

# 3DiScene: Editing Any Scene via Language-guided Disentangled Gaussian Splatting

QIHANG ZHANG, The Chinese University of Hong Kong, HKSAR  
YINGHAO XU, Stanford University, United States of America  
CHAOYANG WANG, Snap Inc., United States of America  
HSIN-YING LEE, Snap Inc., United States of America  
GORDON WETZSTEIN, Stanford University, United States of America  
BOLEI ZHOU, University of California Los Angeles, United States of America  
CEYUAN YANG, ByteDance, China

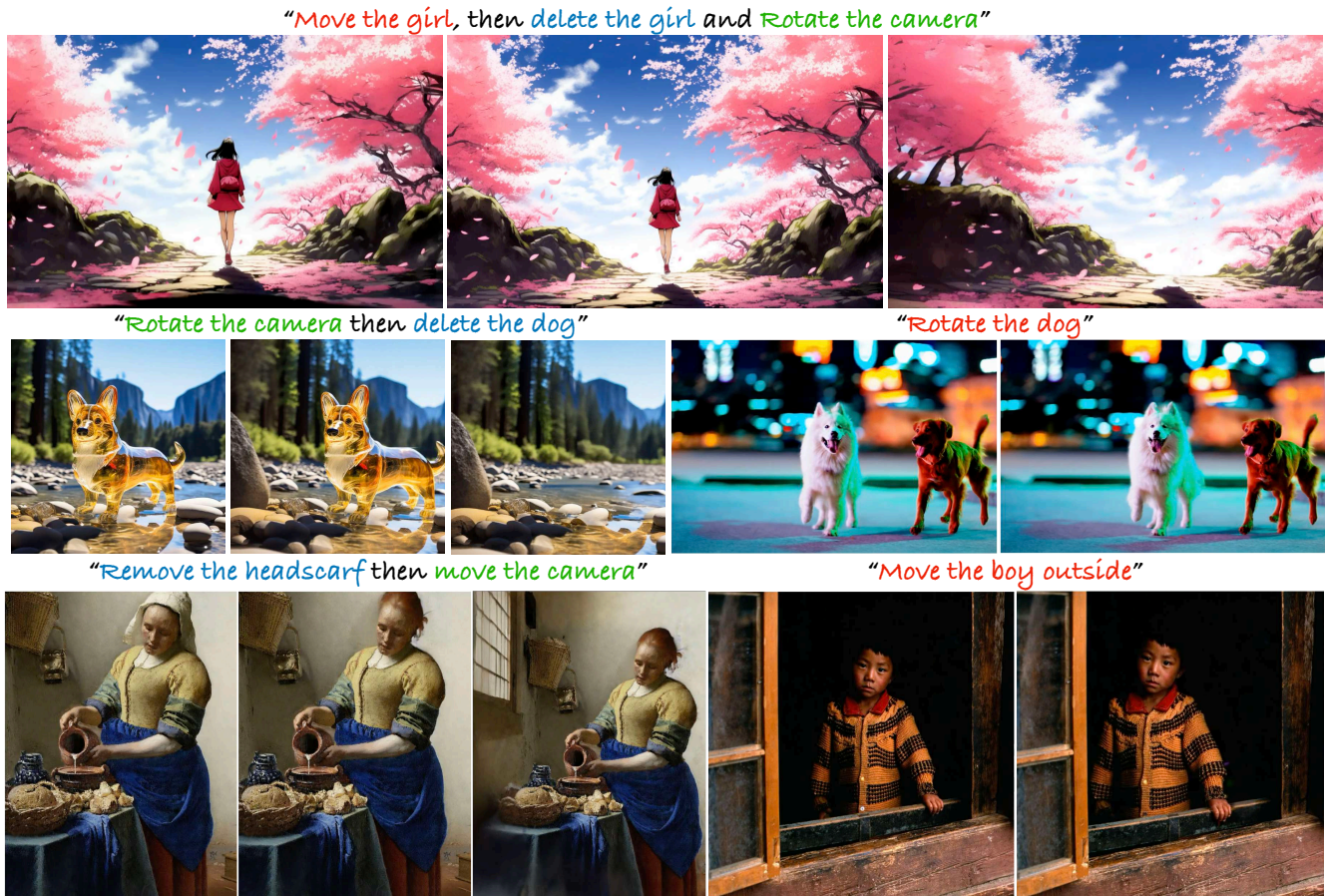


Fig. 1. Image pairs edited by 3DiScene. Our method has the capability to simultaneously handle diverse types of editing across both 2D and 3D dimensions.

Scene image editing is crucial for entertainment, photography, and advertising design. Existing methods solely focus on either 2D individual object or 3D global scene editing. This results in a lack of a unified approach to effectively control and manipulate scenes at the 3D level with different levels of granularity. In this work, we propose 3DiScene, a novel and unified scene editing framework leveraging language-guided disentangled Gaussian Splatting that enables seamless editing from 2D to 3D, allowing precise control over scene composition and individual objects. We first incorporate 3D Gaussians that are refined through generative priors and optimization

techniques. Language features from CLIP then introduce semantics into 3D geometry for object disentanglement. With the disentangled Gaussians, 3DiScene allows for manipulation at both the global and individual levels, revolutionizing creative expression and empowering control over scenes and objects. Experimental results demonstrate the effectiveness and versatility of 3DiScene in scene image editing. Code and online demo can be found at our project homepage: <https://zqh0253.github.io/3DiScene/>.

CCS Concepts: • Computing methodologies → Computer vision.

Additional Key Words and Phrases: Image editing, 3D scene generation

## 1 INTRODUCTION

Editing scene images is of great importance in various fields, ranging from entertainment, professional photography and advertising design. Content editing allows to create immersive and captivating experiences for audiences, convey the artistic vision effectively and achieve the desired aesthetic outcomes. With the rapid development of deep generative modeling, many attempts have been made to edit an image effectively. However, they have encountered limitations that hindered their potential.

Previous methods primarily concentrate on scene editing in 2D image space. They commonly rely on generative priors, such as GANs and Diffusion Models (DM), and employ techniques like modification of cross-attention mechanisms [Hertz et al. 2022, 2023], and optimization of network parameters [Chen et al. 2023a; Gal et al. 2022; Kawar et al. 2023; Kim et al. 2022; Ruiz et al. 2023] to edit the appearance and object identity within scene images. While some efforts have been made to extend these methods to 3D editing, they ignore 3D cues and pose a challenge in maintaining 3D consistency, especially when changing the camera pose. Moreover, these approaches typically focus on global scenes and lack the ability to disentangle objects accurately, resulting in limited control over individual objects at the 3D level.

In order to edit any scene images and enable 3D control over both scene and its individual objects, we propose 3Di tScene, a novel scene editing framework which leverage a new scene representation, language-guided disentangled Gaussian Splatting. Concretely, the given image is first projected into 3D Gaussians which are further refined and enriched through 2D generative prior [Poole et al. 2022; Rombach et al. 2022]. We thus obtain a comprehensive 3D scene representation that naturally enables novel view synthesis for a given image. In addition, language features from CLIP are distilled into the corresponding 3D Gaussians to introduce semantics into 3D geometry. These semantic 3D Gaussians help disentangle individual objects out of the entire scene representation, resulting in language-guided disentangled Gaussians for scene decomposition. They also allow for a more user-friendly interaction *i.e.*, users could query specific objects or interest via text. To this end, our 3Di tScene enables seamless editing from 2D to 3D and allow for modifications at both the global and individual levels, empowering creators to have precise control over scene composition, and object-level edits.

We dub our pipeline as 3Di tScene. Different from previous works that focus on addressing a single type of editing, 3Di tScene integrates editing requirements within a unified framework. Our teaser figure demonstrates the versatility of 3Di tScene by showcasing its application to diverse scene images. We have conducted evaluations of 3Di tScene under various settings, and the results demonstrate significant improvements over baseline methods.

## 2 RELATED WORK

**Image Editing with Generative Models.** The field of 2D image synthesis has advanced significantly with the development of generative models such as GANs [Karras et al. 2021, 2019] and diffusion models [Ho et al. 2020; Rombach et al. 2022; Song et al. 2020]. Many studies capitalize on the rich prior knowledge embedded in generative models for image editing. Some endeavors utilize

GANs for various image editing tasks, including image-to-image translation, latent manipulation [Jahani et al. 2019; Shen et al. 2020; Xu et al. 2021; Yang et al. 2021; Zhu et al. 2020], and text-guided manipulation [Patashnik et al. 2021]. However, due to limitations in training on large-scale data, GANs often struggle to perform well on real-world scene images. As diffusion models make notable progress, the community is increasingly focusing on harnessing the potent text-to-image diffusion model for real image editing [Chen et al. 2023a; Gal et al. 2022; Hertz et al. 2022, 2023; Kawar et al. 2023; Kim et al. 2022; Meng et al. 2021a; Ruiz et al. 2023; Su et al. 2022]. However, these methods are confined to the 2D domain and are limited in editing objects within a 3D space. Concurrently, other research efforts [Yenphraphai et al. 2024a] attempt to address 3D-aware image editing, but they introduces inconsistency in the editing process, and cannot change the camera perspective of the entire scene. In contrast, our model leverages an explicit 3D Gaussian to convert 2D images into 3D space while disentangling objects with language guidance. This approach enables our model not only to consistently perform 3D-aware object editing but also facilitates scene-level novel-view synthesis.

**Single-view 3D Scene Synthesis.** Among 3D scenes generation [Chen et al. 2023b,c; Chung et al. 2023; Epstein et al. 2024; Höllein et al. 2023; Mao et al. 2023; Zhang et al. 2023b], conditional generation on a single-view presents a unique challenge. Previous approaches address this challenge by training a versatile model capable of inferring a 3D representation of a scene based on a single input image [Flynn et al. 2019; Han et al. 2022; Hong et al. 2023; Hu et al. 2021; Li et al. 2021; Tucker and Snavely 2020; Wiles et al. 2020; Yu et al. 2021]. However, these methods demand extensive datasets for training and tend to produce blurry textures when confronted with significant changes in camera viewpoints. Recently, several works have embraced diffusion priors [Chan et al. 2023; Gu et al. 2023; Liu et al. 2023; Qian et al. 2023; Tang et al. 2023; Xu et al. 2023] to acquire a probabilistic distribution for unseen views, leading to better synthesis results. Nevertheless, these methods often concentrate on object-centric scenes or lack 3D consistency. Our approach connect 2D images and 3D scenes with explicit 3D Gaussians and incorporate diffusion knowledge, which overcome the mentioned challenges.

## 3 METHOD

Our target is to propose a 3D-aware scene image editing framework that allows simultaneous control over the camera and objects. To accomplish this, Sec. 3.1 introduces a novel scene representation called language-guided disentangled Gaussian splatting. In order to achieve object-level control, Sec. 3.2 further distills language features into the Gaussian splatting representation, achieving disentanglement at the object level. We elaborate the optimization process in Sec. 3.3 and demonstrate the flexible user control enabled by our framework during inference in Sec. 3.4.

### 3.1 3D Gaussian Splatting from Single Image

**Preliminary.** 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023] has been proved effective in both reconstructive [Luiten et al. 2023; Yang et al. 2023] and generative setting [Tang et al. 2023; Zou et al.

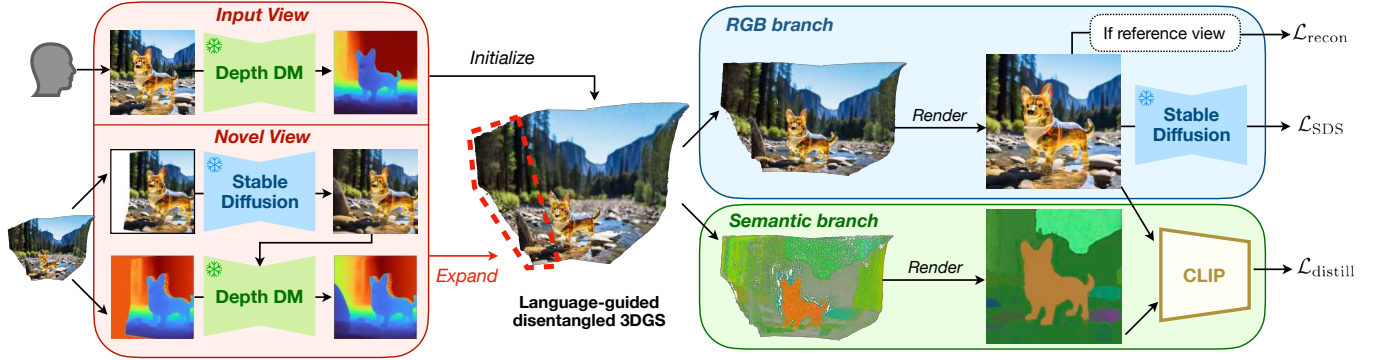


Fig. 2. **3DitScene training pipeline.** Given input view, we first initialize 3DGS by lifting pixels to 3D space and then expand it over novel views by RGB and depth inpainting. Semantic features are then distilled into 3D Gaussians to achieve object-level disentanglement.

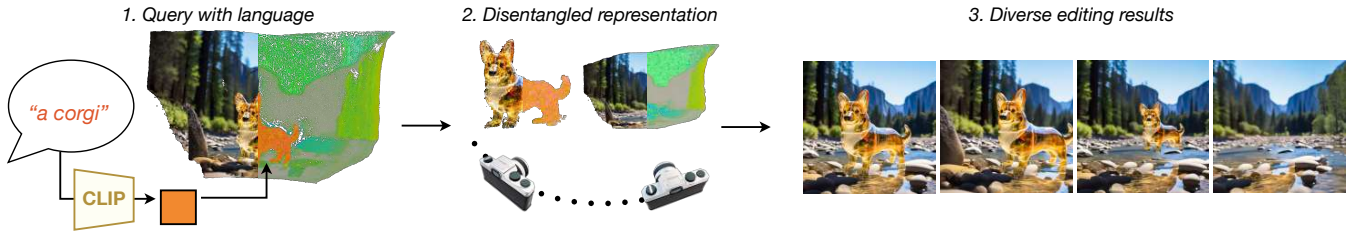


Fig. 3. **3DitScene Inference pipeline.** User can query object of interest via language prompt. Enabled by the disentangled 3D representation, user can change camera viewpoint, and manipulate the object of interest in a flexible manner.

2023]. It represents a 3D scene via a set of explicit 3D Gaussians. Each 3D Gaussian describes its location by a center vector  $\mathbf{x} \in \mathbb{R}^3$ , a scaling factor  $\mathbf{s} \in \mathbb{R}^3$ , a rotation quaternion  $\mathbf{q} \in \mathbb{R}^4$ , and also stores an opacity value  $\alpha \in \mathbb{R}$  and spherical harmonics (SH) coefficients  $\mathbf{c} \in \mathbb{R}^k$  ( $k$  represents the degrees of freedom of SH) for volumetric rendering. All the above parameters can be denoted as  $\Theta = \{\mathbf{x}_i, \mathbf{s}_i, \mathbf{q}_i, \alpha_i, \mathbf{c}_i | i \in [0, \dots, N-1]\}$ , where  $N$  is the number of 3D Gaussians. A tile-based rasterizer is used to render these Gaussians into 2D image.

**Image-to-3DGS initialization.** Given an input image  $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$ , an off-the-shelf depth prediction model is applied to estimate its depth map  $\mathbf{D} \in \mathbb{R}^{H \times W}$ . Then, we could transform image pixels into 3D space, forming the corresponding 3D point clouds:

$$\mathcal{P} = \phi_{2 \rightarrow 3}(\mathbf{I}, \mathbf{D}, \mathbf{K}, \mathbf{T}), \quad (1)$$

where  $\mathbf{K}$  and  $\mathbf{T}$  are camera intrinsic and extrinsic matrices respectively. Such point clouds  $\mathcal{P}$  are then used to initialize the 3DGS by directly copying the location and color values, with other GS-related parameters randomly initialized. To refine the 3DGS’s appearance, we adopt a reconstruction loss:

$$\mathcal{L}_{\text{recon}} = \|\mathbf{I} - f(\mathcal{P}, \mathbf{K}, \mathbf{T})\|_2^2, \quad (2)$$

where  $f$  is the rendering function.

We further enhance the rendered quality by leveraging prior knowledge from image generative foundation model, namely Stable Diffusion [Rombach et al. 2022]. It provides update direction to

the images rendered by the current 3DGS in the form of Score Distillation Sampling [Poole et al. 2022] loss, denoted as  $\mathcal{L}_{\text{SDS}}$ .

**3DGS expansion by inpainting.** When camera perspectives changes, rendered views will contain holes due to occlusion or new region outside the original view frustum. We use Stable Diffusion to inpaint the uncovered regions. Then, the newly added pixels need to be accurately transformed into 3D space to align seamlessly with the existing 3D Gaussians.

Previous methods [Chung et al. 2023; Höllein et al. 2023; Yu et al. 2023] first predict the depth values, and then use heuristic methods to adjust the values to align with the existing 3D structure. However, relying on heuristic methods often overlooked various scenarios, leading to artifacts such as depth discontinuities or shape deformations.

Instead, we propose a novel method to lifted novel contents to 3D while ensuring seamless alignment without any heuristic procedures. The key insight is to treat the problem as an image inpainting task, and utilize state-of-the-art diffusion-based depth estimation models [Fu et al. 2024; Yang et al. 2024] as a prior to solve the task. During denoising steps, rather than using models to predict the noise over the entire image, we employ the forward diffusion process to determine the value of fixed areas [Meng et al. 2021b]. This approach guarantees the final result, after denoising, adheres to the depth of original fixed parts, ensuring smooth expansion.

After smooth 3DGS expansion via depth inpainting, we take the imagined novel views as reference views, and apply reconstruction loss  $\mathcal{L}_{\text{recon}}$  to supervise the updated 3DGS. SDS loss  $\mathcal{L}_{\text{SDS}}$



is adopted for views rendered from camera perspectives that are interpolated between the user-provided viewpoint and the newly imagined views.

### 3.2 Language-guided Disentangled Gaussian Splatting

Based on the 3DGS built from single input image, users can generate novel views. In this section, we further distill CLIP [Radford et al. 2021] language feature to 3D Gaussians. This introduces semantics into 3D geometry, which helps disentangle individual objects out of the entire scene representation.

**Language feature distillation.** We augment each 3D Gaussian with a language embedding  $\mathbf{e} \in \mathbb{R}^C$ , where  $C$  denotes the number of the channels. Similar to RGB image  $\mathbf{I}$ , a 2D semantic feature map  $\mathbf{E} \in \mathbb{R}^{C \times H \times W}$  can also be rendered by the rasterizer. To learn the embedding, we first use Segment Anything Model (SAM) [Kirillov et al. 2023; Zhang et al. 2023a] to get semantic masks  $\mathbf{M}_i$ . Then, we can obtain embedding of each object  $\mathbf{I} \odot \mathbf{M}_i$  and supervise the corresponding region on rendered feature map  $\mathbf{E}$ , according to the distillation loss:

$$\mathcal{L}_{\text{distill}} = \sum_i \left\| (\mathbf{E} - g(\mathbf{I} \odot \mathbf{M}_i)) \odot \mathbf{M}_i \right\|_2^2, \quad (3)$$

where  $g$  is the CLIP’s image encoder, and  $\odot$  denotes element-wise multiplication. Following LangSplat [Qin et al. 2024], we additionally train an autoencoder to compress the embedding space to optimize the memory consumption of language embedding  $\mathbf{e}$ .

**Scene decomposition.** After distillation, we can decompose the scene into different objects. This enables user to query and ground specific object, and perform editing over single object (e.g. translation, rotation, removal, re-stylizing).

It is worth noting that such scene decomposition property not only enables more flexible edits during inference stage, but also provides augmentation over scene layouts during the optimization process. Since now we can query and render each object independently, we apply random translation, rotation, and removal over objects. This augmentation over the scene layout leads to a significant improvement in the appearance of occluded regions, ultimately enhancing the overall quality of the edited views (see Sec. 4.4).

### 3.3 Training

The overall training objective can be expressed as:

$$\mathcal{L} = \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} + \lambda_{\text{SDS}} \mathcal{L}_{\text{SDS}} + \lambda_{\text{distill}} \mathcal{L}_{\text{distill}}, \quad (4)$$

where  $\lambda_{\text{recon}}$ ,  $\lambda_{\text{SDS}}$  and  $\lambda_{\text{distill}}$  are coefficients that balance each loss term.

### 3.4 Inference

Due to the disentangled nature of our representation, users can now interact with and manipulate objects in a flexible manner. Here, we mainly discuss prompting objects via two different modalities:

**Text prompt.** Users can query an object through text prompts as shown in Fig. 3. Following LERF [Kerr et al. 2023] and LangSplat [Qin et al. 2024], we calculate the relevancy score  $\text{score}$  between the language embedding  $\mathbf{e}$  in the 3D Gaussians and the embedding of

Table 1. **User study result.** We report the percentage of favorite users for the consistency and quality of images edited by each method

		AnyDoor	Object 3DIT	Image Sculpting	Ours
Consistency	Human	5.1	16.8	12.7	<b>65.4</b>
	GPT4-v	0.0	6.7	31.3	<b>62.0</b>
Quality	Human	10.4	0.5	25.1	<b>64.0</b>
	GPT4-v	6.7	13.3	39.2	<b>40.8</b>

the text prompt  $\mathbf{e}_l$  as:

$$\text{score} = \min_i \frac{\exp(\mathbf{e} \cdot \mathbf{e}_l)}{\exp(\mathbf{e} \cdot \mathbf{e}_l) + \exp(\mathbf{e} \cdot \mathbf{e}_{\text{canon}}^i)}, \quad (5)$$

where  $\mathbf{e}_{\text{canon}}^i$  is the CLIP embeddings of canonical phrases including “object”, “things”, “stuff”, and “texture”. Gaussians that have relevance scores below a predefined threshold are excluded. The remaining part is identified as the object of user interest.

**Bounding box.** Users can also select an object by drawing an approximate bounding box around it on the input image. 3D Gaussians within the bounding box are first identified, followed by K-Means clustering based on their language embeddings  $\mathbf{e}$ . Assuming the object is the most significant one within the bounding box, clusters whose number of Gaussians does not exceed a threshold proportion will be discarded.

In the meantime, user can also adjust the camera viewpoint by specifying intrinsic and extrinsic parameters.

## 4 EXPERIMENTS

### 4.1 Settings

**Implementation details.** To lift an image to 3D, we use GeoWizard [Fu et al. 2024] to estimate its relative depth. Stable Diffusion [Rombach et al. 2022]’s inpainting pipeline is adopted to generate new content for 3DGS’s expansion. We leverage MobileSAM [Zhang et al. 2023a] and OpenCLIP [Ilharco et al. 2021] to segment and compute rendered views’ feature maps, which are further leveraged to supervise the language embedding of 3D Gaussians. We use Stable Diffusion to perform Score Distillation Sampling [Poole et al. 2022] during optimization. Given the already decent image quality at the start of optimization benefited from explicit 3DGS initialization, we adopt a low classifier-free guidance [Ho and Salimans 2022] scale.

**Baselines.** We compare our method with following scene image editing works: (1) AnyDoor [Chen et al. 2023a] is a 2D diffusion-based model that can teleport target objects into given scene images. It leverages Stable Diffusion’s powerful image generative prior by finetuning upon it. (2) Object 3DIT [Michel et al. 2024] is designed for 3D-aware object-centric image editing via language instructions. It finetunes Stable Diffusion over a synthetic dataset containing pairs of original image, language instruction, and edited image. (3) Image Sculpting [Yenphraphai et al. 2024b] is also designed for 3D-aware object-centric image editing. It estimates a 3D model from an object in the input image to enable precise 3D control over the geometry. It also uses Stable Diffusion to refine the edited image quality. (4) AdaMPI [Han et al. 2022] focuses on the control over camera perspective. It leverages monocular depth estimation and color

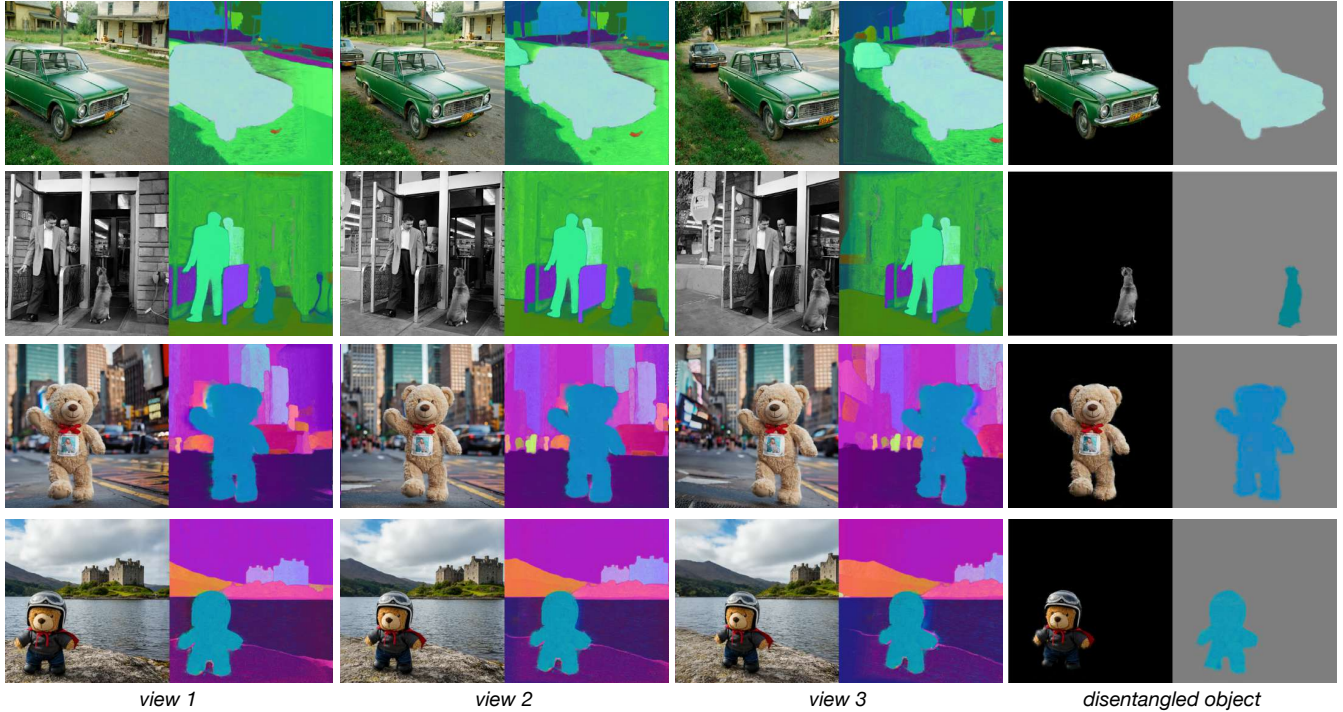


Fig. 4. **Visualization of rendered images and feature maps.** For each sample, we show three views of rendered images and feature maps. To demonstrate the disentangled scene representation, we use the language embedding to select a foreground object and render it exclusively.

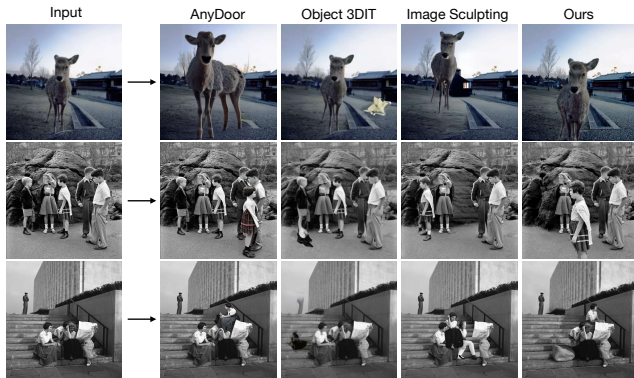


Fig. 5. **Comparison results of object-centric manipulation.** We apply translation, resizing, and removal over foreground objects. Two different edit results are shown for each method.

inpainting with learned adaptive layered depth representations. (5) LucidDreamer [Chung et al. 2023] tackles novel view synthesis by querying Stable Diffusion’s inpainting pipeline with dense camera trajectories.

## 4.2 Quantitative results

We conduct a user study to compare the edited results by our method with the established baselines. We generate 20 samples for each method and request users to vote for their preferred method

based on consistency with the original image and quality for each sample. We collect feedback from 25 users, and report the result in Tab. 1. Our method consistently outperforms previous baselines in terms of both consistency and image quality. As recommended in a previous study [Wu et al. 2024], GPT4-v has the ability to evaluate 3D consistency and image quality. Therefore, we include GPT4-v as an additional criterion. The preference of GPT4-v is well aligned with human preference, which once again demonstrates the superiority of 3DiTScene.

## 4.3 Qualitative results

Fig. 4 showcases the generated novel views with their respective feature maps produced by our framework. The feature maps demonstrate remarkable accuracy in capturing the semantic content of the images. This ability to distinctly separate semantic information plays a crucial role in achieving precise object-level control. In the following, we demonstrate flexible editing over scene images enabled by our framework, and also compare with baseline methods. **Object manipulation.** Since different methods define object manipulation, particularly translation operations, in different coordinate systems<sup>1</sup>, it becomes challenging to evaluate them under a unified and fair setting. Therefore, we evaluate each method under its own specific setting to achieve the best possible result. As shown in Fig. 5, AnyDoor struggles to maintain object identity and 3D consistency

<sup>1</sup>AnyDoor, Object 3DIT and Image Sculpting respectively employs 2D masks, language prompts, and image coordinates for control. We use coordinates in 3D space instead.



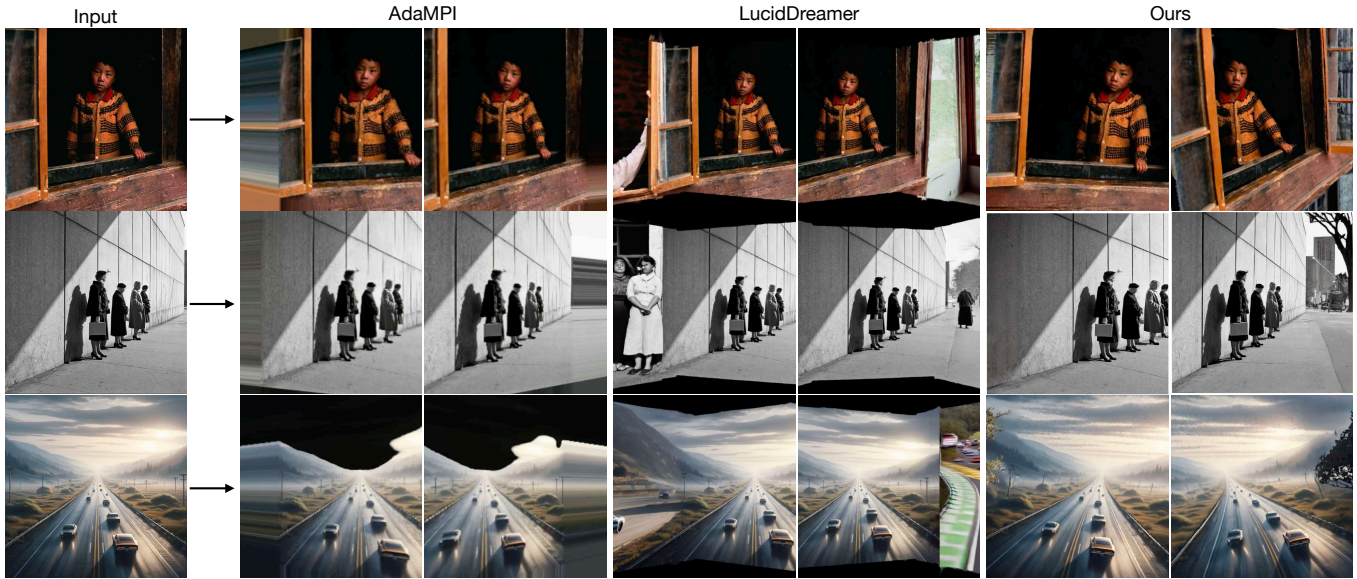


Fig. 6. **Comparison results of camera control.** We show two views with different camera perspectives for each method.

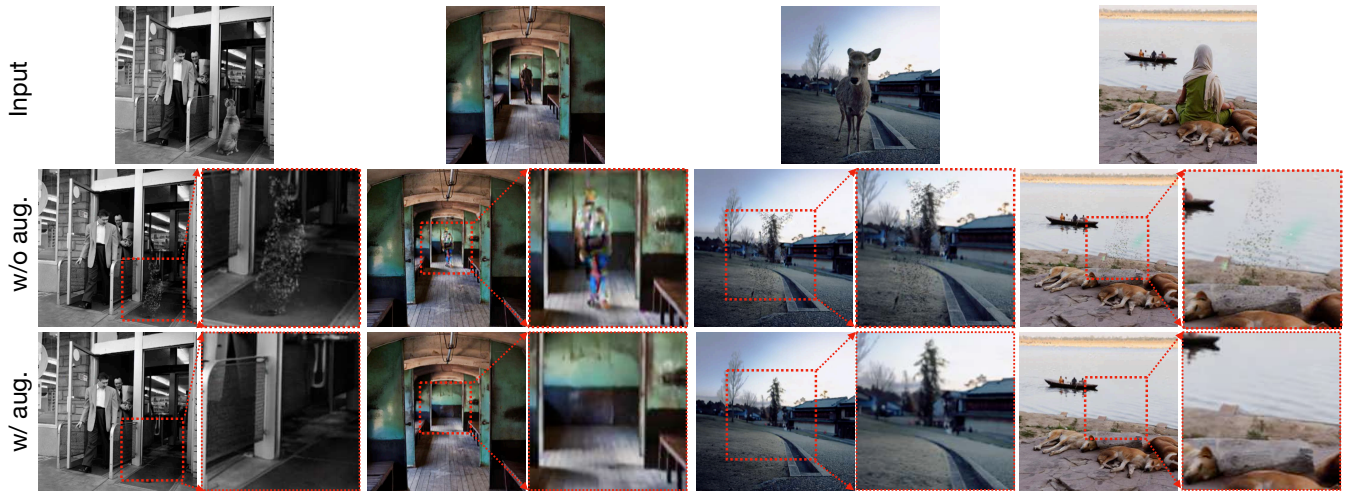


Fig. 7. **Ablation results for layout augmentation during optimization.** To evaluate the degree of object-level disentanglement, we conduct object removal for each sample. The top row displays the input image, while the next two rows showcase the edited scene

when manipulating object layouts, primarily due to the absence of 3D cues. Object 3DIT, trained on synthetic datasets, exhibits limited generalization ability to real images. By leveraging a 3D model derived from the input image, Image Sculpting achieves better results. Nonetheless, it encounters issues with inconsistency when manipulating objects. This arises from the fact that they solely rely on the 3D model for providing rough guidance, resulting in the loss of finer details during optimization.

In contrast, our method delivers satisfactory 3D-aware object-level editing results. It maintains accurate 3D consistency of edited objects after rearranging their layout. Additionally, it preserves

occlusion relationships within the scene, such as moving the girl to be partially occluded by a foreground object in the last row example. **Camera control.** We compare our methods with AdaMPI and LucidDreamer for camera control. As illustrated in Fig. 6, AdaMPI only focuses on scenarios where the camera zooms in, and does not consider novel view synthesis. Therefore, this approach is not suitable for 3D-aware image editing when large camera control is required. LucidDreamer also leverages Stable Diffusion’s inpainting capacity for novel view synthesis. However, it suffers from sudden transitions in the content within the frame (see sample in the bottom line). It also requires dense camera poses. In contrast, our



Fig. 8. **Ablation results for loss terms.** We show rendered novel views under different loss settings. The left column lists the input image. In right columns, two views are shown for each configuration. The quality degrades when reconstruction or SDS loss term is discarded

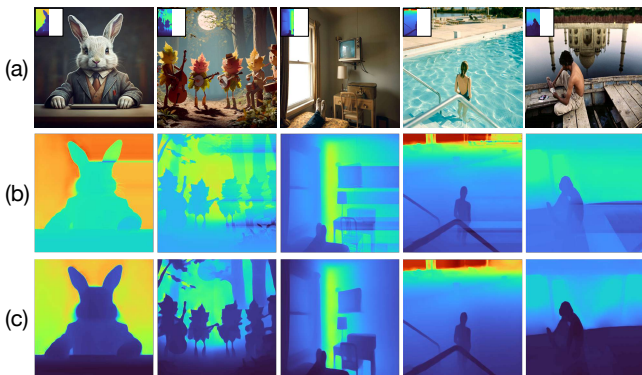


Fig. 9. **Ablation results for depth inpainting.** Row (a): images and corresponding depth map (only available in the left half); Row (b): depth map predicted by heuristic alignment; Row (c): depth map predicted by our depth inpainting method.

method only needs as few as three camera poses and enables smooth transitions from the input view to novel views, enhancing user control over the camera perspective.

#### 4.4 Ablation study

**Layout augmentation during optimization.** As our representation disentangles at object level, we could perform layout augmentation during optimization. Here, we investigate whether disentanglement property benefits the optimization process. We use the task of removing objects to evaluate the degree of disentanglement.

As illustrated in Fig. 7, when layout augmentation is disabled during optimization, floating artifacts can be observed. We discover that these Gaussians lie inside the object. They are occluded by Gaussians at the surface. As they do not contribute to the rendering result, they are consequently not updated by gradient descent during optimization, leaving their language embeddings unsupervised.

In contrast, when applying layout augmentation during optimization, such Gaussians will be exposed when the foreground object is moved away, and hence updated. With this ablation, it is concluded that the disentanglement property of the proposed representation

not only enables more flexible inference, but also contributes to the optimization process.

**Loss terms.** During optimization, we adopt three loss terms:  $\mathcal{L}_{\text{recon}}$ ,  $\mathcal{L}_{\text{SDS}}$ , and  $\mathcal{L}_{\text{distill}}$ .  $\mathcal{L}_{\text{distill}}$  plays a critical role in distilling language embedding into 3D. The remaining two terms focus on enhancing the visual quality of images. Here, we investigate the contribution of these two items through an ablation study. As input images can provide guidance of overall structure and detailed appearance, there is no need of applying a large classifier free guidance (CFG) value for SDS loss. Thus, by default, we choose 5 as the CFG value.

As illustrated in Fig. 8, the image quality degrades severely without  $\mathcal{L}_{\text{recon}}$  or  $\mathcal{L}_{\text{SDS}}$ . Without  $\mathcal{L}_{\text{recon}}$ , the image is only refined by the SDS loss, which creates discrepancies with the original image. When the CFG value is set low, 5 as default, the image appears lacking in details and exhibits unusual texture patterns. Increasing the CFG value introduces more details, yet leads to inconsistencies with the original image, while the issue of strange texture patterns persists. Additionally, only applying  $\mathcal{L}_{\text{recon}}$  results to floating artifacts and blurriness across the entire image. In conclusion, both SDS and reconstruction loss are crucial for achieving decent image quality.

**Depth inpainting.** When expanding 3DGS at novel views, we need to estimate the depth map of unseen regions. Here, we compare our inpainting-based depth estimation with heuristic-based method. Fig. 9 show images with depth map available in the left part. The task is to predict the depth map of the right part. Method relying on heuristic alignment results to artifacts like depth discontinuity. In contrast, our proposed method is capable of producing accurate depth maps that align well with the left known part.

## 5 CONCLUSION AND DISCUSSION

We present a novel framework, 3Di tScene, for scene image editing. Our primary objective is to facilitate 3D-aware editing of both objects and the entire scene within a unified framework. We achieve this by leveraging a new scene representation, language-guided disentangled scene representation. This representation is learnt by distilling CLIP’s language feature into 3D Gaussians. The semantic 3D Gaussians effectively disentangle individual objects out of the entire scene, thereby enabling localized object editing. We test 3Di tScene under different settings and prove its superiority compared to previous methods.



## REFERENCES

- Eric R Chan, Koki Nagano, Matthew A Chan, Alexander W Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. 2023. Generative novel view synthesis with 3d-aware diffusion models. *ICCV (2023)*.
- Dave Zhenyu Chen, Haoxuan Li, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. 2023b. Scenetex: High-quality texture synthesis for indoor scenes via diffusion priors. *arXiv preprint arXiv:2311.17261 (2023)*.
- Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. 2023a. Anydoor: Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481 (2023)*.
- Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. 2023c. Scenedreamer: Unbounded 3d scene generation from 2d image collections. *arXiv preprint arXiv:2302.01330 (2023)*.
- Jaeyoung Chung, Suyoung Lee, Hyeonjin Nam, Jaerin Lee, and Kyoung Mu Lee. 2023. Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384 (2023)*.
- Dave Epstein, Ben Poole, Ben Mildenhall, Alexei A Efros, and Aleksander Holynski. 2024. Disentangled 3D Scene Generation with Layout Learning. *arXiv preprint arXiv:2402.16936 (2024)*.
- John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. 2019. Deepview: View synthesis with learned gradient descent. In *CVPR*.
- Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuxin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and Xiaoxiao Long. 2024. GeoWizard: Unleashing the Diffusion Priors for 3D Geometry Estimation from a Single Image. *arXiv preprint arXiv:2403.12013 (2024)*.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bernano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618 (2022)*.
- Jiatao Gu, Alex Trevithick, Kai-En Lin, Joshua M Susskind, Christian Theobalt, Lingjie Liu, and Ravi Ramamoorthi. 2023. Nerfdiff: Single-image view synthesis with nerf-guided distillation from 3d-aware diffusion. In *ICML*.
- Yuxuan Han, Ruicheng Wang, and Jiaolong Yang. 2022. Single-view view synthesis in the wild with learned adaptive multiplane images. In *ACM SIGGRAPH Conference Proceedings*.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2022. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626 (2022)*.
- Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. 2023. Style aligned image generation via shared attention. *arXiv preprint arXiv:2312.02133 (2023)*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. (2020).
- Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598 (2022)*.
- Lukas Höllein, Ang Cao, Andrew Owens, Justin Johnson, and Matthias Nießner. 2023. Text2room: Extracting textured 3d meshes from 2d text-to-image models. *arXiv preprint arXiv:2303.11989 (2023)*.
- Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. 2023. LRM: Large Reconstruction Model for Single Image to 3D. *arXiv preprint arXiv:2311.04400 (2023)*.
- Ronghang Hu, Nikhila Ravi, Alexander C Berg, and Deepak Pathak. 2021. Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. In *ICCV*.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. 2021. *OpenCLIP*. <https://doi.org/10.5281/zenodo.5143773> If you use this software, please cite it as below..
- Ali Jahani, Lucy Chai, and Phillip Isola. 2019. On the "steerability" of generative adversarial networks. *arXiv preprint arXiv:1907.07171 (2019)*.
- Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2021. Alias-Free Generative Adversarial Networks.
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *CVPR*.
- Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. 2023. Imagic: Text-based real image editing with diffusion models. In *CVPR*.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (2023).
- Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. 2023. Lrf: Language embedded radiance fields. In *CVPR*. 19729–19739.
- Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. 2022. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *CVPR*.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. *arXiv preprint arXiv:2304.02643 (2023)*.
- Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. 2021. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *ICCV*.
- Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 2023. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9298–9309.
- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. 2023. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713 (2023)*.
- Weijia Mao, Yan-Pei Cao, Jia-Wei Liu, Zhongcong Xu, and Mike Zheng Shou. 2023. ShowRoom3D: Text to High-Quality 3D Room Generation Using 3D Priors. *arXiv preprint arXiv:2312.13324 (2023)*.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2021a. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073 (2021)*.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2021b. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073 (2021)*.
- Oscar Michel, Anand Bhattad, Eli VanderBilt, Ranjay Krishna, Aniruddha Kembhavi, and Tanmay Gupta. 2024. Object 3dit: Language-guided 3d-aware image editing. *Advances in Neural Information Processing Systems* 36 (2024).
- Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *CVPR*.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988 (2022)*.
- Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. 2023. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843 (2023)*.
- Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. 2024. LangSplat: 3D Language Gaussian Splatting. In *CVPR*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *CVPR*.
- Natanuel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*.
- Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. 2020. InterFaceGAN: Interpreting the Disentangled Face Representation Learned by GANs. *IEEE TPAMI (2020)*.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2020. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456 (2020)*.
- Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. 2022. Dual diffusion implicit bridges for image-to-image translation. *arXiv preprint arXiv:2203.08382 (2022)*.
- Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. 2023. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653 (2023)*.
- Richard Tucker and Noah Snavely. 2020. Single-view view synthesis with multiplane images. In *CVPR*.
- Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. 2020. Synsin: End-to-end view synthesis from a single image. In *CVPR*.
- Tong Wu, Guandao Yang, Zhibing Li, Kai Zhang, Ziwei Liu, Leonidas Guibas, Dahua Lin, and Gordon Wetzstein. 2024. GPT-4V (ision) is a Human-Aligned Evaluator for Text-to-3D Generation. *arXiv preprint arXiv:2401.04092 (2024)*.
- Yinghao Xu, Yujun Shen, Jiapeng Zhu, Ceyuan Yang, and Bolei Zhou. 2021. Generative Hierarchical Features from Synthesizing Images. In *CVPR*.
- Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein, Zexiang Xu, et al. 2023. Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. *arXiv preprint arXiv:2311.09217 (2023)*.
- Ceyuan Yang, Yujun Shen, and Bolei Zhou. 2021. Semantic hierarchy emerges in deep generative representations for scene synthesis. *IJCV (2021)*.
- Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth anything: Unleashing the power of large-scale unlabeled data. *arXiv preprint arXiv:2401.10891 (2024)*.
- Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. 2023. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101 (2023)*.
- Jiraphon Yenphraphai, Xichen Pan, Sainan Liu, Daniele Panozzo, and Saining Xie. 2024a. Image Sculpting: Precise Object Editing with 3D Geometry Control. *arXiv preprint arXiv:2401.01702 (2024)*.
- Jiraphon Yenphraphai, Xichen Pan, Sainan Liu, Daniele Panozzo, and Saining Xie. 2024b. Image Sculpting: Precise Object Editing with 3D Geometry Control. *arXiv preprint arXiv:2401.01702 (2024)*.



- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021. pixelnerf: Neural radiance fields from one or few images. In *CVPR*.
- Hong-Xing Yu, Haoyi Duan, Junhwa Hur, Kyle Sargent, Michael Rubinstein, William T Freeman, Forrester Cole, Deqing Sun, Noah Snively, Jiajun Wu, et al. 2023. WonderJourney: Going from Anywhere to Everywhere. *arXiv preprint arXiv:2312.03884* (2023).
- Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, Sung-Ho Bae, Seungkyu Lee, and Choong Seon Hong. 2023a. Faster Segment Anything: Towards Lightweight SAM for Mobile Applications. *arXiv preprint arXiv:2306.14289* (2023).
- Qihang Zhang, Chaoyang Wang, Aliaksandr Siarohin, Peiye Zhuang, Yinghao Xu, Ceyuan Yang, Dahua Lin, Bolei Zhou, Sergey Tulyakov, and Hsin-Ying Lee. 2023b. Scenewiz3d: Towards text-guided 3d scene composition. *arXiv preprint arXiv:2312.08885* (2023).
- Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. 2020. In-domain gan inversion for real image editing. In *ECCV*.
- Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. 2023. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. *arXiv preprint arXiv:2312.09147* (2023).

“Remove the sheep and Rotate the camera”



“Move the boy, then Rotate the camera”



“Rotate the camera, and Replace the rabbit with cat”



Fig. 10. Image pairs edited by 3DitScene.



“Move the toy bear closer, and Rotate the camera”



“Replace the chicken made of clay with a ball of yarn, then Rotate the camera”



“Rotate the camera, and Remove the woman”

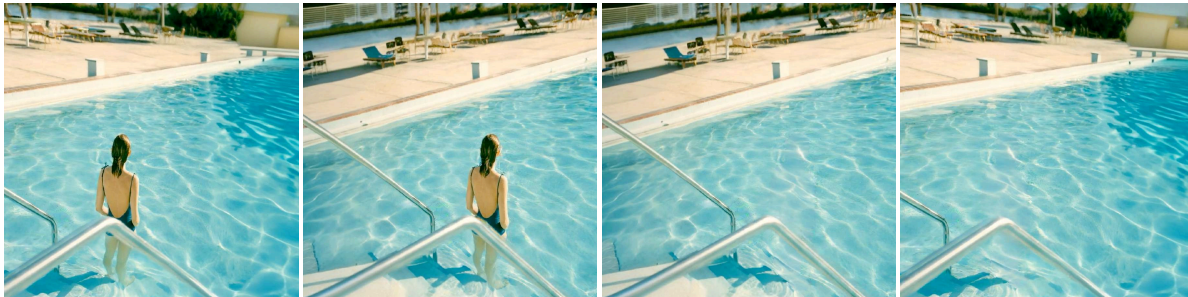


Fig. 11. Image pairs edited by 3DitScene.