

SAGD: Boundary-Enhanced Segment Anything in 3D Gaussian via Gaussian Decomposition

Xu Hu^{1,2}, Yuxi Wang^{2,3}, Lue Fan^{3,4}, Chuanchen Luo⁶, Junsong Fan^{2,3}, Zhen Lei^{2,3,4},

Qing Li^{1†}, Junran Peng^{5†}, and Zhaoxiang Zhang^{2,3,4†}

Abstract—3D Gaussian Splatting has emerged as an alternative 3D representation for novel view synthesis, benefiting from its high-quality rendering results and real-time rendering speed. However, the 3D Gaussians learned by 3D-GS have ambiguous structures without any geometry constraints. This inherent issue in 3D-GS leads to a rough boundary when segmenting individual objects. To remedy these problems, we propose SAGD, a conceptually simple yet effective boundary-enhanced segmentation pipeline for 3D-GS to improve segmentation accuracy while preserving segmentation speed. Specifically, we introduce a Gaussian Decomposition scheme, which ingeniously utilizes the special structure of 3D Gaussian, finds out, and then decomposes the boundary Gaussians. Moreover, to achieve fast interactive 3D segmentation, we introduce a novel training-free pipeline by lifting a 2D foundation model to 3D-GS. Extensive experiments demonstrate that our approach achieves high-quality 3D segmentation without rough boundary issues, which can be easily applied to other scene editing tasks. Our code is publicly available at <https://github.com/XuHu0529/SAGS>.

Index Terms—3D gaussian splatting, 3D segmentation, boundary issues.

I. INTRODUCTION

3D scene understanding is a challenging and crucial task in computer vision and computer graphics, which involves scene reconstruction from images or videos and the perception of a given 3D real-world environment. Researchers have conducted extensive studies in scene reconstruction and 3D scene perception in recent years. For instance, Neural Radiance Fields (NeRF) [10], [33], [56], [78] have significantly contributed to the progress of 3D scene reconstruction by representing scenes in an implicit way. In the field of scene perception, continuous research is being conducted on 3D detection and semantic segmentation, based on the representation of range images [23], [55], point clouds [22], [29], [39], [63], [76], and Bird’s Eye View (BEV) [42], [61], [62], [69], [91]. Although current methods have attained noteworthy success in 3D scene understanding, the time-consuming nature of NeRF and the high costs associated with 3D data collection pose challenges for scaling up these approaches.

Recently, 3D Gaussian Splatting (3D-GS) [35] is emerging as a prospective method for modeling static 3D scenes. 3D-

GS characterizes intricate scenes by employing numerous colored 3D Gaussians, rendering them into camera views through splatting-based rasterization. Through differentiable rendering and gradient-based optimization, the positions, sizes, rotations, colors, and opacities of these Gaussians can be finely tuned to accurately represent the 3D scene, availing for the comprehension of a 3D environment. Segmentation in 3D-GS is quite a natural pathway to scene understanding for its explicit representation. Recent works [6], [40], [92], [110] have been proposed to achieve segmentation in 3D-GS via lifting 2D foundation models in a learnable way. They share the same ideas to distill other features or identity encodings generated from 2D foundation models to 3D Gaussian fields.

However, the Gaussians learned by 3D-GS are ambiguous without any geometry constraint, leading to the issue of boundary roughness. A single Gaussian might correspond to multiple objects, complicating the task of accurately segmenting individual objects (shown in Fig. 1 (a)). As shown in Fig. 1 (b), direct segmentation will leave rough edges at the boundary. This is because these Gaussians across multiple objects will be left if no additional processing is performed. In addition, even if SAGA [6] uses filtering and growing post-processing for segmented 3D Gaussians, its feature matching method also has the limitation of incomplete boundary, as shown in Fig. 1 (c).

We present our boundary-enhanced segmentation method, an interactive training-free pipeline for efficient and effective segmentation in 3D-GS without rough boundary issues. By leveraging the 2D foundation model SAM [38], our method can generate a segmented mask from a single input view based on the given prompts. Starting from the obtained mask, our method automatically generates multi-view masks and achieves consistent 3D segmentation via the proposed assignment strategy. Specifically, to resolve the inherent boundary issues, we incorporate a simple but effective Gaussian Decomposition scheme, ingeniously utilizing the special structure of 3D Gaussian. Since the boundary roughness issue results from the non-negligible spatial sizes of 3D Gaussians located at the boundary, our key insight is to find and then decompose the original boundary Gaussians according to our proposed principles. Consequently, our method eliminates the inherent issues and achieves more complete segmentation quickly and efficiently, as shown in Fig. 1 (d).

In summary, the contributions of this paper are as follows:

- We propose a simple yet effective training-free pipeline for segmentation in 3D Gaussians without any learnable

[†] Corresponding author.

¹ The Hong Kong Polytechnic University

² Center for Artificial Intelligence and Robotics, HKISI, CAS

³ Institute of Automation, Chinese Academy of Sciences

⁴ University of Chinese Academy of Sciences

⁵ University of Science and Technology Beijing

⁶ Shandong University

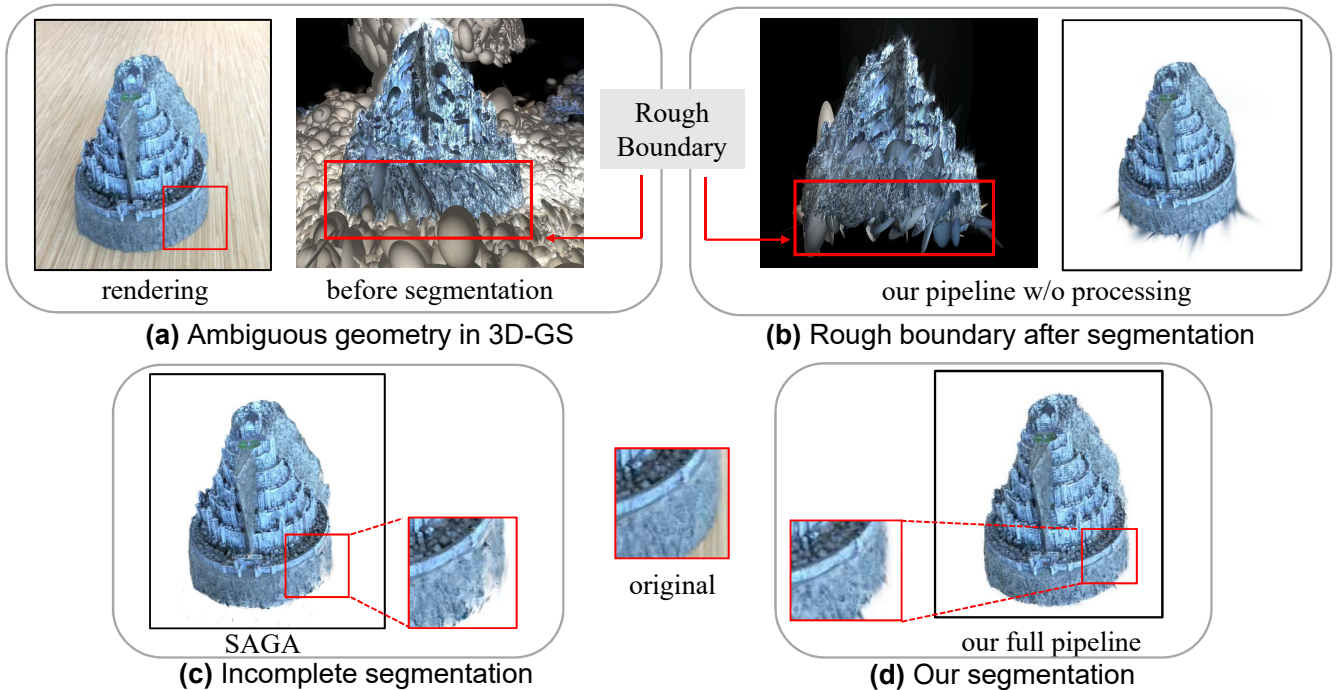


Fig. 1. (a)The training of 3D-GS doesn’t consider the structure of objects, leading to the ambiguous geometry; (b)Direct segmentation without Gaussian Decomposition processing will result in rough boundary segmentation; (3) The recent SAGA also has incomplete segmentation caused by the same issue; (d) Our full pipeline considers this issue and can achieve better segmentation.

parameters;

- We incorporate the Gaussian Decomposition module to mitigate the boundary roughness issues in 3D segmentation resulting from inherent 3D-GS geometry;
- Extensive segmentation experiments on a considerable amount of 3D scenes and editing applications demonstrate the effectiveness of our proposed method.

II. RELATED WORK

A. Radiance Fields and 3D Gaussian Splatting

Novel view synthesis (NVS) involves rendering unseen viewpoints of a scene from a given set of images. One popular approach is Neural Radiance Fields (NeRF), which uses a Multilayer Perceptron (MLP) to represent 3D scenes and map from a 3D coordinate to properties of the scene at the corresponding location. It leverages the differentiable volume rendering technique to translate a 3D scene’s continuous representation into 2D images. Several works have been proposed to enhance NeRF’s performance by addressing aspects such as speed [10], [24], [25], [59], [60], [72], [78], [79], [94] and adapting it to other tasks [74], [77], [85], [105]. Though some breakthroughs have been achieved, the reliance on low-efficient volume rendering still hinders real-time rendering.

Recently, 3D Gaussian Splatting [35] has been proposed as a new technique for novel view synthesis. It leverages an ensemble of anisotropic 3D Gaussian splats to represent the scene and employs differentiable splatting for rendering. 3D Gaussian Splatting has been shown to be an alternative 3D representation of NeRF, benefiting from both its high-quality rendering results and real-time rendering speed. Recent research on this technique involves applying 3D Gaussian Splatting to large-scale scene reconstruction [49], [50], the

dynamic scenes [1], [31], [41], [45], [46], [51], [52], [75], [80], [87], [95], [98], [112], 3D object and scene generation [16], [43], [81], [89], [90], [93], [103], [109], [111], surface reconstruction [20], [27], [30], [97] and other tasks [9], [44], [75], [86].

B. 2D and 3D Perception

Detection and segmentation are fundamental computer vision tasks. Numerous studies [5], [8], [64]–[68], [101], [102], [106] have deeply explored various sub-fields. Especially, Grounding-DINO [47] enhances the grounding results by introducing the language instructions. Traditionally, the segmentation includes three major tasks: semantic [12]–[15], [84], instance [4], [28], [48], and panoptic [37] segmentation. Due to the operation consistency of the above tasks at the pixel level, many studies have tried to use a unified framework, such as K-net [104], MaskFormer [18], and Mask2Former [17]. A significant breakthrough in 2D segmentation is the Segment Anything Model (SAM) [38]. SAM seeks to unify the 2D segmentation task using a prompt-based segmentation approach. Its introduction has sparked a new wave of research, with many studies focused on enhancing its functionality. These improvements include efficient fine-tuning techniques [34], [70], [71] and distillation-based acceleration methods [99], [107]. Additionally, SAM has been adapted for various fields, such as medical image analysis [21], [53], [54], [88], concealed object detection [32], [82], image editing [96], remote sensing [11], and 3D bird’s eye view (BEV) sensing [100].

C. Segmentation in Radiance Fields

Neural Radiance Fields (NeRFs) are a popular way of implicitly representing 3D scenes with neural networks. Many

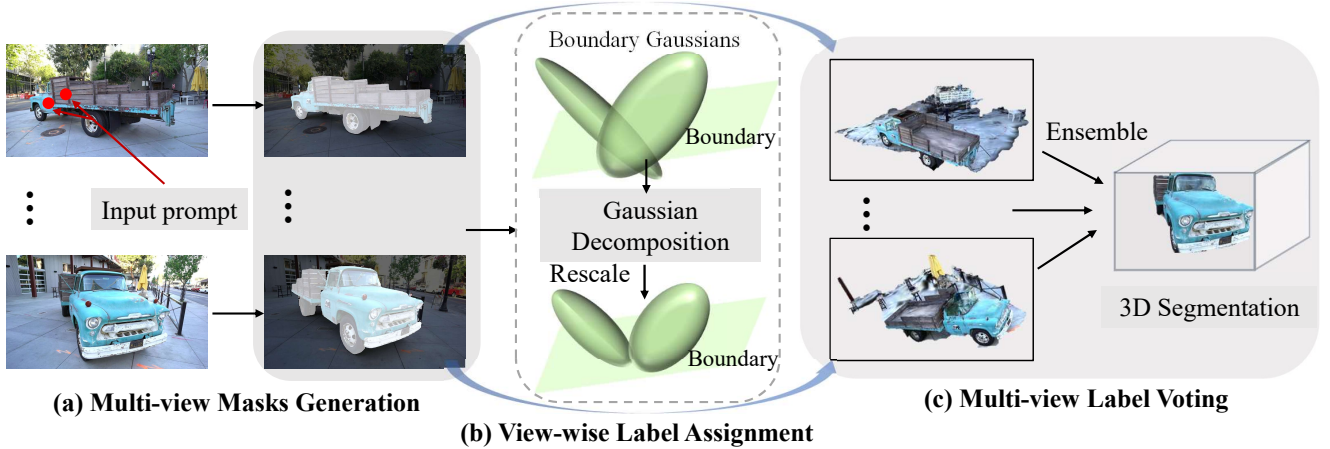


Fig. 2. Pipeline of our proposed method. (a) Given a set of clicked points on the 1st reference view, we utilize SAM to generate masks for corresponding objects under every view automatically; (b) For every view, Gaussian Decomposition is first performed to address the issue of boundary roughness and then label propagation is implemented to assign binary labels to each 3D Gaussian; (c) Finally, with assigned 3D labels from all views, we adopt a simple yet effective voting strategy to determine the segmented Gaussians.

researchers have explored how to segment objects in 3D using NeRFs for various applications such as novel view synthesis, semantic segmentation, 3D inpainting, and language grounding. Some methods, such as Semantic-NeRF [108], NVOS [73], and SA3D [7], use different types of inputs to guide the segmentation, such as semantic labels, user scribbles, or 2D masks. Other methods, such as N3F [83], DFF [83], LERF [36], and ISRF [26], learn additional feature fields that are aligned with NeRFs, and use 2D visual features from pre-trained models or language embeddings to query the 3D features. These methods usually require modifying or retraining the original NeRF models or training other specific parameters to obtain the 3D segmentation. However, limited by the representation and rendering speed of NeRF, it is challenging to apply it to more practical applications, such as scene editing, collision analysis, etc. Inspired by the real-time 3D-GS, recent works [6], [40], [92], [110] have been proposed to achieve segmentation in 3D-GS in a learnable way. SAGA [6] and Featree 3DGS [110] distill the knowledge embedded in the SAM decoder into the feature field of 3D-GS, and Gaussian-Grouping [92] supervises the Identity Encodings during the differentiable rendering by leveraging the 2D mask predictions by SAM. However, no research has been proposed to solve the rough boundary issues in segmentation resulting from the inherent properties of 3D Gaussians. In this work, we are the first to propose a novel Gaussian Decomposition scheme to solve this issue, incorporated into our effective training-free segmentation pipeline.

III. METHOD

In this section, we first present the preliminary in 3D Gaussian Splatting and Segment Anything Model (SAM) in Sec. III-A for clear understanding and then define the problem and task in Sec. III-B. Our basic training-free pipeline consists of Sec. III-C and Sec. III-E. The details of Gaussian Decomposition are described in Sec. III-D. The pipeline of our method can be seen in Fig. 2.

A. Preliminary

a) 3D Gaussian Splatting: 3D Gaussian Splatting (3D GS) [35] is an emerging method for real-time radiance field rendering. It has been proven effective in Novel View Synthesis with high rendering quality as NeRF and real-time rendering speed. 3D GS represents scenes with a set of 3D Gaussians. Specifically, each 3D Gaussian is parameterized by a position $\mu \in \mathbb{R}^3$, a covariance matrix Σ consisting of a scaling factor $s \in \mathbb{R}^3$, and a rotation quaternion $q \in \mathbb{R}^4$, an opacity value α , and spherical harmonics (SH). Each 3D Gaussian is characterized by:

$$G(x) = \frac{1}{(2\pi)^{3/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (1)$$

To render an image, it uses the splatting rendering pipeline, where 3D Gaussians are projected onto the 2D image plane. The projection transforms 3D Gaussians into 2D Gaussians in the image plane. All 2D Gaussians are blended together by the α -blending algorithm to generate the color:

$$\mathbf{c} = \sum_{i \in N} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2)$$

During the α -blending process, for each 2D Gaussian, only the 2D points with probability density larger than a certain threshold are calculated. This means a 2D Gaussian and 3D Gaussian can be intuitively regarded as a 2D **ellipse** and a 3D **ellipsoid** respectively. Empirically, for an axis of the ellipse, its length is set to 3σ , where σ is the square root of the variance in the axis.

b) Segment Anything Model (SAM): SAM [38] takes an image I and a set of prompts P as input to output the corresponding segmentation mask M_{SAM} :

$$M_{SAM} = SAM(I, P) \quad (3)$$

where the prompts P can be points, bounding boxes, or texts.

B. Problem Definition

Given a set of trained 3D Gaussians $\mathbb{G} = \{\mathbf{g}_0, \mathbf{g}_1, \dots, \mathbf{g}_n\}$ and a random initial view \mathbf{v}_0 , users could offer 2D point prompt set $\mathbb{P}_{2D} = \{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_m\}$ to specify a 2D object in view \mathbf{v}_0 . Our algorithm is supposed to segment the corresponding 3D object \mathbb{O} in \mathbb{G} according to the human prompts, where \mathbb{O} is a subset of \mathbb{G} .

Let \mathbf{m}_i denote the projected binary mask of \mathbb{O} in i -th view, an accurate segmentation \mathbb{O} means \mathbf{m}_i equals to \mathbf{m}_i^* for $\forall i \in \{0, 1, \dots, n\}$. Here \mathbf{m}_i^* is the ground truth mask of \mathbb{O} in i -th view. Different from the conventional 2D segmentation mask in the image or 3D segmentation task in the point cloud, there is no ground truth for 3D Gaussians. Thus, our algorithm is designed to minimize the difference between \mathbf{m}_i and \mathbf{m}_i^* .

C. Segment 3D Gaussians with 2D Mask

3D Prompts for Multiview Masks Generation. By the definition in Sec. III-B, users are given the first rendered view to specify the target object. However, only the first view is far from sufficient to segment the target object in 3D space. So here we first generate multiview masks to aid the 3D segmentation. With multiple masks from different views, a 3D object can be segmented by the intersection of the corresponding frustum of these masks.

The core of obtaining masks is to generate 2D prompt points in each view. Denoting the i -th 2D prompt point in the first given view \mathbf{v}_0 as \mathbf{p}_i^0 , we define the corresponding 3D prompt \mathbf{p}_i^{3D} as:

$$\arg \min_{\mu} \{d(\mu), d(\mu) > 0 \mid \mu \in \mathbb{G}, \|\mathbf{P}_0 \mu - \mathbf{p}_i^0\|_1 < \epsilon\}, \quad (4)$$

where $d(\mu)$ is the depth of Gaussian center μ and \mathbf{P}_0 is the projection for initial view \mathbf{v}_0 . Thus $\mathbf{P}_0 \mu$ is the position of μ in view \mathbf{v}_0 . Eq. 4 indicates that the corresponding 3D prompt of \mathbf{p}_i is the center of a certain 3D Gaussian. This center meets two requirements: (1) it has a similar projected position with \mathbf{p}_i with a Manhattan distance less than ϵ , and (2) if there are multiple 3D Gaussian centers satisfying the first requirement, the one with the smallest positive depth is selected as the 3D prompt.

For all the 2D prompt points in the first view, we could get a set of 3D prompts by Eq. 4. Then for another view \mathbf{v}_i , we project these 3D prompts into the 2D plane, resulting in 2D prompts in the view \mathbf{v}_i . In this way, we obtain 2D prompts in all views, and all masks are obtained by prompting the SAM. **View-wise Label Assignment.** With all the masks, we proceed to assign binary labels to each 3D Gaussian. In particular, we maintain a matrix \mathbf{L} , whose element \mathbf{L}_{ij} is defined by

$$\mathbf{L}_{ij} = \begin{cases} 1 & \text{if } \mathbf{P}_j \mu_i \in \mathbf{m}_j, \\ 0 & \text{if } \mathbf{P}_j \mu_i \notin \mathbf{m}_j, \end{cases} \quad (5)$$

where μ_i is the i -th Gaussian center of the scene and \mathbf{m}_j is the foreground mask in j -th view. \mathbf{P}_j stands for the projection matrix of j -th view.

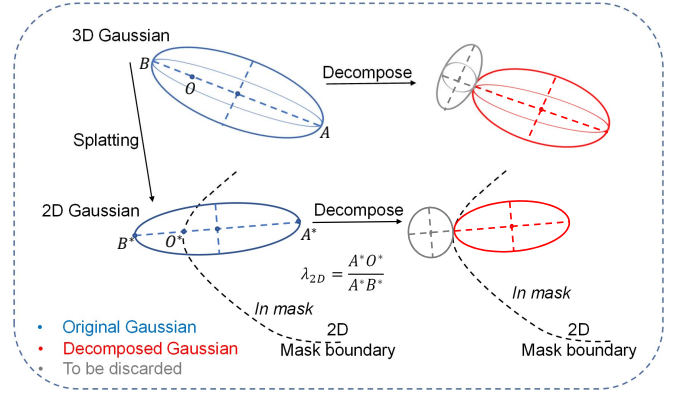


Fig. 3. Illustration of the Gaussian Decomposition process. It involves two basic steps: first, to find out the boundary Gaussians and then decompose these Gaussians.

D. Gaussian Decomposition

After obtaining the prompt points of each view, we could assign labels to 3D Gaussians by the projected position of its center. However, 3D Gaussian has non-negligible spatial volume, so those Gaussians projected to the mask boundary usually have a part out of the boundary, greatly increasing the roughness of boundaries. A straightforward solution is directly removing the Gaussians across mask boundaries. However, such a solution greatly damages the 3D structures of the object.

To address this issue, we propose the *Gaussian Decomposition* to mitigate the boundary roughness while maintaining the 3D structures as complete as possible. Fig. 3 illustrates the basic idea of Gaussian decomposition. It has two basic steps to achieve the Gaussian decomposition:

(1) For each 3D Gaussian, we obtain the corresponding 2D Gaussian by projection and mark it as a boundary Gaussian if one of its long-axis endpoints is outside of the 2D mask. Then, as shown in Fig. 3, we denote the two ends of the long axis of a 2D Gaussian as A^* and B^* , and the intersection of A^*B^* and mask boundary is O^* . A , B , and O are the corresponding 3D points in the long axis of the 3D Gaussian. Assuming A^* is in the mask while B^* is out of the mask, we define

$$\lambda_{3D} = \frac{OA}{AB}, \quad (6)$$

$$\lambda_{2D} = \frac{O^*A^*}{A^*B^*}. \quad (7)$$

(2) After obtaining the boundary Gaussians, we solve the ratio λ_{3D} by first calculating λ_{2D} and then decompose the original 3D Gaussian according to λ_{3D} .

However, the transformation from λ_{2D} to λ_{3D} is not straightforward because the perspective projection from 3D space to 2D space is not affine, which means the ratios are different. Fortunately, 3D Gaussian Splatting leverages local affine approximation to simplify the rendering process. It projects 3D Gaussian to the 2D plane by

$$\Sigma' = \mathbf{J} \mathbf{W} \Sigma \mathbf{W}^T \mathbf{J}^T \quad (8)$$

where Σ , Σ' are the covariance matrix of the 3D and 2D gaussian distribution, respectively. \mathbf{W} is the projection transformation and \mathbf{J} is the Jacobian of affine approximation derived in EWA algorithm [113]. Eq. 8 indicates that Gaussian

Splatting simplifies the perspective projection to an affine projection, thus decomposition scaling ratio λ_{3D} in 3D space is equivalent to the ratio λ_{2D} in the 2D plane.

Let \mathbf{g} denote a 3D Gaussian across the boundary. Its scale in the long axis and the 3D center are defined as s and μ , respectively. We have

$$s' = \lambda_{2D}s, \quad (9)$$

$$\mu' = \mu + \frac{1}{2}(s - \lambda_{2D}s)\mathbf{e}, \quad (10)$$

where \mathbf{e} is the unit vector pointing from the 3D Gaussian center to the in-mask endpoint of the long axis. The decomposed Gaussian \mathbf{g}' adopt μ' and s' as the new center and long-axis scale, maintaining other properties unchanged. Another decomposed Gaussian outside of the mask is removed.

E. Multiview Label Voting

So far, every 3D Gaussian including the decomposed one has a list of binary labels \mathbf{L}_i , using the label assignment in Sec. III-C. Leveraging the assigned labels, here we adopt a simple yet effective heuristic rule to determine if a 3D Gaussian \mathbf{g}_i belongs to the target 3D object. In particular, we first define the confidence score s_i of \mathbf{g}_i as

$$s_i = \frac{1}{N} \sum_{j=0}^{N-1} L_{ij}, \quad (11)$$

where N is the number of views. First, $s_i > 0.5$ is performed to vote out the object 3D Gaussian. Then, to reduce the background bias due to occlusion, we adopt a threshold τ and an Object-ID \mathbf{O}_i for each 3D Gaussian \mathbf{g}_i is determined by:

$$\mathbf{O}_i = \begin{cases} 1 & \text{if } s_i > \tau, \\ 0 & \text{if } s_i < \tau, \end{cases} \quad (12)$$

F. Multi-object Segmentation

There are numerous objects in the 3D scenes. Our method can also support segmenting all objects simultaneously by merely extending the binary masks to multi-label masks. Specifically, the element \mathbf{L}_{ij} in the matrix \mathbf{L} is extended from binary values to multiple values, where $\mathbf{L}_{ij} \in \{0, 1, \dots, C-1\}$ and C is the number of objects. The similar multi-view label voting strategy will determine an Object-ID \mathbf{O}_i for each 3D Gaussian \mathbf{g}_i :

$$\mathbf{O}_i = \text{Mode}(\mathbf{L}_{ij}), j \in \{0, 1, \dots, N-1\}, \quad (13)$$

where $\text{Mode}(\cdot)$ is the function of finding the element that appears the most, and N is the number of views.

IV. EXPERIMENTS

A. Datasets

We choose different datasets to testify our method, including LLFF [57], Mip-NeRF 360 [3], LERF [36], and some test scenes from the 3D Gaussian Splatting [35]. These datasets contain both small indoor objects and large outdoor scenes, which are very complex and challenging. For quantitative experiments, because there is no existing benchmark that can be used in 3D Gaussian space, we use the SPIn-NeRF [58] dataset with 2D ground truth for evaluation.

B. Implementation Details

As the implementation of 3D Gaussian Splatting [35] for each scene, we follow the official code with default parameters, and each scene is trained with 30000 iterations. For generating single object segmentation masks, we don't limit the number of clicked points on a single reference view for SAM. Binary masks of other views can be automatically generated via the proposed method. This is reasonable since users can refine their input prompts to help SAM generate a 2D mask as accurately as possible from the reference view. For scene segmentation to generate multi-label masks, we first employ SAM to produce instance masks for individual views under automatic segmentation mode. To ensure 2D mask consistency across views, we use a pre-trained zero-shot tracker [19], [92] to propagate and associate masks. For the selection of hyper-parameters, we only have one to control, which is the confidence score threshold τ in the Multiview Label Voting method III-E. By default, we set the value to 0.7 in all experiments. Also, in practical use, users can manually set this value according to the complexity of different scenes. As for the number of views used in experiments, we follow SA3D [7] to select all view images of each scene to finish our segmentation process.

C. Quantitative results

Point-Guided Segmentation: We first conduct experiments on the SPIn-NeRF [58] dataset for quantitative analysis. Given a set of images of the scene, we follow the process described in Section III to obtain the segmentation of the target object in 3D Gaussian Space. The segmented 3D Gaussians are used to render 2D masks in other views. Finally, we calculate the IoU and Accuracy between these rendered and the ground-truth masks. Results can be seen in Table I.

In the comparison, it is noteworthy that both the MVSeg [58] and the SA3D [7] require additional parameters and a computation-costly training process. While SAGA [6] shows nearly one-thousandth of the time compared with them, it needs extra training time (10 minutes per scene) to distill the knowledge embedded in the SAM decoder into the feature field. By contrast, the ‘‘Single view’’ [7] refers to mapping the 2D masks to the 3D space based on the corresponding depth information, which does not need additional training process. In this sense, our method is the same as the ‘‘Single view’’ that does not incorporate any additional training or model parameters. Results show that our approach can achieve comparable segmentation quality to the MVSeg and SA3D methods and better results than SAGA. The time cost of ours is much less than MVSeg and SA3D. Considering the training time of SAGA, the average segmentation cost of a single object is basically the same as our method. In some 360 outdoor scenes, such as ‘‘Truck’’, our approach even outperforms the SA3D and the MVSeg. When compared with the training-free method ‘‘Single view’’, our approach achieves a significant promotion of +15.9% IoU and +3.5% Acc. These results demonstrate our approach is very efficient in obtaining high-quality segmentation masks.

TABLE I
QUANTITATIVE RESULTS ON SPIN-NeRF DATASET. *Single view* DENOTES PROJECTING THE 2D SEGMENTATION RESULT TO 3D SIMPLY.

Scenes	Single View [7]		MVSeg [58]		SA3D [7]		SAGA [6]		Ours	
	IoU	Acc	IoU	Acc	IoU	Acc	IoU	Acc	IoU	Acc
Orchids	79.4	96.0	92.7	98.8	83.6	96.9	-	-	85.4	97.5
Ferns	95.2	99.3	94.3	99.2	97.1	99.6	-	-	92.0	98.9
Room	73.4	96.5	95.6	99.4	88.2	98.3	-	-	86.5	98.1
Horns	85.3	97.1	92.8	98.7	94.5	99.0	-	-	91.1	98.4
Fortress	94.1	99.1	97.7	99.7	98.3	99.8	-	-	95.6	99.5
Fork	69.4	98.5	87.9	99.5	89.4	99.6	-	-	83.4	99.3
Pinecone	57.0	92.5	93.4	99.2	92.9	99.1	-	-	92.6	99.0
Truck	37.9	77.9	85.2	95.1	90.8	96.7	-	-	93.0	97.9
Lego	76.0	99.1	74.9	99.2	92.2	99.8	-	-	90.2	99.7
Mean	74.1	95.2	90.5	98.8	91.9	98.8	88.0	98.5	90.0	98.7

TABLE II
QUANTITATIVE RESULTS WITH TEXT PROMPTS ON FOUR SCENES OF SPIN-NeRF DATASET.

Scene	room	truck	fortress	pinecone
Text	“the table”	“the truck”	“the fortress”	“the pinecone”
IoU	86.3	93.7	92.8	86.5
Acc	97.9	97.8	98.8	97.0

Text-Guided Segmentation: We replace the point prompts with text prompts corresponding to the objects. Given a text input, we prompt the Grounding DINO [47] to generate target bounding boxes, which serve as input prompts for the SAM to obtain 2D segmentation masks. Then, following the segmentation process in Section III, we can transform the 2D masks into 3D Gaussian masks corresponding to the text input. Similar to the above evaluation process, we also calculate the IoU and Acc between 2D rendered masks and given ground-truth masks. Quantitative results are shown in Table II. For both room and truck scenes, the IoU and Acc values can achieve a similar level as using clicked points as input prompts. These results demonstrate our method can combine multi-modal prompts as input.

D. Boundary Segmentation Analysis

The main contribution of our work is to address the boundary issues of segmentation in 3DGS, as mentioned in Sec. I and Fig. 1. To validate the effectiveness of our work, we evaluate metrics, including IoU, AP, and F1-score, on region nearby boundaries. Specifically, we extract parts of 3-pixel width alongside boundaries of ground-truth masks and compare results with previous work SAGA [6]. Experiments are conducted on Scene *pinecone* and *fortress* from the LLFF [57] dataset. Better results of ours are indicated in Tab. III. Benefiting from our proposed Gaussian Decomposition for ambiguous geometry around boundaries, we can achieve more complete and accurate boundary segmentation results.

E. Qualitative results

We conduct four kinds of tasks to demonstrate the potential application and qualitative performance of our approach, including point-guided segmentation, text-guided segmentation, scene editing, and collision detection.

TABLE III
EVALUATION ON BOUNDARY MASKS. WE EXTRACT PARTS OF 3-PIXEL WIDTH ALONGSIDE BOUNDARIES OF GROUND TRUTH AND SEGMENTATION MASKS AND EVALUATE ON **IOU**, **AP**, AND **F1-SCORE** METRICS.

Scene	IoU		AP		F1-score	
	SAGA [6]	Ours	SAGA	Ours	SAGA	Ours
pinecone	13.8	21.8	15.7	23.5	24.0	35.6
fortress	16.7	23.3	31.8	60.2	28.5	37.2

Point-Guided Segmentation: We first conduct visualization experiments on 3D object segmentation guided by one or a set of point prompts clicked on the first-view rendered image. The results are shown in Fig. 4. The first row shows the segmentation results of SA3D, SAGA and ours on LERF-figurines scene. The second row compares with SAGA, which distills the knowledge embedded in the SAM decoder into the feature field. Though the boundary issues in 3D-GS can be alleviated, incomplete segmentation occurs even after post-processing. In contrast, ours can solve the boundary problems while achieving more complete segmentation. More segmentation results are shown in the third row.

Text-Guided Segmentation: We further conduct 3D object segmentation experiments guided by text prompts. We follow the same process as in Table II experiments to conduct visualization experiments. In order to compare our method with SA3D under the same setting, we choose the same garden scene from Mip-NeRF 360 dataset [3], using three text prompts including “The table”, “The vase” and “The bonsai”. Results in Fig. 5 show that our approach can achieve accurate segmentation results by simply providing object names, demonstrating our method’s potential in combining with language models. Compared with the SA3D in the “The table” case, our approach can segment the complete object with desired table legs, while SA3D only gives the tabletop.

Comparison between 2D and 3D Segmentation: Our method improves the segmentation quality of SAM. Due to the sensitivity of SAM to the scene viewpoint (e.g. illumination), the 2D segmentation results might be incomplete under some views, but our 3D segmentation can still guarantee complete segmentation results via a multi-view ensemble (Fig. 7, 1st Row). Furthermore, we can obtain detailed 3D segmentation results under noisy predictions from SAM. As shown in Fig. 7,

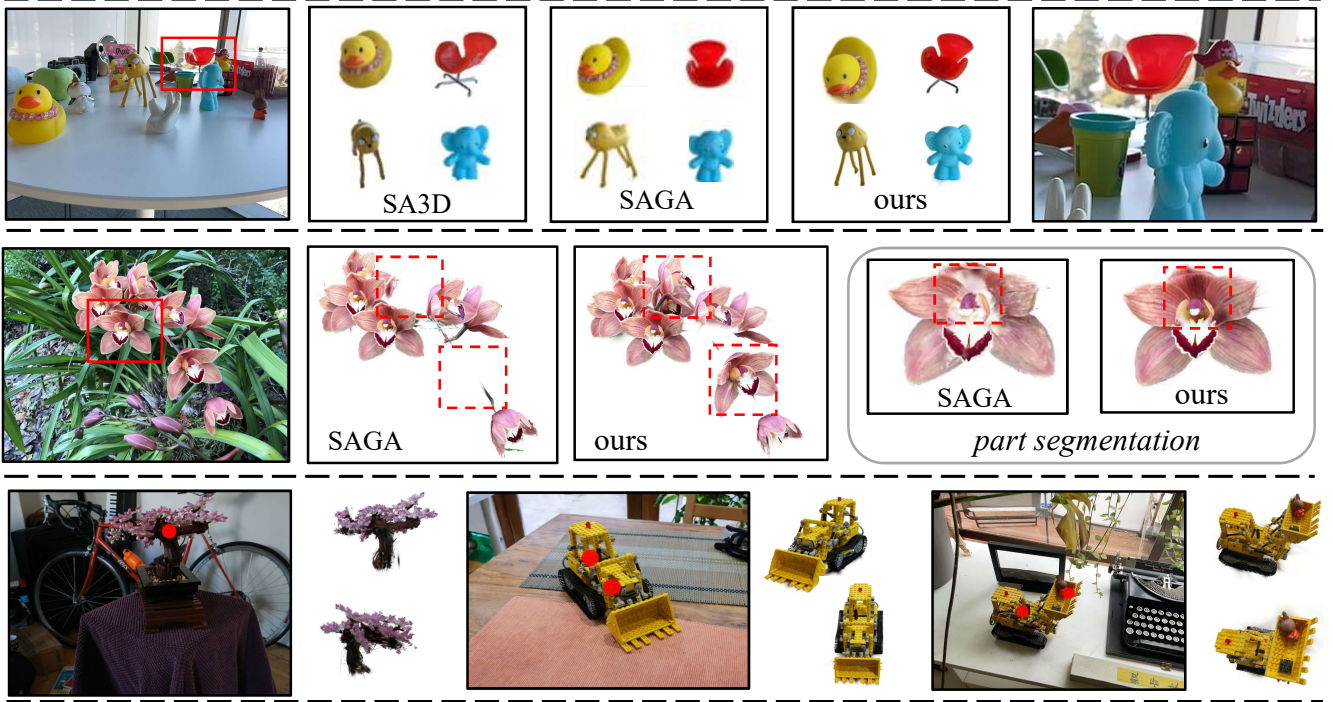


Fig. 4. Qualitative results compared with SA3D [7] and SAGA [6] in different scenes (LERF-figurines [36], SPIn-NeRF-Orchids [58], LERF-dozer-nerfgunwald [36]). We enlarge the boxed area on the right for a better visualization.

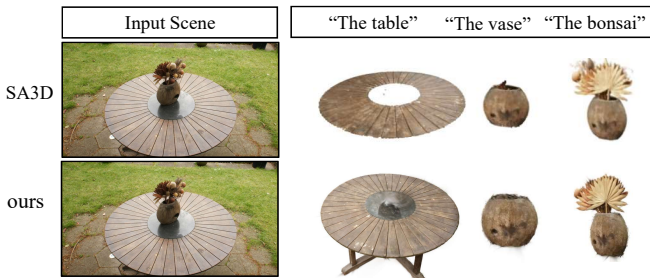


Fig. 5. 3D segmentation with text prompts in Mip-NeRF360-garden [3].

2nd Row, the multi-view voting ensures that object Gaussians are retained correctly and background Gaussians obtained from the noisy boundaries of SAM predictions are removed.

Scene Editing: Scene editing is a basic application for 3D reconstructed scenes. However, it’s quite difficult to do this without being able to locate the specific objects. This task demonstrates the ability of our approach to help edit the 3D scenes. Specifically, after segmenting the objects in the 3D Gaussian space, we can manipulate the objects by removing, translating, and rotating them. Thanks to the simplicity of the explicit Gaussian representations, without bells and whistles, our approach obtains satisfactory scene editing results, as shown in Fig. 6. The instances in 1st column with red bounding boxes are objects to be segmented. It can be seen that with objects segmented in 3D Gaussians, they can be translated and rotated in any direction in the scene. After the removal of segmented objects, original scenes can still keep intact.

Collision Detection: Collision detection is an indispensable component in practical 3D applications, such as games, movies, and simulators. In this task, we demonstrate our approach can directly help in revealing the collision body of target objects in the 3D Gaussian space. We choose two scenes from SPIn-NeRF dataset. Following the process of our segmentation method, we can obtain the corresponding segmented 3D Gaussian points in the scene. To this end, we use the Quickhull [2] algorithm to build the collision mesh upon our segmented objects. The results in the last column of Fig. 6 show that our segmented objects can successfully derive correct convex hulls for downstream applications.

F. Time Cost Analysis

We conduct the time cost analysis compared with previous 3DGS segmentation methods, namely SAGA [6] and Gaussian Grouping [92]. We select the *figurines* scene from the LERF dataset [36], which includes 20+ objects, using a single NVIDIA A100 GPU for fair evaluation. All methods utilize SAM to extract image features and then perform mask generation. Thus, they share a similar time cost, roughly 2 minutes to process 300 images (*SAM Extraction* in Tab. IV). This procedure can be regarded as a pre-processing and can be accelerated by image-batching or deploying Efficient-SAM. For extra processing time, both baseline methods require gradient descent optimization over 30,000 iterations to distill SAM features [6] or ID embedding [92] into 3D space linked with each 3D Gaussian, resulting in considerable additional training time for re-optimizing 3D scenes. On the contrary, our method only needs to perform multi-view and multi-label

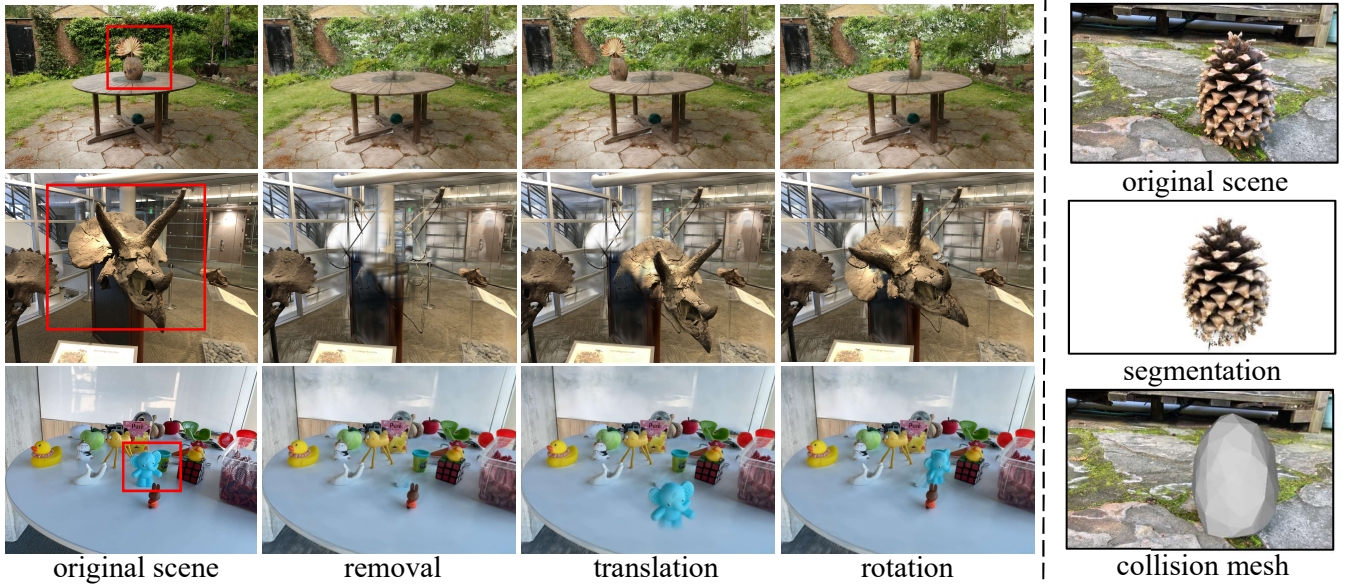


Fig. 6. Visualization examples of scene editing after the object segmentation. We offer three scene editing examples: removal, translation, and rotation. We further provide collision mesh computed after segmentation, as shown in the last column.

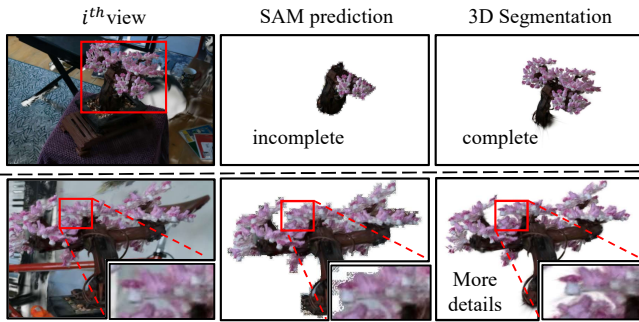


Fig. 7. Comparison between 2D and 3D Segmentation. Noisy predictions exist in SAM’s segmentation, such as incomplete objects and background regions. These issues are solved in our 3D segmentation.

TABLE IV

TIME COST COMPARISONS ON *figurines* SCENE WITH 20+ OBJECTS TO BE SEGMENTED SIMULTANEOUSLY.

Time	SAM Extraction	Extra Processing	Optimization Steps	Forward Segmentation
SAGA [6]		10 min	30000	0.5s
Ga-Grouping [92]	2 min	9 min	30000	0.3s
Ours		12s	1	0.4ms

assignments, the time cost of which is merely 12 seconds. As for segmenting a single object, our method merely requires an efficient voting process, taking just 0.4 milliseconds, which is much faster than other methods as they both demand network forwarding.

G. Ablation Study

Gaussian Decomposition: Gaussian Decomposition is proposed to address the issue of roughness boundaries of 3D segmented objects, which results from the non-negligible

TABLE V

ABLATION ON GAUSSIAN DECOMPOSITION (GD) ON SPIN-NeRF DATASET. THE 2ND AND 3RD COLUMNS COMPARE THE RESULTS BEFORE AND AFTER USING THE GAUSSIAN DECOMPOSITION (GD). THE LAST COLUMN REPRESENTS THE RESULTS OF DIRECTLY REMOVING THE GAUSSIANS ACROSS MASK BOUNDARIES.

Scene	w/ GD		w/o GD		Delete	
	IoU	Acc	IoU	Acc	IoU	Acc
Orchids	85.4	97.5	82.2	96.8	78.2	95.4
Ferns	92.0	98.9	89.2	98.4	89.8	98.5
Room	86.5	98.1	81.3	97.2	85.4	97.9
Horns	91.1	98.4	83.2	96.5	84.8	97.6
Fortress	96.5	99.4	88.5	98.1	82.3	96.7
Fork	83.4	99.3	81.8	99.2	79.9	99.1
Pinecone	92.6	98.9	91.6	98.9	82.9	97.5
Truck	93.0	97.9	93.4	97.8	91.4	96.8
Lego	90.2	99.7	88.4	99.6	82.9	99.4
mean	90.0	98.7	86.6	98.1	84.2	97.7

spatial sizes of 3D Gaussian located at the boundary. We conduct both quantitative and qualitative experiments to verify the effectiveness of our approach. For quantitative ablation experiments, Table V shows the results. We compare our proposed Gaussian Decomposition scheme with two other processing methods, one for segmentation without any special handling (the 3rd column), and the other for directly removing such Gaussians (the last column). The experiments are conducted on SPIN-NeRF dataset, following the same evaluation process as Table I. Compared with the others, our proposed approach outperforms them on all scenes with +3.4% and +5.8% IoU, respectively, demonstrating the effectiveness of Gaussian Decomposition. Besides, it can be seen that directly removing these Gaussians cannot decrease the roughness of mask boundaries, leading to even 2.4% lower IoU values.

Fig. 9 shows the visualization results for comparison. We select two representative scenes in which large-scale Gaussian

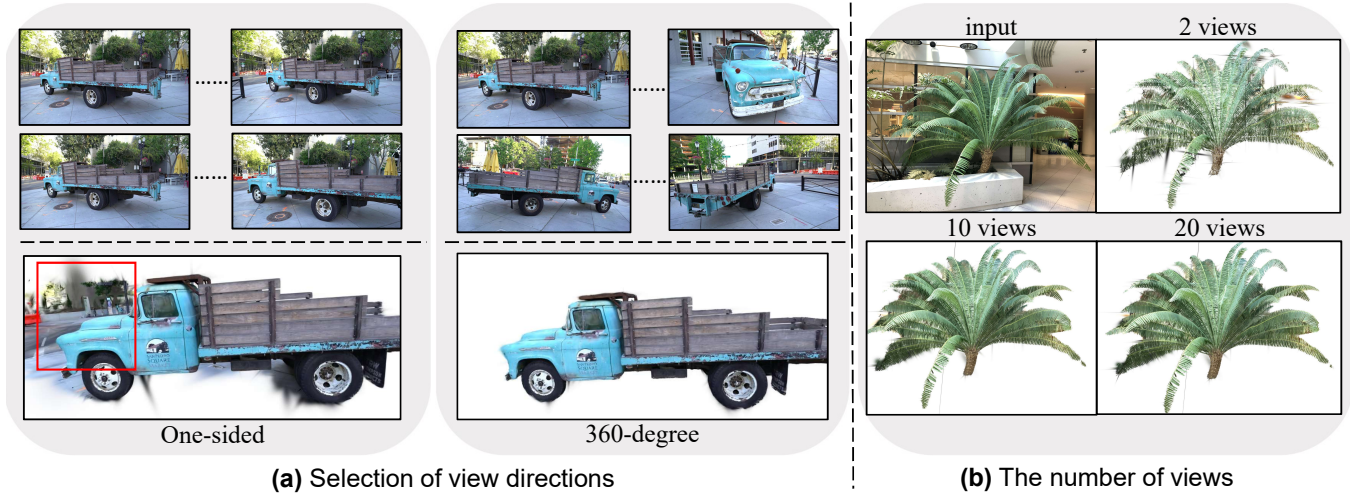


Fig. 8. Ablation on view selection, including view directions and the number of views.

TABLE VI

ABLATION ON DIFFERENT NUMBER OF VIEWS FOR 3D SEGMENTATION. NUMBERS IN PARENTHESSES REPRESENT THE USED VIEW PERCENTAGE OF THE TOTAL TRAINING VIEW.

Number of views	5(10%)	10(20%)	21(50%)	42(100%)
IoU on Fortress	91.16	92.11	93.82	95.6
Number of views	25(10%)	50(20%)	125(50%)	251(100%)
IoU on truck	90.27	90.97	92.11	93.0

TABLE VII

ABLATION ON CONFIDENCE SCORE THRESHOLD τ ON THE *truck* SCENE.

Scene / τ	0.6	0.65	0.7	0.8
pinecone	90.8	92.0	92.6	91.4
fortress	94.0	95.4	95.6	96.2
lego	89.7	89.8	90.2	90.2

points exist at the junction/contact (can be seen in the 2nd column). This issue is alleviated after utilizing the proposed Gaussian Decomposition strategy. Though this idea is simple, the improvement is obvious.

Selection of Views: In this section, we study the influence of selected 2D views on 3D Gaussian segmentation. Fig. 8 (b) demonstrates the segmentation quality using different numbers of views. It is noteworthy that the leaves in the “fern” scene are very tiny and challenging. Even so, with only two sparse views, our approach can achieve decent results, and the segmentation quality quickly improves when increasing the number of views. Table VI also draws a similar conclusion: as the number of views increases, the IoU values will also increase. It is worth noting that even with sparse views (10% percentage), we can obtain relatively decent results for both two types of scenes, which also demonstrates the robustness and effectiveness of our approach. In practical applications, using sparse views (below 10% of total views) will greatly improve efficiency, although under 100% views, our method can still be completed within one minute.

Fig. 8 (a) study the influence of view directions on final

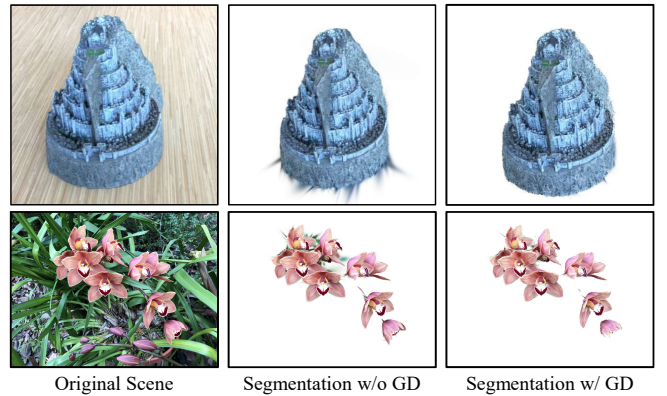


Fig. 9. Ablation results on Gaussian Decomposition (GD) on the LERF-pinecone.

segmentation results. We choose the truck with 360-degree views. The upper row shows the result of using roughly a single direction. In this case, the background is mistakenly incorporated into the segmentation. By contrast, with the same number of views but variant directions, the approach achieves high-quality clean segmentation results. This result reveals the importance of choosing non-monotone view directions in practice usage with 360 scenes.

Hyper-Parameters: In the process of our method, we set a hyper-parameter τ to control the threshold of label confidence score in Sec. III-E. To evaluate the robustness of the selection of the hyper-parameter under different scenes, we set different values to conduct experiments on Scene *pinecone*, *fortress*, and *lego* from LLFF [57] dataset. As indicated in Tab. VII, our method is not sensitive to the threshold selection. Robust results are observed in multiple scenes when setting the threshold between 0.6 and 0.8. Values within this range can obtain fine results. In all of our experiments, the value of the hyper-parameter τ is set to 0.7 by default.

V. DISCUSSIONS AND LIMITATIONS

Through experiments, we found that our method does not perform well in objects where 3D Gaussians are very sparse,

such as the LLFF-room [57] scene. The Gaussians of the table are notably sparse; even worse, the Gaussians representing the table surface are across different objects. Though our Gaussian Decomposition can somewhat remedy the boundary issue, it still suffers from extremely sparse points and leaves holes in 3D segmentation. We believe this limitation can be alleviated by future research in a more structured 3D-GS representation, yielding more accurate results.

VI. CONCLUSION

We address the issue that rough/incomplete boundaries in segmenting 3D-GS and propose a novel Boundary-enhanced segmentation pipeline. Given input prompts on the first rendering view, our approach automatically generates multi-view masks and achieves consistent 3D segmentation via the proposed assignment strategy. Our Gaussian Decomposition module can effectively mitigate the boundary roughness issue of segmented objects resulting from the inherent geometry structure in 3D-GS. Extensive segmentation experiments show that our method can effectively obtain high-quality 3D object segmentation without boundary issues, and different scene-editing tasks demonstrate that our method can be easily applied to downstream applications. Overall, we hope our work can inspire more future work in the area of 3D Gaussian representation.

REFERENCES

- [1] Jeongmin Bae, Seoha Kim, Youngsik Yun, Hahyun Lee, Gun Bang, and Youngjung Uh. Per-gaussian embedding-based deformation for deformable 3d gaussian splatting. *arXiv preprint arXiv:2404.03613*, 2024.
- [2] C Bradford Barber, David P Dobkin, and Hannu Huhdanpaa. The quickhull algorithm for convex hulls. *ACM Transactions on Mathematical Software (TOMS)*, 22(4):469–483, 1996.
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [4] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9157–9166, 2019.
- [5] Xingyuan Bu, Junran Peng, Junjie Yan, Tieniu Tan, and Zhaoxiang Zhang. Gaia: A transfer learning system of object detection that fits your needs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 274–283, 2021.
- [6] Jiazhong Cen, Jiemin Fang, Chen Yang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Segment any 3d gaussians. *CoRR*, abs/2312.00860, 2023.
- [7] Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Wei Shen, Lingxi Xie, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs. *arXiv preprint arXiv:2304.12308*, 2023.
- [8] Qing Chang, Junran Peng, Lingxi Xie, Jiajun Sun, Haoran Yin, Qi Tian, and Zhaoxiang Zhang. Data: Domain-aware and task-aware self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9841–9850, 2022.
- [9] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19457–19467, 2024.
- [10] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Proc. ECCV*, pages 333–350. Springer, 2022.
- [11] Keyan Chen, Chenyang Liu, Hao Chen, Haotian Zhang, Wenyuan Li, Zhengxia Zou, and Zhenwei Shi. Rsprompter: Learning to prompt for remote sensing instance segmentation based on visual foundation model. *IEEE TGRS*, 2024.
- [12] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014.
- [13] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- [14] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [15] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018.
- [16] Zilong Chen, Feng Wang, Yikai Wang, and Huaping Liu. Text-to-3d using gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21401–21412, 2024.
- [17] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- [18] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021.
- [19] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1316–1326, 2023.
- [20] Pinxuan Dai, Jiamin Xu, Wenxiang Xie, Xinguo Liu, Huamin Wang, and Weiwei Xu. High-quality surface reconstruction using gaussian surfels. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024.
- [21] Ruining Deng, Can Cui, Quan Liu, Tianyuan Yao, Lucas W Remedios, Shunxing Bao, Bennett A Landman, Lee E Wheless, Lori A Coburn, Keith T Wilson, et al. Segment anything model (sam) for digital pathology: Assess zero-shot segmentation on whole slide imaging. *arXiv preprint arXiv:2304.04155*, 2023.
- [22] Lue Fan, Feng Wang, Naiyan Wang, and ZHAO-XIANG ZHANG. Fully sparse 3d object detection. *Advances in Neural Information Processing Systems*, 35:351–363, 2022.
- [23] Lue Fan, Xuan Xiong, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Rangedet: In defense of range view for lidar-based 3d object detection. In *Proc. ICCV*, pages 2918–2927, 2021.
- [24] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5501–5510, 2022.
- [25] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14346–14355, 2021.
- [26] Rahul Goel, Dhawal Sirikonda, Saurabh Saini, and PJ Narayanan. Interactive segmentation of radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4201–4211, 2023.
- [27] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024.
- [28] Abdul Mueed Hafiz and Ghulam Mohiuddin Bhat. A survey on instance segmentation: state of the art. *International journal of multimedia information retrieval*, 9(3):171–189, 2020.
- [29] Qingyong Hu, Bo Yang, Sheikh Khalid, Wen Xiao, Niki Trigoni, and Andrew Markham. Towards semantic segmentation of urban-scale 3d point clouds: A dataset, benchmarks and challenges. In *Proc. CVPR*, pages 4977–4987, 2021.
- [30] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024.
- [31] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4220–4230, 2024.

- [32] Ge-Peng Ji, Deng-Ping Fan, Peng Xu, Bowen Zhou, Ming-Ming Cheng, and Luc Van Gool. Sam struggles in concealed scenes—empirical study on” segment anything”. *SCIS*, 2023.
- [33] Haian Jin, Isabella Liu, Peijia Xu, Xiaoshuai Zhang, Songfang Han, Sai Bi, Xiaowei Zhou, Zexiang Xu, and Hao Su. Tensorir: Tensorial inverse rendering. In *Proc. CVPR*, pages 165–174, 2023.
- [34] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. In *NeurIPS*, 2023.
- [35] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- [36] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lurf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023.
- [37] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollar. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [38] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- [39] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proc. CVPR*, pages 8500–8509, 2022.
- [40] Kun Lan, Haoran Li, Haolin Shi, Wenjun Wu, Yong Liao, Lin Wang, and Pengyuan Zhou. 2d-guided 3d gaussian segmentation. *arXiv preprint arXiv:2312.16047*, 2023.
- [41] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8508–8520, 2024.
- [42] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *Proc. ECCV*, pages 1–18. Springer, 2022.
- [43] Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6517–6526, 2024.
- [44] Zhihao Liang, Qi Zhang, Ying Feng, Ying Shan, and Kui Jia. Gsir: 3d gaussian splatting for inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21644–21653, 2024.
- [45] Youtian Lin, Zuozhuo Dai, Siyu Zhu, and Yao Yao. Gaussian-flow: 4d reconstruction with dynamic 3d gaussian particle. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21136–21145, 2024.
- [46] Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8576–8588, 2024.
- [47] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- [48] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8759–8768, 2018.
- [49] Yang Liu, Chuanchen Luo, Lue Fan, Naiyan Wang, Junran Peng, and Zhaoxiang Zhang. Citygaussian: Real-time high-quality large-scale scene rendering with gaussians. In *European Conference on Computer Vision*, pages 265–282. Springer, 2025.
- [50] Yang Liu, Chuanchen Luo, Zhongkai Mao, Junran Peng, and Zhaoxiang Zhang. Citygaussianv2: Efficient and geometrically accurate reconstruction for large-scale scenes. *arXiv preprint arXiv:2411.00771*, 2024.
- [51] Zhicheng Lu, Xiang Guo, Le Hui, Tianrui Chen, Min Yang, Xiao Tang, Feng Zhu, and Yuchao Dai. 3d geometry-aware deformable gaussian splatting for dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8900–8910, 2024.
- [52] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023.
- [53] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 2024.
- [54] Maciej A Mazurkowski, Haoyu Dong, Hanxue Gu, Jichen Yang, Nicholas Konz, and Yixin Zhang. Segment anything model for medical image analysis: an experimental study. *MedIA*, 2023.
- [55] Gregory P Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *Proc. CVPR*, pages 12677–12686, 2019.
- [56] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Proc. ECCV*, 2020.
- [57] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019.
- [58] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G Derpanis, Jonathan Kelly, Marcus A Brubaker, Igor Gilitschenski, and Alex Levinstein. Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20669–20679, 2023.
- [59] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022.
- [60] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H Mueller, Chakravarty R Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. In *Computer Graphics Forum*, volume 40, pages 45–59. Wiley Online Library, 2021.
- [61] Mong H Ng, Kaahan Radia, Jianfei Chen, Dequan Wang, Ionel Gog, and Joseph E Gonzalez. Bev-seg: Bird’s eye view semantic segmentation using geometry and semantic point cloud. *arXiv preprint arXiv:2006.11436*, 2020.
- [62] Cong Pan, Yonghao He, Junran Peng, Qian Zhang, Wei Sui, and Zhaoxiang Zhang. Baeformer: Bi-directional and early interaction transformers for bird’s eye view semantic segmentation. In *Proc. CVPR*, pages 9590–9599, 2023.
- [63] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Gao Huang. 3d object detection with pointformer. In *Proc. CVPR*, pages 7463–7472, 2021.
- [64] Junran Peng, Xingyuan Bu, Ming Sun, Zhaoxiang Zhang, Tieniu Tan, and Junjie Yan. Large-scale object detection in the wild from imbalanced multi-labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9709–9718, 2020.
- [65] Junran Peng, Qing Chang, Haoran Yin, Xingyuan Bu, Jiajun Sun, Lingxi Xie, Xiaopeng Zhang, Qi Tian, and Zhaoxiang Zhang. Gaia-universe: Everything is super-netify. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):11856–11868, 2023.
- [66] Junran Peng, Ming Sun, ZHAO-XIANG ZHANG, Tieniu Tan, and Junjie Yan. Efficient neural architecture transformation search in channel-level for object detection. *Advances in neural information processing systems*, 32, 2019.
- [67] Junran Peng, Ming Sun, Zhaoxiang Zhang, Tieniu Tan, and Junjie Yan. Pod: Practical object detection with scale-sensitive network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9607–9616, 2019.
- [68] Junran Peng, Lingxi Xie, Zhaoxiang Zhang, Tieniu Tan, and Jingdong Wang. Accelerating deep neural networks with spatial bottleneck modules. *arXiv preprint arXiv:1809.02601*, 2018.
- [69] Lang Peng, Zhirong Chen, Zhangjie Fu, Pengpeng Liang, and Erkang Cheng. Bevsegformer: Bird’s eye view semantic segmentation from arbitrary camera rigs. In *Proc. WACV*, pages 5935–5943, 2023.
- [70] Zelin Peng, Zhengqin Xu, Zhilin Zeng, Lingxi Xie, Qi Tian, and Wei Shen. Parameter efficient fine-tuning via cross block orchestration for segment anything model. In *CVPR*, 2024.
- [71] Zelin Peng, Zhengqin Xu, Zhilin Zeng, Xiaokang Yang, and Wei Shen. Sam-parser: Fine-tuning sam efficiently by parameter space reconstruction. In *AAAI*, 2024.
- [72] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14335–14345, 2021.

- [73] Zhongzheng Ren, Aseem Agarwala, Bryan Russell, Alexander G Schwing, and Oliver Wang. Neural volumetric object selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6133–6142, 2022.
- [74] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Nerf for outdoor scene relighting. In *European Conference on Computer Vision*, pages 615–631. Springer, 2022.
- [75] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Efficient 4d portrait editing with text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4556–4567, 2024.
- [76] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointtrcn: 3d object proposal generation and detection from point cloud. In *Proc. CVPR*, pages 770–779, 2019.
- [77] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. *Advances in neural information processing systems*, 34:12278–12291, 2021.
- [78] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proc. CVPR*, pages 5459–5469, 2022.
- [79] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Improved direct voxel grid optimization for radiance fields reconstruction. *arXiv preprint arXiv:2206.05085*, 2022.
- [80] Jiakai Sun, Han Jiao, Guangyuan Li, Zhanjie Zhang, Lei Zhao, and Wei Xing. 3dstream: On-the-fly training of 3d gaussians for efficient streaming of photo-realistic free-viewpoint videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20675–20685, 2024.
- [81] Jiayang 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.
- [82] Lv Tang, Haoke Xiao, and Bo Li. Can sam segment anything? when sam meets camouflaged object detection. *arXiv preprint arXiv:2304.04709*, 2023.
- [83] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In *2022 International Conference on 3D Vision (3DV)*, pages 443–453. IEEE, 2022.
- [84] Yuxi Wang, Junran Peng, and ZhaoXiang Zhang. Uncertainty-aware pseudo label refinery for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9092–9101, 2021.
- [85] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition*, pages 16210–16220, 2022.
- [86] Christopher Wewer, Kevin Raj, Eddy Ilg, Bernt Schiele, and Jan Eric Lenssen. latentsplat: Autoencoding variational gaussians for fast generalizable 3d reconstruction. *arXiv preprint arXiv:2403.16292*, 2024.
- [87] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20310–20320, 2024.
- [88] Junde Wu, Rao Fu, Huihui Fang, Yuanpei Liu, Zhaowei Wang, Yanwu Xu, Yueming Jin, and Tal Arbel. Medical sam adapter: Adapting segment anything model for medical image segmentation. *arXiv preprint arXiv:2304.12620*, 2023.
- [89] Dejia Xu, Ye Yuan, Morteza Mardani, Sifei Liu, Jiaming Song, Zhangyang Wang, and Arash Vahdat. Agg: Amortized generative 3d gaussians for single image to 3d. *arXiv preprint arXiv:2401.04099*, 2024.
- [90] Chen Yang, Sikuang Li, Jiemin Fang, Ruofan Liang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Gaussianobject: Just taking four images to get a high-quality 3d object with gaussian splatting. *arXiv preprint arXiv:2402.10259*, 2024.
- [91] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. Bevformer v2: Adapting modern image backbones to bird’s-eye-view recognition via perspective supervision. In *Proc. CVPR*, pages 17830–17839, 2023.
- [92] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. *arXiv preprint arXiv:2312.00732*, 2023.
- [93] Taoran Yi, Jiemin Fang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. *arXiv preprint arXiv:2310.08529*, 2023.
- [94] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021.
- [95] Heng Yu, Joel Julin, Zoltán Á Milacski, Koichiro Niinuma, and László A Jeni. Cogs: Controllable gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21624–21633, 2024.
- [96] Tao Yu, Runsen Feng, Ruoyu Feng, Jinming Liu, Xin Jin, Wenjun Zeng, and Zhibo Chen. Inpaint anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023.
- [97] Zehao Yu, Torsten Sattler, and Andreas Geiger. Gaussian opacity fields: Efficient and compact surface reconstruction in unbounded scenes. *arXiv preprint arXiv:2404.10772*, 2024.
- [98] Yifei Zeng, Yanqin Jiang, Siyu Zhu, Yuanxun Lu, Youtian Lin, Hao Zhu, Weiming Hu, Xun Cao, and Yao Yao. Stag4d: Spatial-temporal anchored generative 4d gaussians. *arXiv preprint arXiv:2403.14939*, 2024.
- [99] Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, Sung-Ho Bae, Seungkyu Lee, and Choong Seon Hong. Faster segment anything: Towards lightweight sam for mobile applications. *arXiv preprint arXiv:2306.14289*, 2023.
- [100] Dingyuan Zhang, Dingkan Liang, Hongcheng Yang, Zhikang Zou, Xiaoqing Ye, Zhe Liu, and Xiang Bai. Sam3d: Zero-shot 3d object detection via segment anything model. *SCIS*, 2023.
- [101] Guowen Zhang, Junsong Fan, Liyi Chen, Zhaoxiang Zhang, Zhen Lei, and Lei Zhang. General geometry-aware weakly supervised 3d object detection. In *European Conference on Computer Vision*, pages 290–309. Springer, 2025.
- [102] Guowen Zhang, Lue Fan, Chenhang He, Zhen Lei, Zhaoxiang Zhang, and Lei Zhang. Voxel mamba: Group-free state space models for point cloud based 3d object detection. *arXiv preprint arXiv:2406.10700*, 2024.
- [103] Shougao Zhang, Mengqi Zhou, Yuxi Wang, Chuanchen Luo, Rongyu Wang, Yiwei Li, Xucheng Yin, Zhaoxiang Zhang, and Junran Peng. Cityx: Controllable procedural content generation for unbounded 3d cities. *arXiv preprint arXiv:2407.17572*, 2024.
- [104] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. *Advances in Neural Information Processing Systems*, 34:10326–10338, 2021.
- [105] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (ToG)*, 40(6):1–18, 2021.
- [106] Zhaoxiang Zhang, Cong Pan, and Junran Peng. Delving into the effectiveness of receptive fields: Learning scale-transferrable architectures for practical object detection. *International Journal of Computer Vision*, 130(4):970–989, 2022.
- [107] Xu Zhao, Wenchao Ding, Yongqi An, Yinglong Du, Tao Yu, Min Li, Ming Tang, and Jinqiao Wang. Fast segment anything. *arXiv preprint arXiv:2306.12156*, 2023.
- [108] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021.
- [109] Mengqi Zhou, Yuxi Wang, Jun Hou, Chuanchen Luo, Zhaoxiang Zhang, and Junran Peng. Scenex: Procedural controllable large-scale scene generation via large-language models. *arXiv preprint arXiv:2403.15698*, 2024.
- [110] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suyu You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. *arXiv preprint arXiv:2312.03203*, 2023.
- [111] Shijie Zhou, Zhiwen Fan, Dejia Xu, Haoran Chang, Pradyumna Chari, Tejas Bharadwaj, Suyu You, Zhangyang Wang, and Achuta Kadambi. Dreamscene360: Unconstrained text-to-3d scene generation with panoramic gaussian splatting. *arXiv preprint arXiv:2404.06903*, 2024.
- [112] Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Drivinggaussian: Composite gaussian

- splatting for surrounding dynamic autonomous driving scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21634–21643, 2024.
- [113] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS'01.*, pages 29–538. IEEE, 2001.