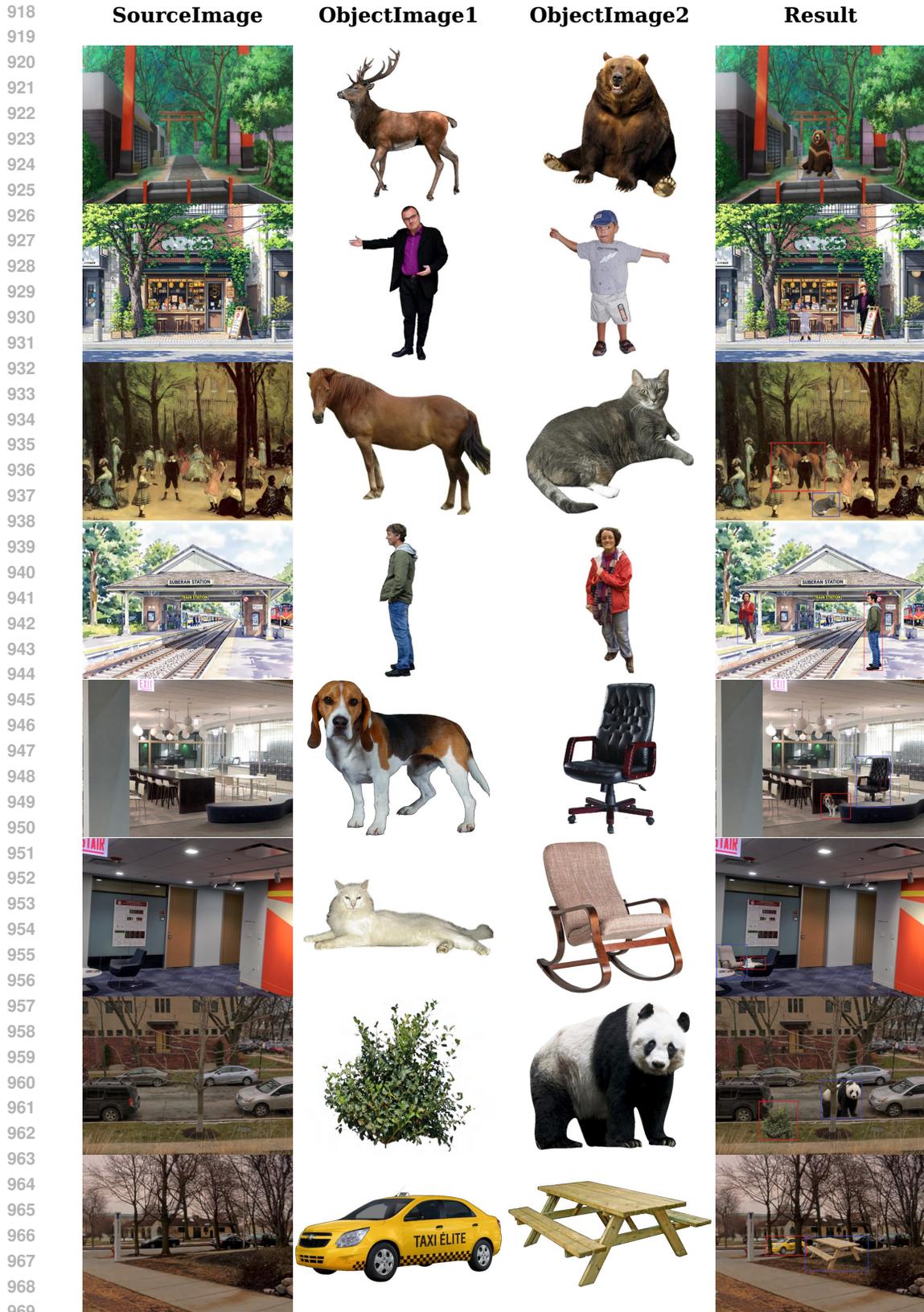


Figure 9: SpatialComposer supports multi-object insertion within a single scene while avoiding the adverse effects of multiple object insertion operations on other regions of the image, zoomed in for detailed observation.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



Figure 10: The first image on the left shows the scene, the second depicts the object to be inserted into the scene, and the third and fourth images demonstrate that through the visualization-based scene-object Gaussian composition process, we can precisely control the placement of the object Gaussian behind the first pillar and the second pillar in the scene, respectively.