

LucidDreamer: Domain-free Generation of 3D Gaussian Splatting Scenes

Jaeyoung Chung^{*} Suyoung Lee^{*} Hyeongjin Nam Jaerin Lee Kyoung Mu Lee
 ASRI, Department of ECE, Seoul National University, Seoul, Korea
 {robot0321, esw0116, namhjsnu28, ironjr, kyoungmu}@snu.ac.kr

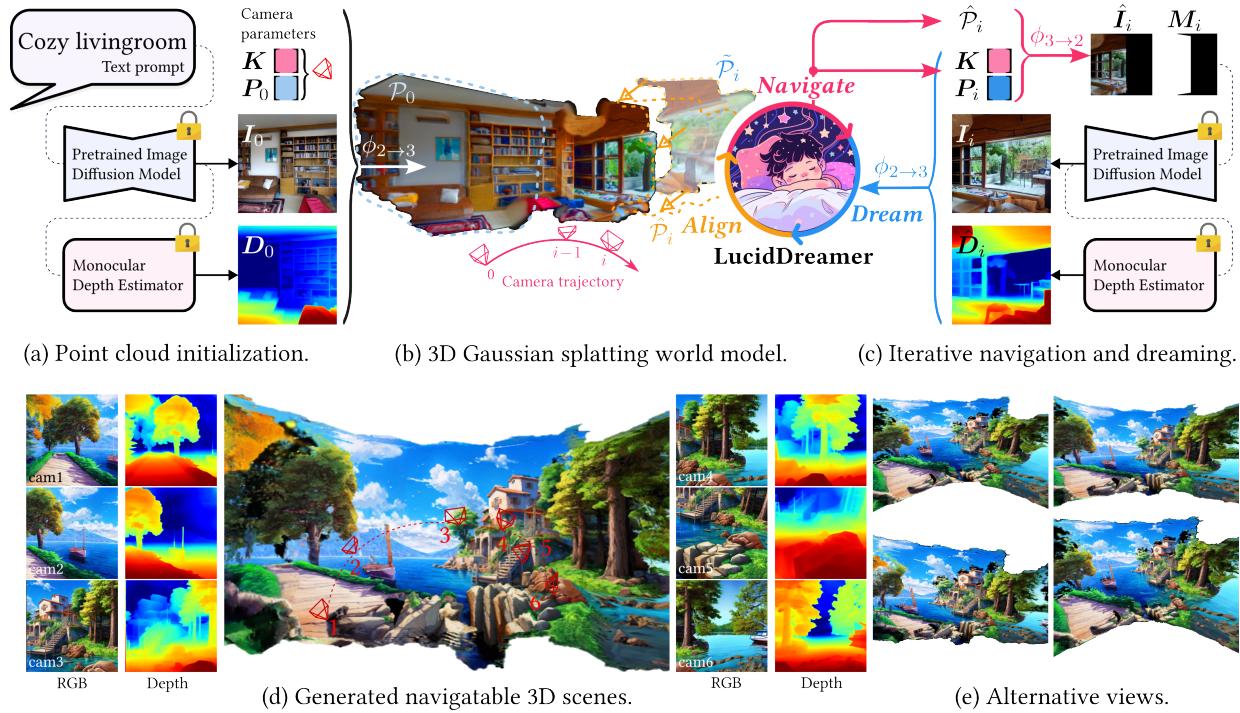


Figure 1. **Introducing LucidDreamer.** We develop *LucidDreamer*, a general framework for generating multiview-consistent and high-quality 3D scenes from various input types: text, RGB, and RGBD. After the initial point cloud is created by lifting the RGBD image, *LucidDreamer* maintains and expands its world model by repeating two operations: *Dreaming* and *Alignment*. The 3D scene is finalized through optimizing a Gaussian splatting representation.

Abstract

With the widespread usage of VR devices and contents, demands for 3D scene generation techniques become more popular. Existing 3D scene generation models, however, limit the target scene to specific domain, primarily due to their training strategies using 3D scan dataset that is far from the real-world. To address such limitation, we propose *LucidDreamer*, a domain-free scene generation pipeline by fully leveraging the power of existing large-scale diffusion-based generative model. Our *LucidDreamer* has two alternate steps: *Dreaming* and *Alignment*. First, to generate multi-view consistent images from inputs, we set the point

cloud as a geometrical guideline for each image generation. Specifically, we project a portion of point cloud to the desired view and provide the projection as a guidance for inpainting using the generative model. The inpainted images are lifted to 3D space with estimated depth maps, composing a new points. Second, to aggregate the new points into the 3D scene, we propose an aligning algorithm which harmoniously integrates the portions of newly generated 3D scenes. The finally obtained 3D scene serves as initial points for optimizing Gaussian splats. *LucidDreamer* produces Gaussian splats that are highly-detailed compared to the previous 3D scene generation methods, with no constraint on domain of the target scene.

^{*} indicates equal contribution.

1. Introduction

With the advent of commercial mixed reality platforms and the rapid innovations in 3D graphics technology, high-quality 3D scene generation has become one of the most important problem in computer vision. This requires the ability to create diverse and photo-realistic 3D scenes from any type of input, such as text, RGB, and RGBD images. There are efforts to use the diffusion model in voxel, point cloud, and implicit neural representation to generate 3D objects and scenes directly [11, 25, 60], but the results show low diversity and quality due to the limitations in training data based on 3D scans. One way to cope with the issue is to leverage the power of a pre-trained image generation diffusion model, such as Stable Diffusion [39], to create diverse high-quality 3D scenes. Such a big model creates plausible images with a data-driven knowledge learned from the large-scale training data, although it does not guarantee multi-view consistency between the generated images [51].

In this work, we propose a pipeline called **LucidDreamer** that utilizes Stable Diffusion [39] and 3D Gaussian splatting[20] to create diverse high-quality 3D scenes from various types of inputs such as text, RGB, and RGBD. Following the pipeline of LucidDreamer, a unified large point cloud is generated by repeating the two processes named Dreaming and Alignment, alternatively. Before beginning the two process, an initial point cloud is generated by the initial image and the corresponding depth map. Dreaming process includes the generation of geometrically consistent images and the lifting of these images into 3D space. We first move the camera along the pre-defined camera trajectory and project a visible region of point cloud in the new camera coordinate to the new camera plane. Then, the projected image is put into the Stable Diffusion-based inpainting network to generate the complete image from the projected one. A new set of 3D points are generated by lifting the inpainted image and the estimated depth map to the 3D space. Then the proposed alignment algorithm seamlessly connects the new 3D points to the existing point cloud by slightly moving the position of the new points in the 3D space. After the large point cloud generated by repeating the above processes a sufficient number of time is obtained, we use it as the initial SfM points to optimize the Gaussian splats. The continuous representation of 3D Gaussian splats removes the holes generated by the depth discrepancy in the point cloud, enabling us to render more photo-realistic 3D scenes than traditional representations. Figure 1 shows the simple process of LucidDreamer and a 3D generation result.

LucidDreamer exhibits significantly more realistic and astonishing results compared to existing models. We compare the generated 3D scenes conditioned with an image from ScanNet [9], NYUDepth [45], and Stable Diffusion, and show better visual results across all datasets. Our model is capable of generating 3D scenes across diverse domains

such as realistic/anime/lego and indoor/outdoor. Not only does our model support various domains, but it also accommodates the simultaneous use of diverse input conditions. For example, by conditioning an image and text together, it generates a 3D scene based on the text but also includes the image. This alleviates the challenges associated with creating the desired scene solely from the text, moving away from generating samples exhaustively. Furthermore, our approach also allows for the change of the input condition while creating the 3D space. These capabilities offer opportunities to create a wide range of 3D scenes, inspiring creativity.

In summary, our contributions are as follows.

- We introduce LucidDreamer, a domain-free high-quality 3D scene generation, achieving better domain generalization in 3D scene generation by leveraging the power of Stable Diffusion, depth estimation, and explicit 3D representation.
- To generate multi-view images from Stable Diffusion, our **Dreaming** process establishes point cloud as geometrical guideline for each image generation. Subsequently, our **Aligning** process harmoniously integrates the generated images to form an unified 3D scene.
- Our model provides users with the ability to create 3D scenes in various ways by supporting different input types, such as text, RGB, and RGBD, allowing the simultaneous use of multiple inputs, and enabling the change of the inputs during the generation process.

2. Related Work

3D Scene Representation. Representative methods for expressing 3D scenes include explicit methods such as point cloud, mesh, and voxel. These are widely used because they allow direct and intuitive control of each element and enable fast rendering through the rasterization pipeline. However, they need a large number of elements for detailed expression because of their simple structure. Complex primitives such as cuboid [52], Gaussian [12], ellipsoid [16], superquadrics [33], convex hull [10], and polynomial surface [56] were developed for more efficient expression. Although primitives have increased expressive power for complex geometry, it is still difficult to express realistic 3D scenes because of simple color representation.

Recently, there have been works to express more detailed 3D scenes using neural networks as implicit expressions. They train a neural network to express the scene creating the desired properties in 3D coordinates, such as signed distance function [32, 49], RGB α [26, 46]. In particular, Neural Radiance Fields [26] showed that it was possible to optimize photorealistic 3D scenes from multiple images through volume rendering, but the scene implicitly stored

in the network form is difficult to handle and slow. To improve this, subsequent studies attempted to use volume rendering in explicit expressions. By utilizing the locality of structures such as sparse voxels[14, 24, 48, 58], featured point clouds[55], Multi-Level Hierarchies [27, 28], tensor [5], infinitesimal networks [15, 38], triplane [4], polygon [7], and Gaussain splats [20] they greatly improve the training and rendering speed. In particular, 3D Gaussian splatting [20] utilizes the concept of Gaussian splats combined with spherical harmonics and opacity to represent complete and unbounded 3D scenes. It supports not only alpha-blending but also differentiable rasterization, resulting in fast, high-quality 3D scene optimization. This structure is essential for our generation method, which cannot determine the bounds of the scene due to sequential image generation, and plays a role in making the scene complete.

3D Scene Generation. Inspired by the early success of generative adversarial network (GAN) [17] in image generation, similar attempts are made in 3D creation. Creating a set of multiview consistent images [3, 31, 42], or directly creating voxel [29, 30, 54] or point cloud [1, 43] were studied. However, they suffer from GAN’s learning instability [36] and memory limitation in 3D representation, limiting the generation quality. Encouraged by the recent success of diffusion [19, 47] in the field of image generation [37, 39], there are many attempts to introduce the diffusion model into 3D representation, such as voxel [60], point cloud [25, 59], triplane [4, 6, 44], implicit neural network [11, 34, 57]. They use object-centric coordinates because of their nature and focus on simple examples. Some generative diffusion models overcome this problem by using a mesh as a proxy and diffusing in the UV space. They create a large portrait scene by continuously building the mesh [13] or create indoor scenes [8, 22] and more realistic objects [35, 50]. However, their performance falls short of foundation models[39] because they involve training a new diffusion model in a different representation space, which is limited by data availability and computational resources. In comparison, our method leverages the power of the foundation model to generate diverse images and creates reliable 3D scenes through depth estimation and optimization.

3. Method

While the range of target scenes of existing scene generation models is strictly restricted due to the limitations in training dataset, LucidDreamer can generate even more realistic, higher-resolution 3D scenes with much more general input conditions. For instance, LucidDreamer can generate a text-relevant scene if only the text prompt is given. Also, the style of the input image is maintained along the scene, while existing models keep producing scenes that are similar to the style of the training dataset, not the input image.

The pipeline of LucidDreamer is broadly divided into two stages: point cloud construction and Gaussian splats optimization. During the first stage, an initial point cloud is formed from the input image, and the area of the point cloud is expanded to create a large scene using Stable Diffusion inpainting and monocular depth estimation. Then, the point cloud and the reprojected images are used to optimize Gaussian splats. By representing the scene with Gaussian splats, we can fill the empty space that appears in the point cloud due to the depth discrepancy.

3.1. Point cloud construction

To generate multi-view consistent 3D point cloud, we create the initial point cloud and aggregate the points by moving back and forth between 3D space and the camera plane while moving the camera. The overall process of point cloud construction is illustrated in Figure 1.

Initialization. A point cloud generation starts from lifting the pixels of the initial image. If the user gives a text prompt as input, the latent diffusion model is used to generate an image relevant to the given text, and the depth map is estimated using the monocular depth estimation model such as ZoeDepth [2]. We denote the generated or received RGB image and the depth map as $\mathbf{I}_0 \in \mathbb{R}^{3 \times H \times W}$ and $\mathbf{D}_0 \in \mathbb{R}^{H \times W}$, where H and W are height and the width of the image. The camera intrinsic matrix and the extrinsic matrix of \mathbf{I}_0 are denoted as \mathbf{K} and \mathbf{P}_0 , respectively. For the case where \mathbf{I}_0 and \mathbf{D}_0 are generated from the diffusion model, we set the values of \mathbf{K} and \mathbf{P}_0 by convention regarding the size of the image.

From the input RGBD image $[\mathbf{I}_0, \mathbf{D}_0]$, we lift the pixels into the 3D space, where the lifted pixels will form a point cloud in a 3D space. The generated initial point cloud using the first image is defined as \mathcal{P}_0 :

$$\mathcal{P}_0 = \phi_{2 \rightarrow 3}([\mathbf{I}_0, \mathbf{D}_0], \mathbf{K}, \mathbf{P}_0), \quad (1)$$

where $\phi_{2 \rightarrow 3}$ is the function to lift pixels from the RGBD image $[\mathbf{I}, \mathbf{D}]$ to the point cloud.

Point cloud aggregation. We sequentially attach points to the original point cloud to create a large 3D scene. Specifically, we set the camera trajectory with length N , where \mathbf{P}_i indicates the position and pose of the camera in the i -th index, then inpaint and lift the missing pixel in each step. Here, the generated points should satisfy two conditions; the images projected from the points should have high perceptual quality and be consistent with image parts produced from the existing points. To achieve the former condition, we borrow the representation power of the Stable Diffusion [39] to the image inpainting task.

Navigation. At step i , we first move and rotate the camera from the previous position (\mathbf{P}_{i-1}) to \mathbf{P}_i . We change the

coordinate from the world to the current camera and project to the camera plane using \mathbf{K} and \mathbf{P}_i .

Dreaming. We denote the projected image at camera \mathbf{P}_i as $\hat{\mathbf{I}}_i$. Since the position and the pose of the camera are changed, there would be some regions that cannot be filled from the existing point cloud. We define the mask \mathbf{M}_i to discriminate the region that is filled by existing points in $\hat{\mathbf{I}}_i$. Specifically, the value of \mathbf{M}_i is one if the corresponding pixel is already filled or 0 otherwise. The Stable Diffusion inpainting model (\mathcal{S}) is executed to generate a realistic image, \mathbf{I}_i , from the incomplete image ($\hat{\mathbf{I}}_i$) and the mask (\mathbf{M}_i). The corresponding depth map ($\hat{\mathbf{D}}_i$) is estimated using the monocular depth estimation network (\mathcal{D}). Here, the monocular depth estimation model can only estimate the relative depth, and the depth coefficients from the relative depth to the actual depth can be different between images. If the depth coefficients are different, the lifted 3D point clouds in the two generated images are not connected and are spaced apart. We estimate the optimal depth scale coefficient, d_i , that minimizes the distance between the 3D points of the new image and the corresponding points in the original point cloud, \mathcal{P}_{i-1} . Then the actual depth map, \mathbf{D}_i is calculated by multiplying the coefficient d_i to the estimated depth map, $\hat{\mathbf{D}}_i$.

$$\begin{aligned} \mathbf{I}_i &= \mathcal{S}(\hat{\mathbf{I}}_i, \mathbf{M}_i), \hat{\mathbf{D}}_i = \mathcal{D}(\mathbf{I}_i), \mathbf{D}_i = d_i \hat{\mathbf{D}}_i, \\ d_i &= \operatorname{argmin}_d \left(\sum_{\mathbf{M}_i=1} \left\| \phi_{2 \rightarrow 3} \left([\mathbf{I}_i, d \hat{\mathbf{D}}_i], \mathbf{K}, \mathbf{P}_i \right) - \mathcal{P}_{i-1} \right\|_1 \right). \end{aligned} \quad (2)$$

Here, $\mathbf{M}_i = 1$ implies that the distance of point pairs in the overlapping regions is used for estimating d_i .

Using the image and the corresponding depth map, $[\mathbf{I}_i, \mathbf{D}_i]$, we lift the pixels to 3D space. Here, we note that only inpainted pixels of \mathbf{I}_i are lifted to prevent points overlapping and mitigate the inconsistency problem. The output of dreaming, $\hat{\mathcal{P}}_i$, can be calculated as:

$$\hat{\mathcal{P}}_i = \phi_{2 \rightarrow 3}([\mathbf{I}_i, \mathbf{D}_i | \mathbf{M}_i = 0], \mathbf{K}, \mathbf{P}_i), \quad (3)$$

where $[\mathbf{I}_i, \mathbf{D}_i | \mathbf{M}_i = 0]$ indicates the inpainted region in the RGBD image.

Alignment. Compared to the way that trains a generative model to generate both RGB and the depth map at once, such as RGBD2 [57], the depth map estimated by off-the-shelf depth estimation method is more accurate and generalizable to various situations since off-the-shelf methods are trained on large and various datasets. However, since $\mathbf{D}_0, \mathbf{D}_1, \dots, \mathbf{D}_{i-1}$ is not considered when estimating \mathbf{D}_i , inconsistency problem occurs when we add new points, $\hat{\mathcal{P}}_i$. To overcome the problem, we move the points of $\hat{\mathcal{P}}_i$ in 3D space to attach the two point cloud (\mathcal{P}_{i-1} and $\hat{\mathcal{P}}_i$) smoothly. Specifically, we extract the region where the value of mask

Algorithm 1: Constructing point cloud

Input: A single RGBD image $[\mathbf{I}_0, \mathbf{D}_0]$
Input: Camera intrinsic \mathbf{K} , extrinsics $\{\mathbf{P}_i\}_{i=0}^N$
Output: Complete point cloud \mathcal{P}_N

```

1  $\mathcal{P}_0 \leftarrow \phi_{2 \rightarrow 3}([\mathbf{I}_0, \mathbf{D}_0], \mathbf{K}, \mathbf{P}_0)$ 
2 for  $i \leftarrow 1$  to  $N$  do
3    $\hat{\mathbf{I}}_i, \mathbf{M}_i \leftarrow \phi_{3 \rightarrow 2}(\mathcal{P}_{i-1}, \mathbf{K}, \mathbf{P}_i)$ 
4    $\mathbf{I}_i \leftarrow \mathcal{S}(\hat{\mathbf{I}}_i, \mathbf{M}_i), \hat{\mathbf{D}}_i \leftarrow \mathcal{D}(\mathbf{I}_i)$ 
5    $d_i \leftarrow 1$ 
6   while not converged do
7      $\tilde{\mathcal{P}}_i \leftarrow \phi_{2 \rightarrow 3}([\mathbf{I}_i, d_i \hat{\mathbf{D}}_i], \mathbf{K}, \mathbf{P}_i)$ 
8      $\mathcal{L}_d \leftarrow \frac{1}{\|\mathbf{M}_i=1\|} \sum_{\mathbf{M}_i=1} \left\| \tilde{\mathcal{P}}_i - \mathcal{P}_{i-1} \right\|_1$ 
9     Calculate  $\nabla_d \mathcal{L}_d$ 
10     $d_i \leftarrow d_i - \alpha \nabla_d \mathcal{L}_d$ 
11  end
12   $\mathbf{D}_i \leftarrow d_i \hat{\mathbf{D}}_i$ 
13   $\hat{\mathcal{P}}_i \leftarrow \phi_{2 \rightarrow 3}([\mathbf{I}_i, \mathbf{D}_i | \mathbf{M}_i = 0], \mathbf{K}, \mathbf{P}_i)$ 
14   $\mathcal{P}_i \leftarrow \mathcal{P}_{i-1} \cup \mathcal{W}(\hat{\mathcal{P}}_i)$ 
15 end

```

changes ($|\nabla \mathbf{M}_i| > 0$) to find the corresponding points to that region in both \mathcal{P}_{i-1} and $\hat{\mathcal{P}}_i$. Then, we calculate the displacement vector from $\hat{\mathcal{P}}_i$ to \mathcal{P}_{i-1} . However, moving the points in a naive way may distort the shape of the lifted point cloud and make a misalignment between the point cloud and the inpainted image. We mitigate the issue by giving the restrictions for moving the points and using the interpolation algorithm to preserve the overall shape of the points.

First, we force each point in $\hat{\mathcal{P}}_i$ to move along the ray line from the camera center to the corresponding pixel. We find the closest point to the corresponding point in \mathcal{P}_{i-1} along the ray line and report how much the depth changes are caused by the movement. Using the constraint, we preserve the contents of RGB image (\mathbf{I}_i) although moving the points in 3D space. Next, we assume that the depth does not change at the opposite side of the mask boundary region. Then, for the points that do not have their ground truth counterparts, i.e. $\mathbf{M}_i = 0$, we calculate for each pixel how much the depth value should change using linear interpolation. By interpolating smoothly, the mismatch among the pixels caused by the drastic movement is alleviated. The aligned points are combined with the original one:

$$\mathcal{P}_i = \mathcal{P}_{i-1} \cup \mathcal{W}(\hat{\mathcal{P}}_i), \quad (4)$$

where we denote calculating movement and interpolation as \mathcal{W} . We repeat the process N times to construct the final point cloud, \mathcal{P}_N . By reprojection, \mathcal{P}_N provides high-quality

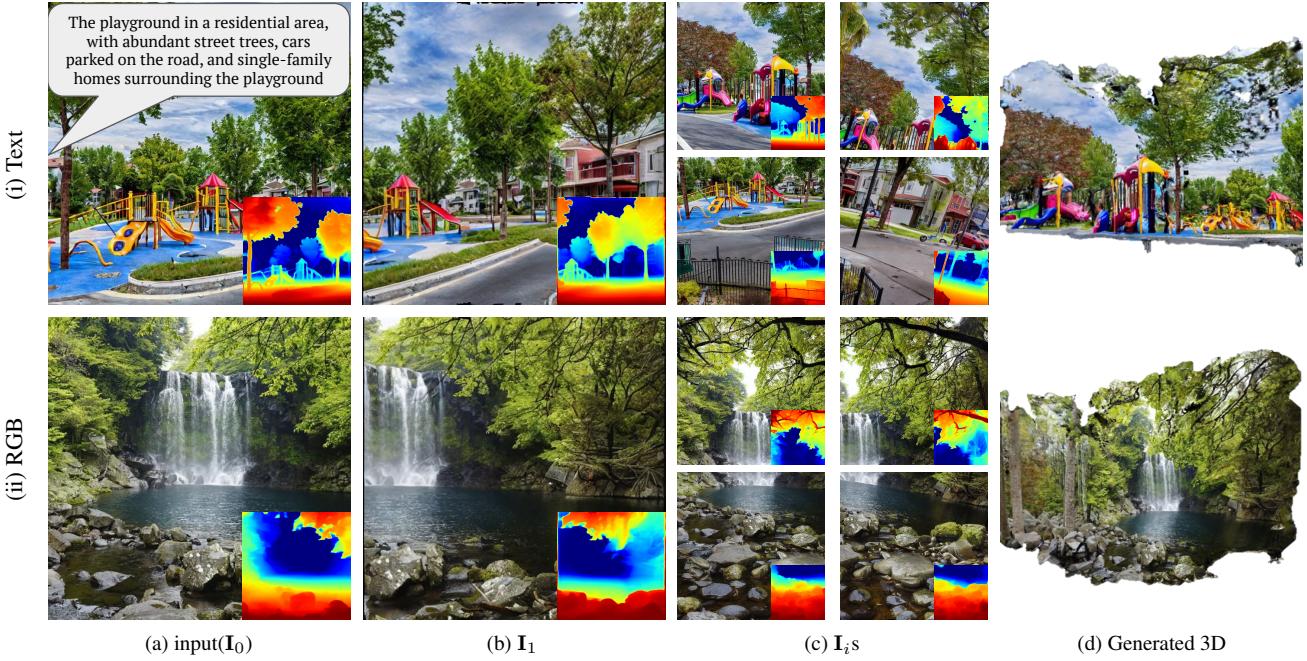


Figure 2. **Intermediate images during point cloud generation and final 3D output between different inputs.** We generate 3D scene from different input types (text and RGB image). The input image in the first row is generated image using Stable diffusion. Our model is capable of generating consistent images high-quality 3D scene regardless of input type.

and multi-view consistent images. The whole process of constructing \mathcal{P}_N from $[\mathbf{I}_0, \mathbf{D}_0]$, \mathbf{K} , and $\{\mathbf{P}_i\}_{i=0}^N$ is written in Algorithm 1.

3.2. Rendering with Gaussian Splatting

After the point cloud is created, we train 3D Gaussian splatting model [21] using the point cloud and the projected images. The center of Gaussian splatting points are initialized by the input point cloud, and the volume and the position of each point are changed by the supervision of input ground truth projected images. We use the generated point cloud (\mathcal{P}_N) as the initial SfM points. Initialization with \mathcal{P}_N will boost up the convergence of the network and encourage the network to focus on generating the details of the representation. For the images to train the model, we use additional M images as well as $(N + 1)$ images for generating the point cloud, since the initial $(N + 1)$ images are not sufficient to train the network for generating the plausible output. The M new images and the masks are generated by reprojecting from the point cloud \mathcal{P}_N by a new camera sequence of length M , denoted as $\mathbf{P}_{N+1}, \dots, \mathbf{P}_{N+M}$.

$$\mathbf{I}_i, \mathbf{M}_i = \phi_{3 \rightarrow 2}(\mathcal{P}_N, \mathbf{K}, \mathbf{P}_i), i = M + 1, \dots, M + N. \quad (5)$$

We note that we do not inpaint \mathbf{I}_i when optimizing Gaussian splats. Instead, when calculating the loss function, we only consider the valid image region where the mask value is 1. It prevents the model from learning the wrong details of the

reprojected images. Since each point is represented as a Gaussian distribution, the missing pixels when training the model are naturally filled, and the rasterized image after the training becomes plausible.

4. Experiments

4.1. Experiment settings

Datasets. Since LucidDreamer is optimized for every input, we do not need training dataset to train the model. For the text input, we randomly generate several text prompts relevant to scene images to generate the first image using Stable Diffusion. We use real or generated high-quality images for the RGB input. For the case of RGBD inputs, we use ScanNet [9] and NYUdepth [45] since the two datasets have ground truth depth maps.

Implementation details. The modules we used to construct LucidDreamer can be either trained using manual design or brought from off-the-shelf models. We use pre-trained large-scale off-the-shelf models to compose the whole network to maximize the generalization capability of the network. Specifically, we adopt Stable Diffusion model [39] to inpaint the masked image. We use the same text prompt input for the Stable Diffusion if the first image is generated from the text. If the input format is a RGB(D) image without text, we use LAVIS [23] to generate the caption according to the image and place it in the diffusion



Figure 3. **Intermediate images during point cloud generation and final 3D output for different text prompt.** We put the different text prompt while having same initial image (I_0) and compare the generation results.

inpainting model to generate consistent content. For the camera trajectory that we use to construct the point cloud ($\{\mathbf{P}_i\}_{i=0}^N$), we create several types of camera trajectory pre-sets in advance, and different types of trajectories were used for different tasks.

4.2. Experiment results

We demonstrate the superiority and high generalizability of LucidDreamer in many aspects. We strongly recommend the readers to watch the video in the supplementary materials where we can entirely show the strength of our model.

Applicability to various input domains and formats. LucidDreamer is capable of generating a consistent and high-quality 3D scene considering the input style. Figure 2 shows the generated realistic images and the 3D scenes. At the top row, we visualize a result of Text-to-3D. We depict an initial image generated from the given text and estimated depth in (a). (b) and (c) present the plausible images and geometry generated through our pipeline involving navigation, dreaming, and alignment. We showcase an overview of the final 3D scene in (d). On the other hand, the bottom row demonstrates an example result of RGB-to-3D. We estimated depth from the given RGB and used it as an initial geometry for the scene. Similar to the top row, we generated

believable images and geometry, resulting in a high-quality 3D scene.

Since our model supports multiple inputs, it can generate 3D scenes in various ways as illustrated in Figure 3. The top and middle rows depict outcomes generated by guaranteeing the inclusion of conditioned RGB during the creation of the 3D scene. Despite different texts, the conditioned RGB is consistently present in the scene. On the other hand, the bottom row displays the outcome of altering the text condition while generating the 3D scene. Through diverse combinations and alterations of conditions, our model facilitates the creation of the desired 3D scene more effortlessly. We illustrate additional example scenes in Figure 5. Our model successfully generates diverse 3D scenes with various styles (e.g. lego, anime) across different camera paths.

Comparison with RGBD2. We qualitatively compare the generation results with RGBD2 [22] and illustrate the result in Figure 4. For fairness of comparison, we compare the results on three images with different domains: generated image, ScanNet, and NYUDepth. For the generated image, the depth map estimated by Zoedepth [2] is considered a ground-truth depth map when processing RGBD2. For ScanNet and NYUDepth, we use the ground truth depth map for both RGBD2 and LucidDreamer when producing a



Figure 4. **Qualitative comparison with RGBD2 [22]** on various image datasets. We compare LucidDreamer with RGBD2 starting from the same input image while changing the datasets. The scene generated by LucidDreamer always shows higher quality than RGBD2, even on ScanNet [9] which RGBD2 is trained on.

Models	CLIP-Score \uparrow [18]	CLIP-IQA [53]		
		Quality \uparrow	Colorful \uparrow	Sharp \uparrow
RGBD2 [22]	0.2035	0.1279	0.2081	0.0126
LucidDreamer	0.2110	0.6161	0.8453	0.5356

Table 1. **Quantitative comparison of generated scenes.** We quantitatively compare the results using CLIP-Score and CLIP-IQA with RGBD2. Our model shows better results on all metrics.



(a) Generated video

(b) 3D whole view

Figure 5. **3D reconstruction results and short video on various styles.** This is a video figure that is best viewed by Adobe Reader.

3D scene. For ScanNet, each scene consists of several images and the corresponding depth maps and camera views. We randomly select one of the given image and depth map pairs and use it as an initial RGBD input. In Figure 4b, we observe that RGBD2 generates (ScanNet-style) images with similar styles regardless of the input image. This remains consistent not only in the initial image but also throughout the following sequence as shown in Figure 4c. We believe the issue arises due to insufficient training data and domain limitations, highlighting the need for a model with sufficient generalization. In contrast, our approach generates high-quality 3D scenes with careful consideration to harmonize well with the input RGB. Moreover, LucidDreamer can generate scenes composed of high-resolution images while RGBD2 can only make 128×128 -sized images, which is too small to use in real applications. We also document the quantitative results evaluated on CLIP-Score [18] and CLIP-IQA [53] in Table 1. We confirm that our model incorporates input conditions well, resulting in the creation of high-quality 3D scenes.

Iters	Source of SfM points	Metrics		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
1000	COLMAP	23.15	0.7246	0.2910
	LucidDreamer	32.59	0.9672	0.0272
3000	COLMAP	30.87	0.9478	0.0353
	LucidDreamer	33.80	0.9754	0.0178
7000	COLMAP	32.52	0.9687	0.0208
	LucidDreamer	34.24	0.9781	0.0164

Table 2. **Reconstruction quality according to the source of initial SfM points.** We use the initial point cloud generated by COLMAP [40, 41] and compare the reconstruction results. Our model consistently shows better reconstruction metrics.

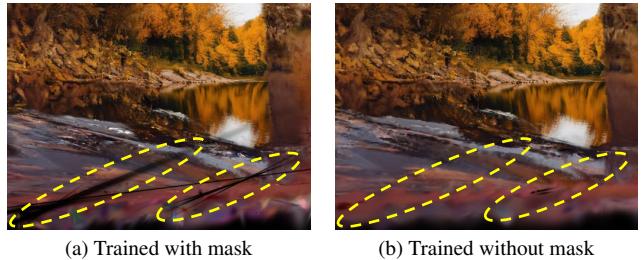


Figure 6. **The effect of the mask during training.** Training with valid masks helps prevent artifacts at the boundaries of the scene.

Ablations on design choices. We demonstrate ablation studies on the design choices of LucidDreamer. In table 2, we compare the effect of COLMAP initialization and the point cloud initialization provided by our method when learning Gaussian splatting. LucidDreamer achieved a high-quality scene within fewer iterations by offering a substantial amount of high-quality initialization points. Figure 6 illustrates the effectiveness of masking the image while learning Gaussian splatting. It removes the black splinter from the black background, resulting in visually pleasing outcomes.

5. Conclusion

In this work, we propose LucidDreamer, a novel pipeline for domain-free 3D scene generation. By fully exploiting the power of large diffusion models, LucidDreamer is capable of generating high-quality scenes without the restriction of the target scene domain. We first generate the point cloud from the input image and repeat ‘Dreaming’ and ‘Alignment’ algorithms to generate the multi-view consistent high-quality image and harmoniously integrate them to the existing point cloud in the 3D space. After the construction is finished, the point cloud is converted to 3D Gaussian splats to enhance the quality of the 3D scene. Extensive experiments show that LucidDreamer can consistently generate high-quality and diverse 3D scenes in various situations.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*, 2018. 3
- [2] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 3, 6
- [3] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021. 3
- [4] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, 2022. 3
- [5] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorrf: Tensorial radiance fields. In *ECCV*, 2022. 3
- [6] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. *arXiv preprint arXiv:2304.06714*, 2023. 3
- [7] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *CVPR*, 2023. 3
- [8] Dana Cohen-Bar, Elad Richardson, Gal Metzer, Raja Giryes, and Daniel Cohen-Or. Set-the-scene: Global-local training for generating controllable nerf scenes. *arXiv preprint arXiv:2303.13450*, 2023. 3
- [9] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 2, 5, 7
- [10] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. Cvxnet: Learnable convex decomposition. In *CVPR*, 2020. 2
- [11] Emilien Dupont, Hyunjik Kim, SM Eslami, Danilo Rezende, and Dan Rosenbaum. From data to functa: Your data point is a function and you can treat it like one. *arXiv preprint arXiv:2201.12204*, 2022. 2, 3
- [12] Rina Foygel and Mathias Drton. Extended bayesian information criteria for gaussian graphical models. 2010. 2
- [13] Rafail Fridman, Amit Abecasis, Yoni Kasten, and Tali Dekel. Scenescape: Text-driven consistent scene generation. *arXiv preprint arXiv:2302.01133*, 2023. 3
- [14] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qin-hong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 3
- [15] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 3
- [16] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *ICCV*, 2019. 2
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NIPS*, 2014. 3
- [18] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: a reference-free evaluation metric for image captioning. In *EMNLP*, 2021. 8
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 3
- [20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM ToG*, 2023. 2, 3
- [21] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM ToG*, 2023. 5
- [22] Jiabao Lei, Jiapeng Tang, and Kui Jia. Rgbd2: Generative scene synthesis via incremental view inpainting using rgbd diffusion models. In *CVPR*, 2023. 3, 6, 7, 8
- [23] Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio Savarese, and Steven C.H. Hoi. LAVIS: A one-stop library for language-vision intelligence. In *ACL*, 2023. 5
- [24] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NIPS*, 2020. 3
- [25] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *CVPR*, 2021. 2, 3
- [26] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 2021. 2
- [27] Thomas Müller, Fabrice Rousselle, Jan Novák, and Alexander Keller. Real-time neural radiance caching for path tracing. *arXiv preprint arXiv:2106.12372*, 2021. 3
- [28] Thomas Müller, Alex Evans, Christoph Schied, and

- Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *TOG*, 2022. 3
- [29] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *ICCV*, 2019. 3
- [30] Thu H Nguyen-Phuoc, Christian Richardt, Long Mai, Yongliang Yang, and Niloy Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. *NeurIPS*, 2020. 3
- [31] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, pages 11453–11464, 2021. 3
- [32] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 2
- [33] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In *CVPR*, 2019. 2
- [34] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3
- [35] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 3
- [36] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 3
- [37] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 3
- [38] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021. 3
- [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 3, 5
- [40] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 8
- [41] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 8
- [42] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *NIPS*, 2020. 3
- [43] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3d point cloud generative adversarial network based on tree structured graph convolutions. In *ICCV*, 2019. 3
- [44] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *CVPR*, 2023. 3
- [45] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012. 2, 5, 7
- [46] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 2020. 2
- [47] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3
- [48] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 3
- [49] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *CVPR*, 2021. 2
- [50] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. 3
- [51] Shitao Tang, Fuyang Zhang, Jiacheng Chen, Peng Wang, and Yasutaka Furukawa. Mvdiffusion: Enabling holistic multi-view image generation with correspondence-aware diffusion. *arXiv preprint arXiv:2307.01097*, 2023. 2
- [52] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *CVPR*, 2017. 2
- [53] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *AAAI*, 2023. 8
- [54] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *NIPS*, 2016. 3
- [55] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *CVPR*, 2022. 3

- [56] Mohsen Yavartanoo, Jaeyoung Chung, Reyhaneh Nezhavar, and Kyoung Mu Lee. 3dias: 3d shape reconstruction with implicit algebraic surfaces. In *ICCV*, 2021. [2](#)
- [57] Tackgeun You, Mijeong Kim, Jungtaek Kim, and Bohyung Han. Generative neural fields by mixtures of neural implicit functions. *arXiv preprint arXiv:2310.19464*, 2023. [3](#), [4](#)
- [58] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. [3](#)
- [59] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. *arXiv preprint arXiv:2210.06978*, 2022. [3](#)
- [60] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *ICCV*, 2021. [2](#), [3](#)