


FlashSplat: 2D to 3D Gaussian Splatting Segmentation Solved Optimally

QiuHong Shen, Xingyi Yang, and Xinchao Wang^{*} 

National University of Singapore
{qiuHong.shen,xyang}@u.nus.edu xinchao@nus.edu.sg

Abstract. This study addresses the challenge of accurately segmenting 3D Gaussian Splatting (3D-GS) from 2D masks. Conventional methods often rely on iterative gradient descent to assign each Gaussian a unique label, leading to lengthy optimization and sub-optimal solutions. Instead, we propose a straightforward yet globally optimal solver for 3D-GS segmentation. The core insight of our method is that, with a reconstructed 3D-GS scene, the rendering of the 2D masks is essentially a linear function with respect to the labels of each Gaussian. As such, the optimal label assignment can be solved via linear programming in closed form. This solution capitalizes on the alpha blending characteristic of the splatting process for single step optimization. By incorporating the background bias in our objective function, our method shows superior robustness in 3D segmentation against noises. Remarkably, our optimization completes within 30 seconds, about $50\times$ faster than the best existing methods. Extensive experiments demonstrate our method’s efficiency and robustness in segmenting various scenes, and its superior performance in downstream tasks such as object removal and inpainting. Demos and code will be available at <https://github.com/florinshen/FlashSplat>.

Keywords: 3D Segmentation · 3D Gaussian Splatting · Neural-based understanding

1 Introduction

The pursuit of understanding and interacting with 3D environments represents a formidable yet pivotal challenge in computer vision. Central to this endeavor is the task of accurately perceiving and segmenting 3D structures [1, 5, 10, 30, 32, 42], a task that becomes increasingly complex as we delve into more sophisticated representations of 3D scenes. Recently, 3D Gaussian Splatting (3D-GS) [17] emerges as a cutting-edge approach that promises to revolutionize how we render and reconstruct 3D spaces. By employing a multitude of colored 3D Gaussians, this method achieves a high fidelity representation of 3D scenes, offering a compelling blend of precision and visual quality that is particularly suited for complex object and scene rendering.

^{*} Corresponding Author.

Despite the potential of 3D-GS, a significant obstacle remains: the segmentation of these 3D Gaussian from mere 2D masks—a process crucial for a range of applications from object recognition to scene manipulation. Existing approaches [3, 53] to this challenge have largely depended on iterative gradient descent methods for labeling 3D Gaussian. However, these methods are marred by their slow convergence speed and frequent entrapment in suboptimal solutions, rendering them impractical for applications requiring real-time performance or high accuracy.

Addressing this gap, our work introduces a simple but globally optimal solver designed explicitly for the segmentation of 3D Gaussian Splatting. Our approach capitalizes on the insight that the process of rendering segmented 2D images from a reconstructed 3D Gaussian Splatting can be simplified to a linear function concerning the accumulated contribution of each Gaussian. This realization allows us to frame the problem as a linear integer programming that can be solved in a closed form, relying solely on the alpha composition term inherent in the splatting process. This breakthrough significantly streamlines the segmentation task, bypassing the need for iterative optimization and directly leading to the optimal label assignment.

Moreover, by integrating a background bias into our objective function, we further enhance the method’s robustness against 2D masks noise in 3D segmentation. This refinement not only bolsters the robustness of our solution but also its applicability across a wider range of scenes and conditions. Impressively, our solver completes its optimization in approximately 30 seconds—an acceleration of $50\times$ faster than existing methods—without sacrificing accuracy or global optimality.

Through extensive experimentation, we have validated the superiority of our approach in efficiently segmenting a variety of scenes, demonstrating its enhanced performance in downstream tasks such as object removal and inpainting. These results underscore the potential of our method to significantly advance the field of 3D scene processing and understanding.

The primary contributions of our work are summarized as follows:

- We introduce a globally optimal solver for 3D Gaussian Splatting segmentation, which significantly enhances the efficiency of lifting 2D segmentation results into 3D space.
- We simplify the 3D-GS rendering process through linearization, turning the 2D to 3D segmentation task into a linear integer optimization problem. This method is effective for both binary and scene segmentation.
- We introduce a background bias in the optimization, demonstrate superior robustness against noises in 3D segmentation, showcasing our method’s robustness and efficiency across a variety of scene segmentation.
- Our method achieves a remarkable optimization speed, completing the process within 30 seconds—approximately 50 times faster than existing methods—while ensuring global optimality for given 2D masks.
- Extensive experiments validate the superiority of our approach in downstream tasks, including object removal and inpainting, thus highlighting its potential to significantly impact 3D data processing and application.

2 Related-work

2.1 3D Gaussian Splatting

3D Gaussian Splatting has emerged as an efficient method for inverse rendering, facilitating the reconstruction of 3D scenes from 2D images through learning explicit 3D Gaussian in the space [9, 18]. Recent advancements have seen a series of studies [28, 51, 52] extending this approach to capture dynamic 3D environments by learning deformation fields for 3D Gaussians. Concurrently, another stream of researches [6, 24, 35, 44, 54, 57] have integrated 3D Gaussian splatting with diffusion-based models [25, 38, 39, 49, 50, 56] to generate static [40, 55, 61] or dynamic [36] 3D objects. A notable portion of these works leverages the explicit representations provided by 3D Gaussian Splatting for rapid optimization. Our method primarily concentrates on elevating corresponding 2D masks onto reconstructed 3D Gaussian scenes. We also exploit the explicit representation of 3D Gaussian Splatting to cast the segmentation of 3D scenes as a Linear Programming optimization problem, thereby enabling instantaneous 3D segmentation.

2.2 3D Neural scene segmentation

Recent advancements in neural 3D scene representations, including NeRF, DVGO, and several others [1, 5, 8, 10, 14, 23, 30, 32, 42, 48], have demonstrated exceptional capabilities in 3D scene reconstruction. Inspired by these developments, the exploration of 3D segmentation within such frameworks has emerged as a significant area of interest. Semantic-NeRF [60] initially showcased NeRF’s aptitude for semantic propagation and refinement. NVOS [37] furthered this by enabling user interaction for 3D object selection via a lightweight MLP trained on specialized 3D features. Concurrently, a series of studies including N3F, DFF, LERF, and ISRF [12, 19, 22, 46] have integrated 2D vision models to facilitate 3D segmentation by elevating 2D visual features, thereby processing embedded 3D features for segmentation. Additionally, various approaches have combined instance and semantic segmentation methods with radiance fields [2, 4, 11, 12, 15, 27, 33, 41, 47, 58] to further explore this domain.

Most closely related methods to our work are SAGA [3], Gaussian Grouping [53], and SAGS [16]. These works concentrate on associating 2D masks, as generated by the Segment Anything Model (SAM) [20], with 3D Gaussian Splatting for mask lifting. Gaussian Grouping employs video tracker to associate 2D masks, subsequently distilling these as object features for 3D Gaussians and segmenting via a classifier trained with 2D identity loss and 3D regularization loss. Similarly, SAGA distills 2D masks into 3D features, requiring additional training features for each 3D Gaussians with a pack of losses. Both of these methods are limited by their substantial training costs required to optimize Gaussian features before 3D segmentation. SAGS introduces a training-free approach by projecting the centers of 3D Gaussians onto 2D masks. This method adopts a intuitive segmentation criterion, where a projected center landing within a foreground mask gets classified as a foreground Gaussian.

Our method distinguishes itself by framing the 2D mask lifting in 3D Gaussian Splatting as a Linear Programming (LP) optimization problem, solved in a single step. Focused solely on lifting associated 2D masks into 3D Gaussians masks, our approach offers a standalone, efficient solution for immediate segmentation without the need for extra training and post-processing.

3 Methods

In this section, we first delve into the rendering process of 3D Gaussian Splatting (3DGS), focusing on the tile-based rasterization and alpha blending. We then describe how this process lends itself to formulating the segmentation of 3DGS as an integer linear programming (ILP) optimization, which we demonstrate can be resolved in a closed form. Recognizing the typically noisy nature of 2D masks, we introduce a softened optimal assignment to mitigate these noises. Beyond binary segmentation, we extend our method to include scene segmentation, enabling the segmentation of all objects within 3D scenes. Finally, we present a method for rendering 2D masks guided by depth information, which projects the 3D segmentation results onto 2D masks from new viewpoints.

3.1 Preliminary: Rasterization of 3DGS

3D Gaussian Splatting [18] (3DGS) stands out as a great method for novel view synthesis. Unlike neural radiance fields [31], it reconstructs a 3D scene as explicit 3D Gaussians with real-time speed. Given a set of captured views with paired camera poses, 3D Gaussian Splatting reconstruct 3D scenes by learning 3D Gaussians $\{G_i\}$. Each 3D Gaussian G_i is parameterized as $G_i = \{m_i, q_i, s_i, o_i, c_i\}$, where $m_i \in \mathbb{R}^3$ is the center position, $q_i \in \mathbb{R}^4$ is a quaternion representing rotations, $s_i \in \mathbb{R}^3$ is the scale of three dimensions, $o_i \in \mathbb{R}$ is the learnt opacity, and $c_i \in \mathbb{R}^{48}$ is the three-order spherical harmonics to represent view dependent color. In its rendering process, 3DGS adopts a rasterization pipeline to achieve super efficiency. Specifically, all 3D Gaussians are first projected onto the image plane as 2D Gaussians. Then the whole image is divided into 16×16 tiles as illustrated in Fig. 1, pixels in each tile B share identical 3D Gaussians subset $\{G_i\}_B \in \{G_i\}$. When rendering each pixels, the traditional alpha composition is applied to blend properties x_i (color or depth) of these 2D Gaussians into pixel space property X by their depth order:

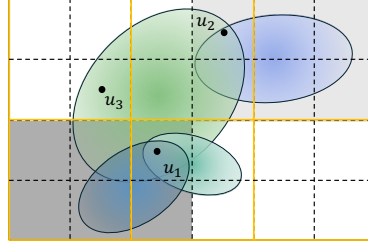


Fig. 1: GS Rasterization. Here we illustrate projected 2D Gaussians in 3DGS rasterization, setting each tile as 2×2 for illustration purposes. Gaussians are shared between different tiles and are likely to be shared among different instances (color blocks).

$$X = \sum_{i \in \{G_i\}_B} x_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) = \sum_{i \in \{G_i\}_B} x_i \alpha_i T_i \quad (1)$$

where α_i is the alpha value when rendering specific pixels, which is the product of o_i and the probability of pixel position drops in the projected 2D Gaussians parameterized as $(m_i^{2D}, \Sigma_i^{2D})$, and $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$ is the transmittance, presenting the fraction that passes through and is not absorbed by the preceding $i - 1$ Gaussians.

3.2 Binary Segmentation as Integer Linear Programming

We start with a reconstructed 3DGS scene, parameterized by $\{G_i\}$. It includes L rendered views with associated 2D binary masks, denoted as $\{M^v\}$. Each mask, $M \in \mathbb{R}^{H \times W}$ has elements of 0 or 1, where 0 represents the background, and 1 denotes the foreground. Our goal is to assign a 3D label, P_i , which can be either 0 or 1, to each 3D Gaussian G_i by projecting the 2D masks into the 3D space.

To do this, we optimize P_i through a differentiable rendering process. We assume P_i as a property x_i in Eq. 1. The aim is to minimize the discrepancy between the rendered masks and the provided binary masks M^v .

Luckily, with $\{G_i\}$ reconstructed, the α_i and the T_i become constant for rendering. As such, the rendering function becomes a linear equation with respect to the blending properties x_i . This gives us a great flexibility to solve for the optimal mask assignment using simple linear optimization.

Formally, our segmentation problem can be formulated as an integer Linear Programming (LP) optimization with mean absolute error:

$$\begin{aligned} \text{Min}_{\{P_i\}} \quad \mathcal{F} &= \sum_v |\mathcal{R}(\{G_i\}, \{P_i\}) - M^v| = \sum_{v \in L} \sum_{M_{jk}^v \in M^v} \left| \sum_i P_i \alpha_i T_i - M_{jk}^v \right| \\ \text{subject to} \quad &P_i \in \{0, 1\}. \end{aligned} \quad (2)$$

where \mathcal{R} denotes the 3DGS renderer. Given $\alpha_i \geq 0$ and $T_i \leq 1$ for any i , and the given light in alpha composition can only be absorbed, the total absorbed light cannot exceed the initial light intensity, which is normalized to 1. Consequently, the sum of absorbed light fractions across all samples is bounded by $0 < \sum_i \alpha_i T_i \leq 1$. This leads us to introduce:

Lemma 1: *In the alpha blending within 3D Gaussian Splatting, for $P_i \in \{0, 1\}$,*

$$0 \leq \sum P_i \alpha_i T_i \leq \sum \alpha_i T_i \leq 1 \quad (3)$$

To solve this LP, we refer *Lemma 1* to rewrite the Eq 2 as:

$$\min \mathcal{F} = \sum_{v,i,j,k} P_i \alpha_i T_i \mathbb{I}(M_{jk}^v, 0) - \sum_{v,i,j,k} P_i \alpha_i T_i \mathbb{I}(M_{jk}^v, 1) + C \quad (4)$$

$$= C + \sum_i P_i \left(\sum_{v,j,k} \alpha_i T_i \mathbb{I}(M_{jk}^v, 0) - \sum_{v,j,k} \alpha_i T_i \mathbb{I}(M_{jk}^v, 1) \right) \quad (5)$$

$$= C + \sum_i P_i (A_0^i - A_1^i) \quad (6)$$

where $C = \sum_{v,j,k} M_{jk}^v$ is a constant value, alongside $\mathbb{I}(\cdot, 0)$ and $\mathbb{I}(\cdot, 1)$, serve as indicator functions signifying background and foreground presence, respectively. **Solving ILP as Majority Vote.** To minimize the function \mathcal{F} , as outlined in Eq. 4, we assign $P_i = 0$ when all the corresponding pixels in masks M^v indicate background (0). Conversely, if there is a contradiction among the masks, we resolve this by as a weighted majority vote. Specifically, we assign the value of each Gaussian G_i based on the most frequent label in the masks, as detailed in Eq 6.

Formally, if the difference $(A_0 - A_1) > 0$, where A_0 and A_1 represent the counts of weighted masks indicating background and foreground respectively, we assign $P_i = 0$ to minimize \mathcal{F} . If $(A_0 - A_1) < 0$, then $P_i = 1$ is assigned. In a more general context, we can formulate the optimal assignment for $\{P_i\}$ as follows:

$$P_i = \arg \max_n A_n, \quad n \in \{0, 1\} \quad (7)$$

where $A_n = \sum_{v,j,k} \alpha_i T_i \mathbb{I}(M_{jk}^v, n)$

Intuitively, this optimal assignment is conceptualized to aggregate the contributions of individual 3D Gaussians during its rendering. Gaussians that significantly contribute to the foreground in all given masks are designated for foreground as $P_i = 1$, conversely, those contributing predominantly to the background are systematically allocated to the background as $P_i = 0$. Additionally, thanks to the simple linear combination form in the objective of Equation 6, we can concurrently assign optimal labels to each 3D Gaussians.

Regularized ILP for Noise Reduction. Practically, the given 2D mask set $\{M^v\}$ is typically predicted and associated by trained 2D vision models, which is likely to introduce noise in certain regions (illustrated in column 1 and column 2 of Fig. 6). This characteristic of the provided 2D masks can lead to noisy 3D segmentation results. For instance, as depicted in Fig.2 (a), background Gaussians misclassified as foreground due to noise can result in a segmented 3D object with sharp, difficult-to-filter edges.

To address this, we refine above optimal assignment. This refinement entails initially normalize overall contributions of Gaussians through $L1$ normalization, represented as $\bar{A}_e = A_e / \sum_t A_t$. We then introduce a background bias, ranging between $\gamma \in [-1, 1]$, to recalibrate $\bar{A}_0 = \bar{A}_0 + \gamma$, thereby adjusting the optimal assignment as $P_i = \arg \max_n \{\bar{A}_0, \bar{A}_1\}$. Employing of $\gamma > 0$ effectively diminish foreground noise in segmentation outcomes; conversely, $\gamma < 0$ results in a

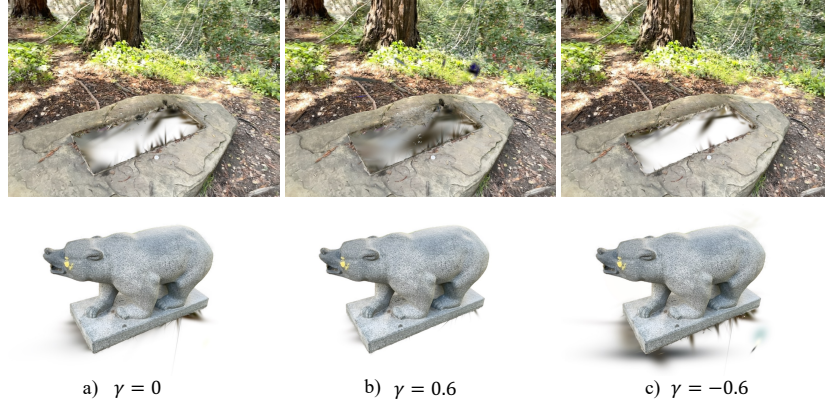


Fig. 2: The effect of the background bias γ . By adjusting the background bias in our optimal assignment, the noisy 3D segmentation caused by 2D masks can be flexibly mitigated for various downstream applications.

cleaner background, as demonstrated in Fig.2 (b) and (c), respectively. This softened form of optimal assignment provides flexibility to produce accurate 3D segmentation results against 2D masks noises for different downstream tasks.

3.3 From Binary to Scene Segmentation

Numerous instances are present across various views in 3D scenes. Segmenting multiple objects within these scenes requires executing binary segmentation multiple times according to above formulation. This process involves repeatedly gathering $\{A_0, A_1\}$, which inherently slows down the pace of scene segmentation. Thus, we extend our methodology from binary segmentation to encompass scene segmentation to address this challenge more efficiently.

This transition to multi-instance segmentation is motivated by two key considerations. Initially, it’s important to recognize that 3D Gaussians do not exclusively belong to a single object. This is exemplified in Figure 1, where pixel u_1 and u_2 are influenced by the same 3D Gaussian, despite belonging to distinct objects (color block). Furthermore, the introduction of multiple instances complicates the constraint in Equation 2 to $P_i \in \{0, 1, \dots, E - 1\}$, with E representing the total number of instances in a scene. Consequently, the set of provided masks M^v becomes $M^v \in \{0, 1, \dots, E - 1\}$ to accommodate multiple segmented 2D instances. Such constraints prevent achieving a global optimum, as the labels across instances are subject to be exchangeable.

To circumvent these challenges, we reinterpret multi-instance segmentation as a combination of binary segmentation, modifying the optimal assignment strategy as outlined in Equation 7. To isolate a specific instance labeled e in a 3D scene, we redefine all other objects within the instance set as the background

to compute A_{others} . This approach is formalized as follows:

$$\begin{aligned}
P_i &= \arg \max_n A_n, \quad n \in \{0, t\} \\
\text{where } A_t &= \sum_{v,j,k} \alpha_i T_i \mathbb{I}(M_{jk}^v, t), \\
A_0 = A_{others} &= \sum_{e \neq t} \sum_{v,j,k} \alpha_i T_i \mathbb{I}(M_{jk}^v, e)
\end{aligned} \tag{8}$$

With this formulation, we only need to accumulate the set $\{A_e\}$ once, and then perform $\arg \max$ on this set to get Gaussians subset $\{G_i\}_e$ for each object e . Consequently, this allows users to selectively remove or modify these 3D Gaussian subsets by specifying object ID. It is also noteworthy that subsets $\{G_i\}_e$ from different instances may overlap under this formulation when $\gamma < 0$, a reflection of the inherent non-exclusive nature of 3D Gaussian Splatting (see more details in supplementary material Sec. 8.3).

3.4 Depth-guided Novel View Mask Rendering

Given above formulation eschews the need for dense optimization, our approach is capable of yielding robust segmentation results using merely approximately 10% of the masked 2D views. This efficiency also empowers our method with the capability to produce 2D masks \hat{M}^v for previously unseen views. For rendering masks in novel views within the context of binary segmentation, we focus exclusively on rendering foreground Gaussians identified by $P_i = 1$. This process involves computing an accumulated alpha value ρ_{jk} for every pixel (j, k) . For novel view mask rendering in binary segmentation, we merely render foreground Gaussians with $P_i = 1$ to produce accumulated alpha value ρ_{jk} for each pixel, then the 2D mask can be obtained by simple quantization on alpha values ρ_{jk} with threshold τ , which is a pre-defined hyper parameter. Specifically, 2D masks are derived through straightforward quantization of these alpha values, where $\hat{M}_{jk}^v = \mathbb{Q}(\rho_{jk}, \tau)$ and τ represents a predetermined threshold, and \mathbb{Q} is the quantization function, when $\rho_{jk} > \tau$, $\mathbb{Q} = 1$, conversely $\mathbb{Q} = 0$.

For novel view mask rendering in scene segmentation, due to the intersections among segmentation results of different objects as indicated in Equation 8, the resultant alpha mask might present ambiguity, as illustrated in 3rd column of Fig. 3. Specifically, rendering each object’s associated 3D Gaussian subset $\{G_i\}_e$ in the same viewpoint might result in several objects meeting the condition $\mathbb{Q}(\rho_{jk}^e, \tau) = 1$. In this case, we introduce depth to determine the final segmentation results. The depth at each pixel location (j, k) is utilized to filter the final 2D mask outcome. The object e satisfying $\mathbb{Q}(\rho_{jk}^e, \tau) = 1$ with minimal depth D_{jk}^e relative to the camera at a given pixel (j, k) is selected as $\hat{M}_{jk} = e$. Figure 3 illustrates the prediction of 2D masks for both binary and multi-object segmentation in novel views.

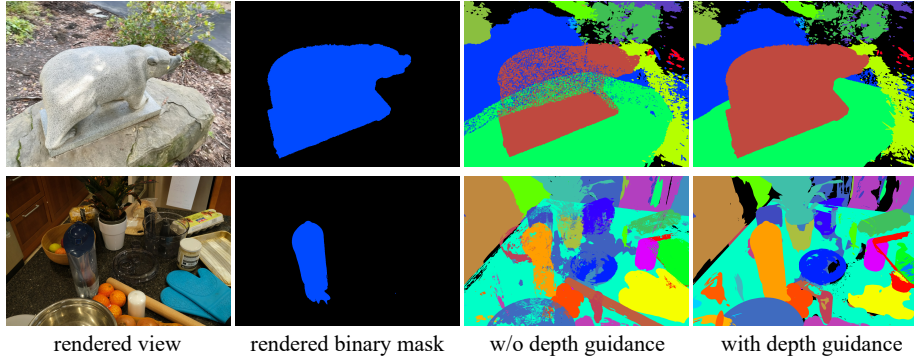


Fig. 3: Novel view mask rendering. Here we showcase mask rendering on novel views for both binary and scene segmentation. Simple alpha mask quantization can generate consistent masks in binary segmentation. With depth guidance, scene segmentation also can generate feasible 2D masks in novel views.

4 Experiments

4.1 Data preparation

Dataset. To assess the efficacy of our approach, we collect 3D scenes from several sources: the MIP-360 [1] dataset, T&T [21] dataset, LLFF [29] dataset, Instruct-NerF2NerF [13], and the LERF [19] dataset, serving as the basis for qualitative analysis. For quantitative analysis, we utilize the NVOS [37] dataset.

2D mask generation and association. In our experimental setup, we utilize the Segment Anything Models (SAM) [20] to extract masks, given that SAM’s segmentation output is inherently semantic-agnostic. It becomes necessary to further associate these 2D masks within our framework. Our approach diverges into two distinct strategies, tailored respectively for binary and scene segmentation. For binary segmentation, where the objective is to isolate a single foreground entity, we initiate by marking point prompts on a single reference view. These point prompts are projected back to the 3D space with reference view camera pose to find their nearest 3D Gaussians with the smallest positive depth. Subsequently, these point prompts are propagated to other views by projecting their corresponding 3D Gaussians’ center. Leveraging these associated point prompts, SAM independently generates a binary mask for each view. For of scene segmentation, our methodology begins with employing SAM to produce instance masks for individual views. To assign each 2D object with a unique ID in the 3D scene, multiple views are treated akin to a video sequence. Utilizing a zero-shot video tracker [7, 53], we ensure the consistent association and propagation of objects across viewpoints.

4.2 Implementation details

We implement the optimal 3D segmentation both in Eq. 7 and Eq. 8 within CUDA kernel functions. The computation of A_i for each Gaussian employs tile-

based rasterization with the renderer, excluding the color component c_i and focusing solely on $G_i = \{m_i, q_i, s_i, o_i\}$. For each pixel in the provided mask set M , we perform alpha composition as defined in Equation 1. The product $\alpha_i T_i$ for each Gaussian is then aggregated into the buffer A_e , corresponding to the pixel label $e = M_{jk}^v$, through atomic operations. Leveraging the rasterization pipeline enables the completion of the set A_e for all objects in under 30 seconds. For binary segmentation, the $\arg \max$ is directly applied to determine the optimal assignment. In the case of scene segmentation, $\arg \max$ is executed iteratively alongside dynamic programming to speedup. This assignment ensures that the computation for both binary and scene segmentation concludes in 1 millisecond, thereby facilitating interactive adjustment of the γ value once the set $\{A_e\}$ is computed. This adjustment allows users to effectively reduce segmentation noise across a variety of downstream tasks.

4.3 3D segmentation results

In Fig.4, we exhibit the results of both binary and scene 3D segmentation. The first row presents the Figurines scene from the LERF dataset [19], and the second row features the Counter scene from the MIP-360 dataset [1]. On both scenes, we apply our scene segmentation approach, rendering 2 views for 5 segmented objects (circled in the ground-truth images) for each scene, demonstrating our method’s capability in conducting scene segmentation with instance masks predicted by SAM [20]. Additionally, binary segmentation is showcased in rows 3, 4 and 5, with row 3 illustrating the Horns scene from the LLFF dataset [29], row 4 displaying the Truck scene from the T&T dataset [21], row 5 displaying the Kitchen scene from the MIP-360 dataset [1]. Two views of the segmented objects are rendered, showing our approach’s capability in segmenting 3D objects.

4.4 Object Removal

3D object removal involves entirely eliminating the 3D Gaussians subset of an object from the scene. Results of such removals are illustrated in Figure 4. It’s important to note that we apply a background bias $\gamma = -0.4$ across all scenes to ensure the background remains clear. For small 3D objects in row 1 and row 2, we simultaneously remove 5 objects from the scene. As the background of these two scenes is likely to be observed in other views, the artifacts in these two scenes are minor even imperceptible for certain objects. For more challenging scenes depicted in rows 3, 4, and 5, the space vacated by the removed objects reveals a noisier background or even black holes, primarily due to the obstruction of background by larger foreground objects. This situation is pronounced in the Horns scene (row 3), which comprises only facing-forward views.

4.5 Object Inpainting

Following 3D object removal, object inpainting [34] aims to correct unobserved region artifacts, ensuring view consistency within the 3D scene. Initially, we render views post-removal and employ Grounding-DINO [26] to identify artifact



Fig. 4: Qualitative result of FlashSplat. FlashSplat is capable of performing both binary segmentation and scene segmentation. With our single step optimal assignment, all 3D segmentation is completed within 30 seconds. These segmentation results show our FlashSplat can robustly segment 3D objects and remove objects in 3D scenes.

regions in each view, which are tracked across views using a video tracker [7]. Pretrained 2D inpainting models [43] then generate inpainted 2D views. Subsequently, the 3DGS parameters are refined with these views by introducing 200K new Gaussians near the original object locations, maintaining background Gaussians frozen. Fine-tuning involves $L1$ loss outside object masks to minimize background impact, and LPIPS loss [59] inside the inpainting masks for scene naturalness and consistency. We display object inpainting results in Fig. 5, rendering three views per scene. Noises and holes are diminished after the object inpainting, demonstrating that our method can effectively separate the foreground from the background in 3D segmentation.

4.6 Quantitative comparison

We perform a quantitative analysis using the NVOS dataset [37], derived from the LLFF dataset [29]. The NVOS dataset comprises 8 facing-forward 3D scenes from LLFF, providing a reference and a target view with 2D segmentation mask

Method	mIOU (%) \uparrow	mAcc (%) \uparrow
NVOS [37]	39.4	73.6
ISRF [12]	70.1	92.0
SGISRF [45]	83.8	96.4
SA3D [4]	90.3	98.2
SAGA [3]	90.9	98.3
FlashSplat (ours)	91.8	98.6

Table 1: Quantitative comparison on the NVOS [37] dataset.

annotations for each scene. Our approach begins with sampling point prompts from the reference view’s annotated mask, which are then propagated to other views to guide the segmentation of SAM [20]. Following this, we apply our binary segmentation approach as depicted in Eq. 2 to isolate the foreground 3D Gaussians. Subsequently, we employ novel view mask rendering with a threshold $\tau = 0.1$ to render 2D masks for the target view, allowing us to calculate the mean Intersection over Union (IoU) and mean accuracy across all 8 scenes. The comparative results, presented in Table 1, benchmark our FlashSplat against previous NeRF-based 3D segmentation approaches such as NVOS [37], ISRF [12], SGIRF [45], SA3D [4], and the 3D-GS method SAGA [3]. Our method demonstrates superior performance in both metrics on this dataset.

4.7 Computation cost

We assess the computational efficiency of FlashSplat in comparison to previous 3DGS segmentation methods, namely SAGA [3] and Gaussian Grouping [53]. This evaluation is conducted on the Figurines scene from the LERF dataset [19] using a single NVIDIA A6000 GPU. Both baseline methods necessitate the utilization of gradient descent optimization over 30,000 iterations to distill 2D masks into object features associated with each 3D Gaussian, thereby incurring significant extra training time for optimized 3D scene. In contrast, our approach merely requires the calculation of the set $\{A_i\}_e$, a process that completes in approximately 26 seconds, making it roughly 50 \times faster than baselines. For segmenting a single 3D object, these baselines demand network forwarding, whereas FlashSplat efficiently employs arg max for optimal assignment determination, taking just 0.4 milliseconds. Furthermore, an analysis of GPU memory usage reveals that our peak memory consumption is only half of that required by the previous method SAGA [3].

	Extra Time	Optimization Steps	Segmentation Time	Peak Memory
SAGA [3]	18 min	30000	0.5 s	15G
Gaussian Grouping [53]	37 min	30000	0.3 s	34G
FlashSplat (ours)	26 s	1	0.4 ms	8G

Table 2: Computation cost comparisons over the Figurines scene.

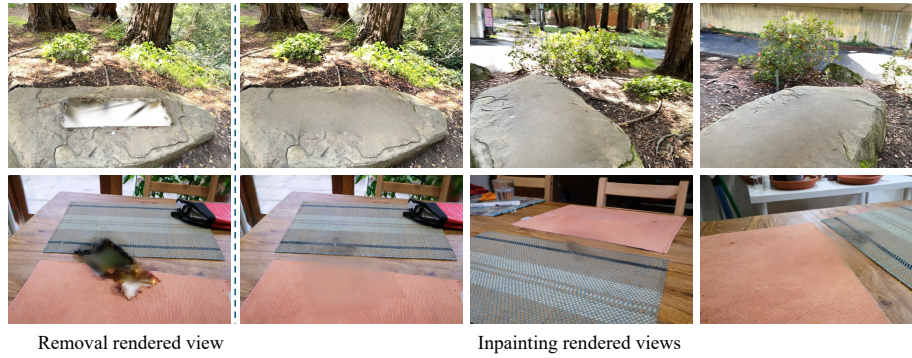


Fig. 5: Object inpainting after removal. Here we showcase the object removal results from 3D scenes, including the kitchen scene from MIP-360 dataset and bear scene from Instruct-NeRF2NeRF [15]. After removal, we inpaint the regions with artifacts by tuning the optimized 3DGS parameters. Artifacts after object removal are diminished by our object inpainting.

4.8 Ablation study

The effect of noise reduction. To further elucidate the previously mentioned noise present in 2D masks, we provide visualizations of the 2D masks generated by SAM [20] in the left column of Fig. 6 for two scenes. Additionally, we render the object mask after 3D segmentation on corresponding views, revealing that the broken regions in the provided 2D masks are remedied. This demonstrates the robustness of our method in generating 3D segmentation despite the presence of noise in the 2D masks.



Fig. 6: Noise reduction with 3D segmentation. Here we show two noisy masks per scene, and the corresponding 2D mask is rendered from the 3D segmentation results. The broken regions in the given 2D masks are remedied from the 3D segmentation.

3D Segmentation with fewer 2D Masks. Since our approach does not rely on multi-step gradient descent optimization, it naturally enables the utilization of fewer 2D masks for segmentation. To validate this capability, we consider

2D mask subsets comprising $1/4$, $1/8$, $1/16$, and $1/32$ of the total number of views in the scene. Subsequently, we visualize the 3D segmentation results by rendering the foreground 3D Gaussians onto unseen views. As depicted in Fig. 7, our method is capable of producing decent segmentation results even with only $1/8$ of the total view masks for these two 360-degree scenes. However, with fewer views, such as $1/16$ and $1/32$, many 3D Gaussians remain unseen in the provided masks, leading to potential artifacts in the segmentation outcomes.

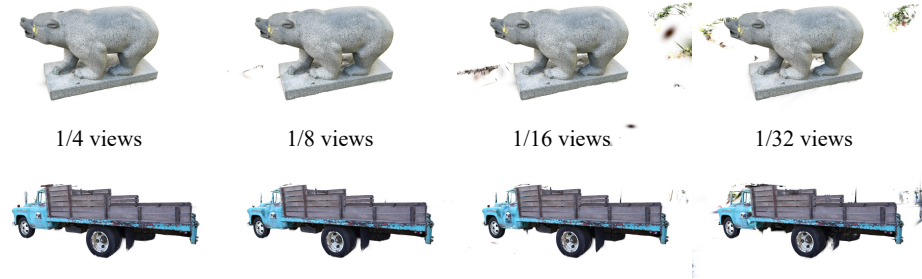


Fig. 7: 3D Segmentation with fewer masks. Here we render 3D segmentation results on novel views with given fewer 2D masks, showing our method is capable of producing decent segmentation results with only 10% of total view numbers.

5 Conclusion

In this work, we introduce an optimal solver for 3D Gaussian Splatting segmentation from 2D masks, significantly advancing the accuracy and efficiency of lifting 2D segmentation into 3D space. By breaking the alpha composition in 3D-GS into overall contributions of each Gaussians, this solver only requires single step optimization to get the optimal assignment. It not only expedit the optimization process by approximately 50 times faster compared to previous methods but also enhanced robustness against noise with a simple background bias. Further, this approach is extended to scene segmentation and capable of rendering masks on novel views. Extensive experiments demonstrate superior performance in scene segmentation tasks, including object removal and inpainting. We hope this work will facilitate 3D scene understanding and manipulating in the future.

Acknowledgement

This project is supported by the National Research Foundation, Singapore, under its Medium Sized Center for Advanced Robotics Technology Innovation.

References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In: CVPR (2022)
2. Bing, W., Chen, L., Yang, B.: Dm-nerf: 3d scene geometry decomposition and manipulation from 2d images. arXiv preprint arXiv:2208.07227 (2022)
3. Cen, J., Fang, J., Yang, C., Xie, L., Zhang, X., Shen, W., Tian, Q.: Segment any 3d gaussians. arXiv preprint arXiv:2312.00860 (2023)
4. Cen, J., Zhou, Z., Fang, J., Yang, C., Shen, W., Xie, L., Jiang, D., Zhang, X., Tian, Q.: Segment anything in 3d with nerfs. In: NeurIPS (2023)
5. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorf: Tensorial radiance fields. In: ECCV (2022)
6. Chen, Z., Wang, F., Liu, H.: Text-to-3d using gaussian splatting. arXiv preprint arXiv:2309.16585 (2023)
7. Cheng, H.K., Oh, S.W., Price, B., Schwing, A., Lee, J.Y.: Tracking anything with decoupled video segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1316–1326 (2023)
8. Fang, J., Yi, T., Wang, X., Xie, L., Zhang, X., Liu, W., Nießner, M., Tian, Q.: Fast dynamic radiance fields with time-aware neural voxels. In: SIGGRAPH Asia (2022)
9. Fei, B., Xu, J., Zhang, R., Zhou, Q., Yang, W., He, Y.: 3d gaussian as a new vision era: A survey. arXiv preprint arXiv:2402.07181 (2024)
10. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022)
11. Fu, X., Zhang, S., Chen, T., Lu, Y., Zhu, L., Zhou, X., Geiger, A., Liao, Y.: Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In: 3DV (2022)
12. Goel, R., Sirikonda, D., Saini, S., Narayanan, P.: Interactive segmentation of radiance fields. arXiv preprint arXiv:2212.13545 (2022)
13. Haque, A., Tancik, M., Efros, A., Holynski, A., Kanazawa, A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. In: ICCV (2023)
14. Hedman, P., Srinivasan, P.P., Mildenhall, B., Barron, J.T., Debevec, P.E.: Baking neural radiance fields for real-time view synthesis. In: ICCV (2021)
15. Hu, B., Huang, J., Liu, Y., Tai, Y.W., Tang, C.K.: Instance neural radiance field. arXiv preprint arXiv:2304.04395 (2023)
16. Hu, X., Wang, Y., Fan, L., Fan, J., Peng, J., Lei, Z., Li, Q., Zhang, Z.: Semantic anything in 3d gaussians. arXiv preprint arXiv:2401.17857 (2024)
17. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics (ToG) (2023)
18. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM TOG **42**(4), 1–14 (2023)
19. Kerr, J., Kim, C.M., Goldberg, K., Kanazawa, A., Tancik, M.: Lerf: Language embedded radiance fields. arXiv preprint arXiv:2303.09553 (2023)
20. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023)
21. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: Benchmarking large-scale scene reconstruction. ACM Trans. Graph. (2017)
22. Kobayashi, S., Matsumoto, E., Sitzmann, V.: Decomposing nerf for editing via feature field distillation. In: NeurIPS (2022)

23. Lindell, D.B., Martel, J.N.P., Wetzstein, G.: Autoint: Automatic integration for fast neural volume rendering. In: CVPR (2021)
24. Ling, H., Kim, S.W., Torralba, A., Fidler, S., Kreis, K.: Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. arXiv preprint arXiv:2312.13763 (2023)
25. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023)
26. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023)
27. Liu, X., Chen, J., Yu, H., Tai, Y., Tang, C.: Unsupervised multi-view object segmentation using radiance field propagation. In: NeurIPS (2022)
28. Luiten, J., Kopanas, G., Leibe, B., Ramanan, D.: Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. arXiv preprint arXiv:2308.09713 (2023)
29. Mildenhall, B., Srinivasan, P.P., Cayon, R.O., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: practical view synthesis with prescriptive sampling guidelines. ACM Trans. Graph. (2019)
30. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020)
31. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020)
32. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. ACM Trans. Graph. (2022)
33. Niemeyer, M., Geiger, A.: GIRAFFE: representing scenes as compositional generative neural feature fields. In: CVPR (2021)
34. Qiu, J., Yang, Y., Wang, X., Tao, D.: Scene essence. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8322–8333 (2021)
35. Ren, J., Pan, L., Tang, J., Zhang, C., Cao, A., Zeng, G., Liu, Z.: Dreamgaussian4d: Generative 4d gaussian splatting. arXiv preprint arXiv:2312.17142 (2023)
36. Ren, J., Xie, K., Mirzaei, A., Liang, H., Zeng, X., Kreis, K., Liu, Z., Torralba, A., Fidler, S., Kim, S.W., et al.: L4gm: Large 4d gaussian reconstruction model. arXiv preprint arXiv:2406.10324 (2024)
37. Ren, Z., Agarwala, A., Russell, B.C., Schwing, A.G., Wang, O.: Neural volumetric object selection. In: CVPR (2022)
38. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR. pp. 10684–10695 (2022)
39. Shen, Q., Yang, X., Wang, X.: Anything-3d: Towards single-view anything reconstruction in the wild. arXiv preprint arXiv:2304.10261 (2023)
40. Shen, Q., Yi, X., Wu, Z., Zhou, P., Zhang, H., Yan, S., Wang, X.: Gamba: Marry gaussian splatting with mamba for single view 3d reconstruction. arXiv preprint arXiv:2403.18795 (2024)
41. Stelzner, K., Kersting, K., Kosiorek, A.R.: Decomposing 3d scenes into objects via unsupervised volume segmentation. arXiv preprint arXiv:2104.01148 (2021)
42. Sun, C., Sun, M., Chen, H.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: CVPR (2022)

43. Suvorov, R., Logacheva, E., Mashikhin, A., Remizova, A., Ashukha, A., Silvestrov, A., Kong, N., Goka, H., Park, K., Lempitsky, V.: Resolution-robust large mask inpainting with fourier convolutions. In: WACV (2022)
44. Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653 (2023)
45. Tang, S., Pei, W., Tao, X., Jia, T., Lu, G., Tai, Y.W.: Scene-generalizable interactive segmentation of radiance fields. In: ACM3D (2023)
46. Tschernezki, V., Laina, I., Larlus, D., Vedaldi, A.: Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In: 3DV (2022)
47. Vora, S., Radwan, N., Greff, K., Meyer, H., Genova, K., Sajjadi, M.S., Pot, E., Tagliasacchi, A., Duckworth, D.: Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. arXiv preprint arXiv:2111.13260 (2021)
48. Wizardwongsa, S., Phongthawee, P., Yenphraphai, J., Suwajanakorn, S.: Nex: Real-time view synthesis with neural basis expansion. In: CVPR (2021)
49. Wu, Z., Zhou, P., Yi, X., Yuan, X., Zhang, H.: Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9892–9902 (2024)
50. Yang, X., Wang, X.: Hash3d: Training-free acceleration for 3d generation. arXiv preprint arXiv:2404.06091 (2024)
51. Yang, Z., Yang, H., Pan, Z., Zhu, X., Zhang, L.: Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. arXiv preprint arXiv:2310.10642 (2023)
52. Yang, Z., Gao, X., Zhou, W., Jiao, S., Zhang, Y., Jin, X.: Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. arXiv preprint arXiv:2309.13101 (2023)
53. Ye, M., Danelljan, M., Yu, F., Ke, L.: Gaussian grouping: Segment and edit anything in 3d scenes. arXiv preprint arXiv:2312.00732 (2023)
54. Yi, T., Fang, J., Wu, G., Xie, L., Zhang, X., Liu, W., Tian, Q., Wang, X.: Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. arXiv preprint arXiv:2310.08529 (2023)
55. Yi, X., Wu, Z., Shen, Q., Xu, Q., Zhou, P., Lim, J.H., Yan, S., Wang, X., Zhang, H.: Mvgamba: Unify 3d content generation as state space sequence modeling. arXiv preprint arXiv:2406.06367 (2024)
56. Yi, X., Wu, Z., Xu, Q., Zhou, P., Lim, J.H., Zhang, H.: Diffusion time-step curriculum for one image to 3d generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9948–9958 (2024)
57. Yin, Y., Xu, D., Wang, Z., Zhao, Y., Wei, Y.: 4dgen: Grounded 4d content generation with spatial-temporal consistency. arXiv preprint arXiv:2312.17225 (2023)
58. Yu, H., Guibas, L.J., Wu, J.: Unsupervised discovery of object radiance fields. In: ICLR (2022)
59. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR (2018)
60. Zhi, S., Laidlow, T., Leutenegger, S., Davison, A.J.: In-place scene labelling and understanding with implicit scene representation. In: ICCV (2021)
61. Zou, Z.X., Yu, Z., Guo, Y.C., Li, Y., Liang, D., Cao, Y.P., Zhang, S.H.: Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10324–10335 (2024)