

COB-GS: Clear Object Boundaries in 3DGS Segmentation Based on Boundary-Adaptive Gaussian Splitting

Jiaxin Zhang, Junjun Jiang*, Youyu Chen, Kui Jiang, Xianming Liu
Harbin Institute of Technology

Abstract

Accurate object segmentation is crucial for high-quality scene understanding in the 3D vision domain. However, 3D segmentation based on 3D Gaussian Splatting (3DGS) struggles with accurately delineating object boundaries, as Gaussian primitives often span across object edges due to their inherent volume and the lack of semantic guidance during training. In order to tackle these challenges, we introduce Clear Object Boundaries for 3DGS Segmentation (COB-GS), which aims to improve segmentation accuracy by clearly delineating blurry boundaries of interwoven Gaussian primitives within the scene. Unlike existing approaches that remove ambiguous Gaussians and sacrifice visual quality, COB-GS, as a 3DGS refinement method, jointly optimizes semantic and visual information, allowing the two different levels to cooperate with each other effectively. Specifically, for the semantic guidance, we introduce a boundary-adaptive Gaussian splitting technique that leverages semantic gradient statistics to identify and split ambiguous Gaussians, aligning them closely with object boundaries. For the visual optimization, we rectify the degraded suboptimal texture of the 3DGS scene, particularly along the refined boundary structures. Experimental results show that COB-GS substantially improves segmentation accuracy and robustness against inaccurate masks from pre-trained model, yielding clear boundaries while preserving high visual quality. Code is available at <https://github.com/ZestfulJX/COB-GS>.

1. Introduction

Understanding and interacting with 3D scenes has long been a critical challenge in computer vision and computer graphics. This task involves reconstructing 3D scenes from collections of images or videos, as well as accurately perceiving and segmenting 3D structures. In recent years, researchers have conducted extensive research on 3D scene representation and perception. Among these advancements, 3D Gaussian Splatting (3DGS) [17], an emerging real-time radiance field rendering technique, has demonstrated comparable

*Correspondence to: Junjun Jiang (junjunjiang@hit.edu.cn)

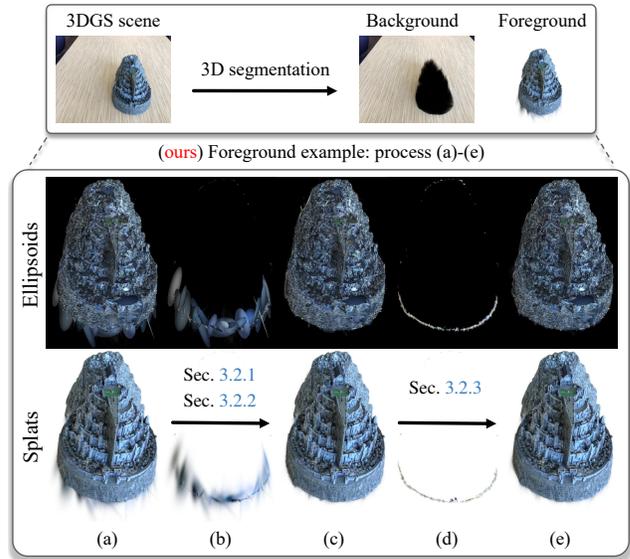


Figure 1. Problems on 3DGS segmentation. Our task is to achieve high-quality 3D segmentation on 3DGS, including foreground and background. Take foreground as an example: (a) unclear segmentation result of existing methods; (b) blurry boundary Gaussians; (c) segmentation results from joint optimization; (d) tiny ambiguous Gaussians due to incorrect masks from pre-trained models; (e) final segmentation results after improving robustness from COB-GS.

novel view rendering quality against Neural Radiance Fields (NeRF) [28], while significantly outperforming NeRF in the speed of optimization and rendering, enabling real-time rendering. As an explicit scene representation, 3DGS opens new avenues for 3D scene perception and interaction.

3D segmentation is fundamental for effective 3D scene perception and interaction, and it is crucial to obtain accurate segmentation results within the neural 3D representations exemplified by 3DGS. Currently, two predominant methodologies exist for executing the segmentation of 3DGS: feature-based methods and mask-based methods. Feature-based methods typically operate in conjunction with 3D scene reconstruction to learn distinctive feature properties for each Gaussian primitive. During the segmentation phase, the similarity between the 3D Gaussian features and the queried feature is computed to specify the Gaussian primi-

tives with the desired semantic [4, 30, 37, 44, 48]. However, these methods encounter challenges for inefficient training and rendering processes, as well as the inherent ambiguity associated with high-dimensional feature representations.

To mitigate these concerns, mask-based post-processing methods leverage the semantic masks of input views from Segment Anything Model (SAM) [19] to learn category labels for each 3D Gaussian in the reconstructed 3DGS scene, filtering these Gaussian primitives with the specified query label to perform 3D segmentation [5, 16, 36]. Despite this advancement, the original scene reconstruction often neglects the semantic information, focusing primarily on visual optimization while overlooking the volumetric characteristics of the Gaussian primitives. This oversight will lead to blurred labels for the boundary Gaussians during scene segmentation, which result in imprecise segmentation results characterized by blurry object edges, as illustrated in Figure 1. Some of existing methods delete the ambiguous boundary Gaussians directly [5, 36]. However, simply removing the Gaussian primitives on the boundary will disturb the visual quality.

To confront these issues, we propose COB-GS, a 3DGS refinement and segmentation method that jointly optimizes semantics and appearance to register semantic masks to Gaussian primitives. Similar to existing approaches, we introduce *mask label* as an additional attribute to each Gaussian for segmentation. Moreover, we reveal a strong correlation between the gradient direction of these labels and the supervising category at the pixel level, which is a strong discriminator of ambiguous Gaussian primitives on the boundary. Specifically, during the mask optimization phase, our approach utilizes gradient statistics of the mask label to identify and split boundary Gaussians, allowing precise alignment with object edges. In the scene optimization phase, we refine the scene texture on the correct boundary structure to maintain visual quality. After scene optimization, 3D segmentation focuses on filtering the mask labels. Additionally, we distinguish between boundary blurring due to the volume of Gaussians and inaccurate masks. By refining tiny boundary Gaussians, we exclusively enhance the robustness of our method against inaccurate masks from the pre-trained model. Finally, we introduce a two-stage mask generation method based on SAM2 [34], significantly simplifying the extraction of region-of-interest masks from 3D reconstruction datasets.

To summarize our contributions in a few words:

- To the best of our knowledge, we are the first 3DGS segmentation method explicitly designed to jointly optimize semantic and visual information, ensuring they enhance one another, aligning Gaussians with object edges to efficiently obtain clear boundaries and improve visual quality.
- We propose a boundary-adaptive Gaussian splitting technique that leverages gradient information from semantics to refine ambiguous boundary Gaussians, along with a boundary-guided scene texture restoration method to preserve scene visual quality on the refined boundary.

- We demonstrate the robustness of our method to inaccurate masks from pre-trained model by extracting and refining the tiny boundary Gaussians during optimization.
- We introduce a two-stage mask generation method using SAM2 based on text prompts, effectively addressing the object continuity issues in long sequence prediction.

2. Related Works

3D Gaussian Splatting. As an emerging real-time inverse rendering technology, 3DGS [17] has been proven to match the high rendering quality of the NeRF [28] in the novel view synthesis, and the speed is much faster than NeRF. Current advancements in 3DGS focus on improving aspects such as reconstruction quality [9, 25, 46], reconstruction speed [12, 26], and storage consumption [11, 21]. Other efforts aim to address special cases, including dynamic scenes [22, 41] and challenging inputs [6, 49]. As an explicit representation, 3DGS also provides more possibilities for 3D editing [7, 31, 42] and 3D generation [23, 39]. Our approach focuses on the 3D segmentation with clear boundaries within 3DGS scenes.

3D Neural Scene Segmentation. Recent advancements in neural 3D scene representation [13, 17, 28, 38] and 2D foundational models [3, 19, 32, 34, 43] have significantly enhanced the ability to perceive and interact with 3D scenes. These methods focus on learning additional attributes for 3D representations leveraging foundational 2D models, expanding beyond color to address a range of tasks. One key area of research is 3D segmentation. Early NeRF [28], as an implicit neural representation, was commonly used for 3D segmentation [5, 14, 15, 18, 24, 29, 33, 47]. However, it faced challenges with decoupling due to the inherent limitations of neural networks. In contrast, 3DGS provides an explicit representation that facilitates better region decoupling in 3D segmentation, offering a more effective alternative to NeRF.

There are two primary methods for 3D segmentation on 3DGS. The feature-based approach [10, 30, 37, 48]: Grouping Gaussian [44] learns identity encodings for each 3D Gaussian and groups them with the same encodings, enforcing spatial consistency in 3D to constrain the identity encoding learning process. SAGA [4] improves segmentation across different scales by incorporating scale-related affinity features into each Gaussian. The closest approach to ours is the mask-based approach [5, 16, 36], in which SAGD [16] employs a cross-view label voting mechanism and Gaussian decomposition to enhance foreground quality. FlashSplat [36] addresses the inverse rendering of 2D masks using linear programming, and introduces the background bias to eliminate boundary Gaussian. Existing methods typically separate scene segmentation from reconstruction, neglecting the volume of 3D Gaussians, leading to inaccurate segmentation and blurry boundaries. In contrast, our method uses joint optimization of semantics and texture to achieve clear object boundaries while preserving visual quality.

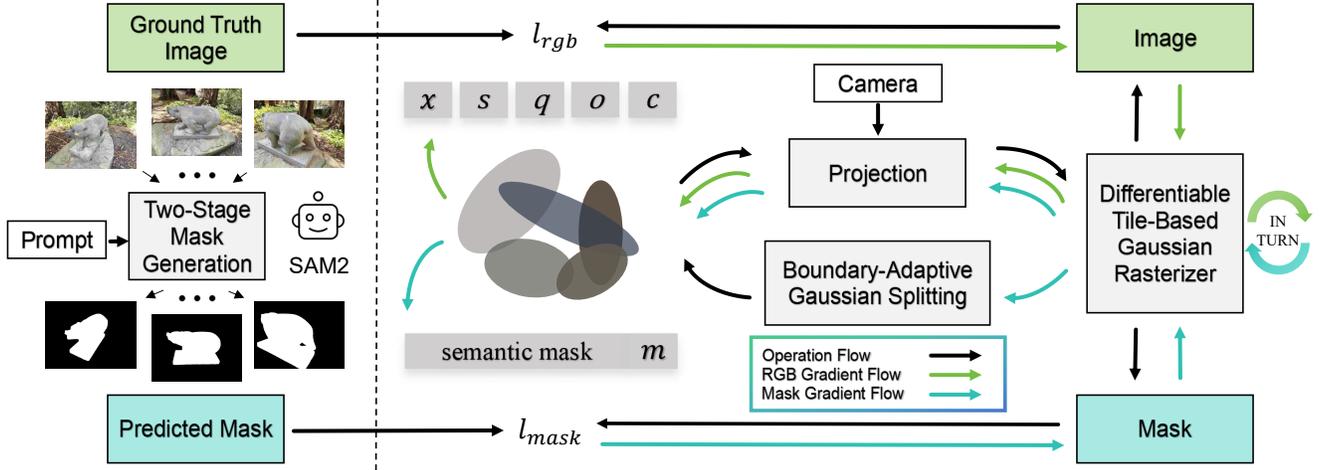


Figure 2. Pipeline of our clear object boundary segmentation method for 3DGS. On the left, we present our two-stage mask generation method, which utilizes SAM2 perform mask prediction on image sequences based on text prompt to obtain masks for regions of interest. Images and masks serve as supervision for 3DGS refinement. On the right, for the reconstructed 3DGS scene, we jointly and alternately optimize the mask and texture. For the mask optimization, boundary-adaptive Gaussian splitting is performed to refine boundary structure.

3. Method

In the section, we first introduce the principles of 3D Gaussian Splatting in Sec 3.1 to clarify the rasterization process. In Sec 3.2, we elaborate on how to jointly optimize semantics and texture, and its core is the boundary-adaptive Gaussian splitting based on mask gradient statistics. Finally, we introduce two-stage mask generation on SAM2 in Sec 3.3. Figure 2 depicts the pipeline of our proposed COB-GS.

3.1. Preliminary: 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) [17] is an emerging real-time radiance field rendering technique that achieves rendering quality comparable to NeRF. As an explicit scene representation, it opens new possibilities for scene segmentation and editing. Given a set of V input views $\{I^v\}$ along with their corresponding camera poses, 3DGS can represent the scene by learning a set of Gaussians $\{G_i\}$. In the original 3DGS, the i -th 3D Gaussian primitive can be parameterized as $G_i = \{x_i, s_i, q_i, o_i, c_i\}$, where $x_i \in \mathbb{R}^3$ is the center position, $s_i \in \mathbb{R}^3$ is the scale, $q_i \in \mathbb{R}^4$ is the rotation, $o_i \in \mathbb{R}$ is the opacity, and $c_i \in \mathbb{R}^{48}$ denotes the color.

During rendering, all 3D Gaussians are first projected onto the image plane as 2D Gaussians. The set of projected 3D Gaussians in the shared space is then accessed in parallel in the form of pixel blocks. Specifically, when rendering a pixel, traditional alpha compositing blends the target attributes p_i (color, depth, or semantic features) of the 2D Gaussians into the pixel space P :

$$P = \sum_{i=1}^N p_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) = \sum_{i=1}^N p_i \alpha_i T_i, \quad (1)$$

where α_i is the product of the opacity of the i -th Gaussian

primitive and the probability of its projection’s distance from the center position of the 2D Gaussians decaying exponentially. $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$ represents the transmittance, indicating the score of penetrating through the previous $i - 1$ Gaussians to the current Gaussian.

3.2. Boundary-Aware Object Segmentation

For efficiency, we introduce a continuous *mask label* $m_i \in (0, 1)$ for each Gaussian, where m_i close to 1 indicates that the i -th Gaussian is necessary for segmenting the 3D object, while m_i close to 0 indicates that the i -th Gaussian belongs to the background. Similar to the color rendering process, the mask labels of the 3D Gaussian primitives are combined through alpha compositing to yield the mask result in the 2D pixel space M_{render} :

$$M_{\text{render}} = \sum_{i=1}^N m_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) = \sum_{i=1}^N m_i \alpha_i T_i. \quad (2)$$

Inspired by SA3D [5], we use a similar loss function to supervise the mask label training process:

$$\mathcal{L}_{\text{mask}} = \sum_{M_{jk}^v} M_{jk}^v \cdot M_{\text{render}}^v + \sum_{M_{jk}^v} (1 - M_{jk}^v) \cdot M_{\text{render}}^v, \quad (3)$$

where $M_{jk}^v \in \{0, 1\}$ is the ground truth mask at pixel location (j, k) for view v from $\{M^v\}$, and the mask generation method is detailed in Sec 3.3. Unlike SA3D, which does not limit the range of m_i , we constrain the mask labels $m_i \in (0, 1)$. Based on the absorption light score of all sampled Gaussians in the alpha compositing formula

$0 < \sum_i \alpha_i T_i < 1$, we can deduce $0 < \sum_i m_i \alpha_i T_i < 1$, ensuring that $\mathcal{L}_{\text{mask}}$ converges. Additionally, we eliminate hyperparameters that determine negative loss in SA3D and instead emphasize learning the background to achieve 3D segmentation rather than just foreground extraction.

3.2.1. Boundary-Adaptive Gaussian Splitting

The original 3DGS relied on RGB supervision, which lacked the object-level semantic information used to shape 3D Gaussians. As a result, the segmentation of freezing geometry and texture will lead to semantically ambiguous boundary Gaussians. Therefore, it is crucial to locate and split these ambiguous Gaussians for obtaining clear object boundaries.

To address this issue, unlike the inefficient forward voting process used in existing methods [16, 36] for each Gaussian, we use the gradient of mask labels in the mask optimization phase for ambiguous Gaussian identification. Specifically, for a pixel location (j, k) at viewpoint v , the derivative of $\mathcal{L}_{\text{mask}}^{vjk}$ with respect to m_i can be expressed as:

$$\frac{d\mathcal{L}_{\text{mask}}^{vjk}}{dm_i} = \begin{cases} -\alpha_i T_i, & \text{if } M_{jk}^v = 1 \\ \alpha_i T_i, & \text{if } M_{jk}^v = 0 \end{cases} \quad (4)$$

During optimization, the viewpoint v is used as the minimum unit of the gradient calculation, and Gaussian has the volume. Therefore the gradient calculation for $\mathcal{L}_{\text{mask}}^v$ with respect to m_i is influenced by multiple pixels, leading to:

$$\frac{d\mathcal{L}_{\text{mask}}^v}{dm_i} = \sum_{j=1}^{N_{v,i}^+} (-\alpha_j T_j) + \sum_{j=1}^{N_{v,i}^-} (\alpha_j T_j), \quad (5)$$

where $N_{v,i}^+$ and $N_{v,i}^-$ represent the number of signals supervising the i -th Gaussian with ground truth masks of 1 and 0, respectively. Thus, the cumulative gradient at a viewpoint v is not effective for distinguishing ambiguous boundary Gaussians, while the sign of the gradient under a pixel reflects the supervised category of the mask labels.

Motivated by the relationship between gradient direction and supervising signals, we introduce a new variable for each Gaussian during backpropagation to capture the consistency strength of supervision signals on the mask label under a viewpoint, using the absolute value of the relative distance:

$$\text{mask_sig}_{v,i} = \left| \frac{N_{v,i}^+ - N_{v,i}^-}{N_{v,i}^+ + N_{v,i}^- + \epsilon} \right|, \quad (6)$$

where ϵ is a small constant. The closer mask_sig is to 0, the stronger the supervision by the inconsistent signal. During the optimization of mask labels, Gaussians with an absolute relative distance below a threshold are identified as the semantically ambiguous boundary Gaussian set $\{G_i\}_B$.

$$\{G_i\}_B = \{G_i \mid i \in \mathcal{I} \wedge \left(\frac{1}{V} \sum_{v=1}^V \text{mask_sig}_{v,i} < \delta \right)\}, \quad (7)$$

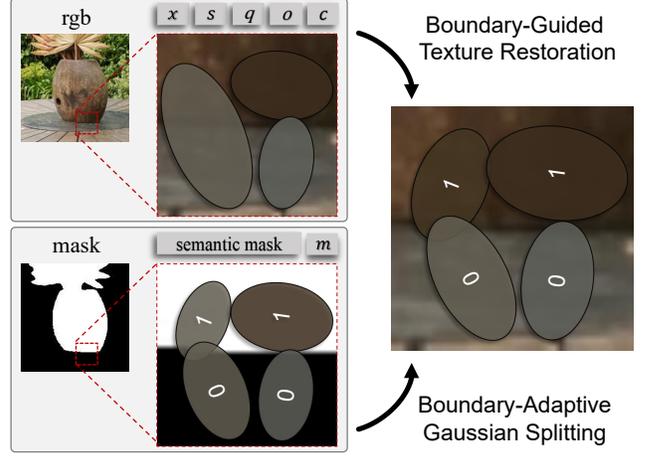


Figure 3. Supervision based solely on texture results in large Gaussians due to the similarity in textures between objects. However, object-level mask supervision facilitates the differentiation of object edges. This allows 3D Gaussians to split along object edges, while also guiding the correct restoration of scene textures.

where δ is the threshold and $\mathcal{I} = \{1, 2, 3, \dots, |\{G_i\}|\}$. During the splitting process, we refer to the original 3D Gaussian Splatting. First, we exclude small-scale Gaussians from $\{G_i\}_B$. For the remaining larger Gaussians, we replace each with two smaller Gaussians, scaling down from the original. We then use the original Gaussian as the probability density function (PDF) to sample their initial positions.

3.2.2. Boundary-Guided Scene Texture Restoration

Existing scene segmentation methods directly remove boundary ambiguous Gaussians [5, 36] or focus only on the scales of foreground Gaussians [16]. These rough methods can compromise visual quality, and object-level semantic conditions are not fully utilized in scene texture learning.

To solve this problem, our unique insight is that accurate object boundaries can enhance 3D segmentation and improve scene texture optimization. As shown in Figure 3, we alternately learn the mask labels and the geometric and texture of the Gaussians. Incorporating object-level semantic information effectively limits the volume of boundary Gaussians, and optimizing texture on accurate boundary structure enhances visual quality for new views.

Specifically, the loss function for learning the geometric and texture information of the Gaussians is consistent with the original Gaussian optimization process:

$$\mathcal{L}_{\text{rgb}} = (1 - \lambda)\mathcal{L}_1 + \lambda\mathcal{L}_{\text{D-SSIM}}, \quad (8)$$

where λ is a hyperparameter. In the alternating optimization process, we first optimize the mask labels by minimizing $\mathcal{L}_{\text{mask}}$ while freezing the geometry and texture of the Gaussians. As mentioned in Sec 3.2.1, we locate and split the semantically ambiguous boundary Gaussians under a certain

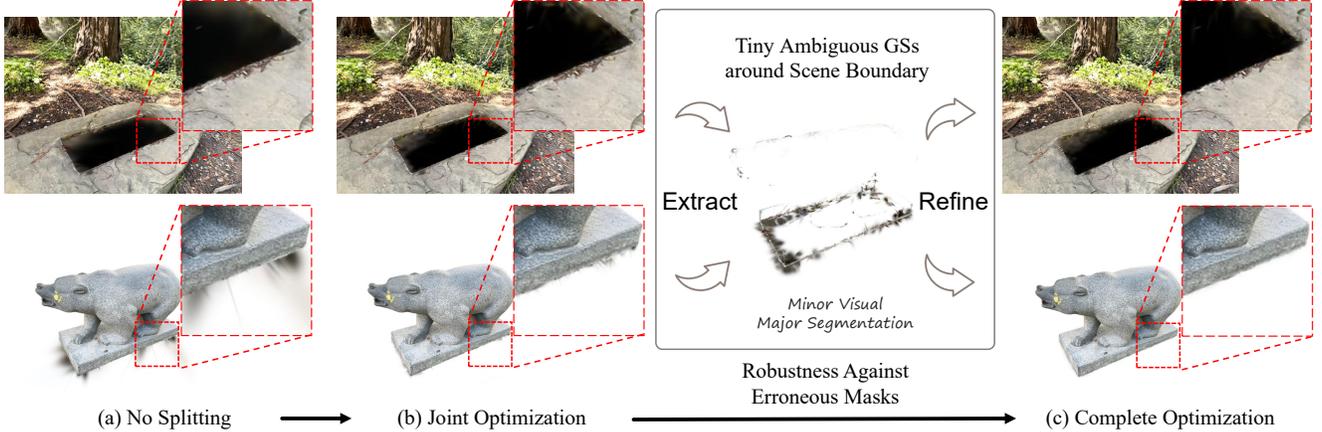


Figure 4. Visualization of different processing phases. (a) Optimizing the mask labels without Gaussian splitting results in unclear boundary segmentation. (b) Jointly optimizing masks and textures with boundary-adaptive Gaussian splitting effectively reduces large ambiguous Gaussians while still leaving tiny ones at the boundaries. (c) Extract and refine tiny ambiguous boundary Gaussians obtained by non-convergent splitting. These tiny Gaussians have little impact on visual quality but affect boundary clarity in 3D segmentation.

number of training views. During this phase, the geometry and texture of the scene are compromised, leading to a significant degradation in visual quality. Therefore, we optimize \mathcal{L}_{rgb} , freezing the mask labels to refine the scene’s geometry and texture details on the accurate boundary structure. The above two stages are performed iteratively and alternately to fuse the information of two modalities to refine 3DGS scene representation, thereby ensuring precise object boundaries and maintaining the visual quality of viewpoint synthesis.

3.2.3. Robustness Against Erroneous Masks

Through alternating optimization of masks and textures, the number of semantically ambiguous Gaussians should progressively decline. In the segmentation phase, we segment the scene based on mask labels. However, experiments reveal that although the quantity of large ambiguous Gaussians reduces, numerous small ambiguous ones remain.

In fact, the binary masks $\{M^v\}$ predicted by the trained 2D vision model exhibit discreteness, which often leads to inaccuracies and inconsistencies in object boundary predictions across different views. This limitation may prevent mask labels of the boundary Gaussians from converging during the optimization process. In contrast to existing methods [36] that roughly address boundary blur caused by inaccurate masks and Gaussian volumes in a combined manner, our approach utilizes the final stage of joint optimization to exclusively enhance robustness against inaccurate masks. Specifically, we identify the tiny boundary Gaussian with ambiguous semantics based on lower values of *mask_sig* and scale *s*, and the visualization result is shown in Figure 4. While these Gaussians minimally affect scene visual quality without 3D segmentation, their removal is essential for achieving clear and complete boundaries for both foreground objects and the background in the 3D segmentation process.

3.2.4. Multi-Object 3D Segmentation

Real 3DGS scene contains multiple objects. Feature-based methods typically require a time-consuming full-scene training process to assign features to each Gaussian for querying. These methods use masks to regularize scene for accurate object boundaries, which limits the granularity of segmentation. For instance, the requirement of clear boundary segmentation of two granularities, “bear’s head” and “bear”, is hard to meet in a trained scene. Similarly, mask-based methods for learning full-scene labels face comparable challenges.

We propose decomposing multi-object segmentation into sequential single-object 3D segmentation. By adding a single integer storage to each Gaussian and utilizing rasterized real-time rendering, we accelerate individual splits, addressing granularity issues. Specifically, for K objects, we define a mask set $\{M^v\}_k$ where $k \in \{1, 2, \dots, K\}$. When optimizing for the k -th object, we perform Gaussian splitting, jointly optimizing with texture to obtain a new 3DGS $\{G_i\}_k$ with clear object boundaries. The optimization process for subsequent objects is conducted on the updated 3DGS $\{G_i\}_k$, iterating until all K objects are optimized. For the specific object to be segmented, a single round of rapid mask label learning is sufficient to achieve clear segmentation results.

3.3. Two-Stage Mask Generation

Mask-based 3DGS segmentation involves generating masks for the target objects based on input images. In our approach, the supervision data consists of V input images $\{I^v\}$ paired with corresponding 2D binary masks $\{M^v\}$. With the development of foundational models such as SAM2 [34], mask prediction across image sequences has become viable, significantly improving inter-frame consistency. However, SAM2 encounters challenges with object continuity over long sequences, potentially failing to infer objects that are heavily

occluded due to discontinuities in visual information.

To address this limitation, we propose a two-stage mask generation approach utilizing text prompts. In the coarse stage, Grounding-DINO [3] is used to extract box prompts from the frame with lower text confidence, which are then applied across the entire sequence for initial mask predictions. In the fine-grained stage, we use Grounding-DINO with higher text confidence to extract box prompts for subsequences where mask prediction was disrupted in the coarse stage. These prompts generate the final mask predictions for the subsequence. See supplementary material for details.

4. Experiments

We demonstrate the effectiveness of our method both quantitatively and qualitatively. For quantification, we utilize the NVOS dataset [35], which is derived from the LLFF dataset and provides ground truth masks with precise object edges. For qualitative evaluation, we employ scenes from various datasets, including LLFF [27], MIP-360 [1], T&T [20], and LERF [18]. These datasets encompass real-world complex scenes, including indoor and outdoor scenes, as well as forward and surrounding scenes. Implementation details and more experiments are provided in supplementary material.

4.1. Quantitative Results

The NVOS dataset contains eight scenes. Each scene includes a reference view and a target view with a clear GT mask. Our method utilizes the annotated masks from the reference view to extract prompt points, which are then passed to the remaining views using SAM2 to generate masks.

Segmentation evaluation. We extract a Gaussian set corresponding to the object and render the 2D mask for the target view. We then calculate the mean IoU and mean accuracy between the GT mask and the rendered mask. Results are shown in Table 1. The NeRF-based method exhibits limited segmentation detail due to inadequate scene representation. In 3DGS segmentation, feature-based approaches result in unsmooth object boundaries because of high-dimensional feature ambiguity. The most similar method to ours is the mask-based 3DGS method. SAGD [16] lacks robustness to erroneous GT masks, leading to small Gaussian artifacts along segmented edges. FlashSplat [36] reduces edge blur but compromises the structural integrity of objects.

Visual evaluation. Due to the lack of reference images for object segmentation results, we use CLIP-IQA [40], a no-reference IQA that evaluates how well an image matches the given text prompt. We set three prompts focusing on boundary quality to comprehensively evaluate the segmentation results, as shown in Table 2. SAGD [16] and FlashSplat [36] process Gaussians roughly, destroying the appearance and sacrificing visual quality for segmentation accuracy. SA3D [5] exhibits obvious ambiguous boundary Gaussians. Ours demonstrates high visual quality across the board.

Table 1. Quantitative segmentation results on NVOS dataset.

Category	Method	mIoU (%)	mAcc (%)
NeRF-based	NVOS [35]	70.1	92.0
	ISRF [14]	83.8	96.4
	SA3D [5]	90.3	98.2
	OmniSeg3D [45]	91.7	98.4
3DGS-based	SAGD [16]	90.4	98.2
	SA3D-GS [5]	90.7	98.3
	SAGA [4]	90.9	98.3
	FlashSplat [36]	91.8	98.6
	COB-GS (ours)	92.1	98.6

Table 2. Quantitative visual results on NVOS dataset.

Method	CLIP-IQA [40] (%) \uparrow		
	<i>Clear / Unclear Boundary</i>	<i>Smooth / Noisy Boundary</i>	<i>Complete / Mutilated Object</i>
SAGD [16]	0.621	0.631	0.788
SA3D-GS [5]	0.658	0.718	0.835
FlashSplat [36]	0.626	0.644	0.829
COB-GS (ours)	0.682	0.731	0.859

4.2. Qualitative Results

We visualize the segmented objects and the background after object removal, including single-object and multi-object segmentation. We compared our results with the current SOTA 3DGS segmentation methods across multiple scenes.

For single-object segmentation, we selected the Truck scene from the T&T dataset and the Kitchen scene from the MIP-360 dataset, as shown in Figure 5. To ensure fair comparison, consistent masks were used across all methods. We rendered the backgrounds and views of the segmented objects on different methods. In contrast, SA3D-GS [5] introduces noticeable Gaussian blur at object boundaries, and its unbalanced loss function degrades the background quality. SAGD [16] leads to tiny edge Gaussians in the foreground. FlashSplat [36] reduces edge Gaussian blur by increasing background bias but indiscriminately removes semantically ambiguous real object regions, such as the truck’s rearview mirror and the Lego bucket in the Truck scene. Our method improves boundary representation accuracy while minimizing background distortion after object removal. Visualization confirms the visual quality assessment results in Table 2.

We further demonstrate the multi-object segmentation capability, using the Figurines scene from the LERF dataset, as illustrated in Figure 6. For segmentation, we employed masks based on text prompts consistent with the mask-based method FlashSplat [36], while the feature-based method SAGA [4] relies on point prompts. Both methods exhibit unclear boundaries and blurred backgrounds. In contrast, our approach achieves significantly clearer delineation of object boundaries while preserving background clarity.

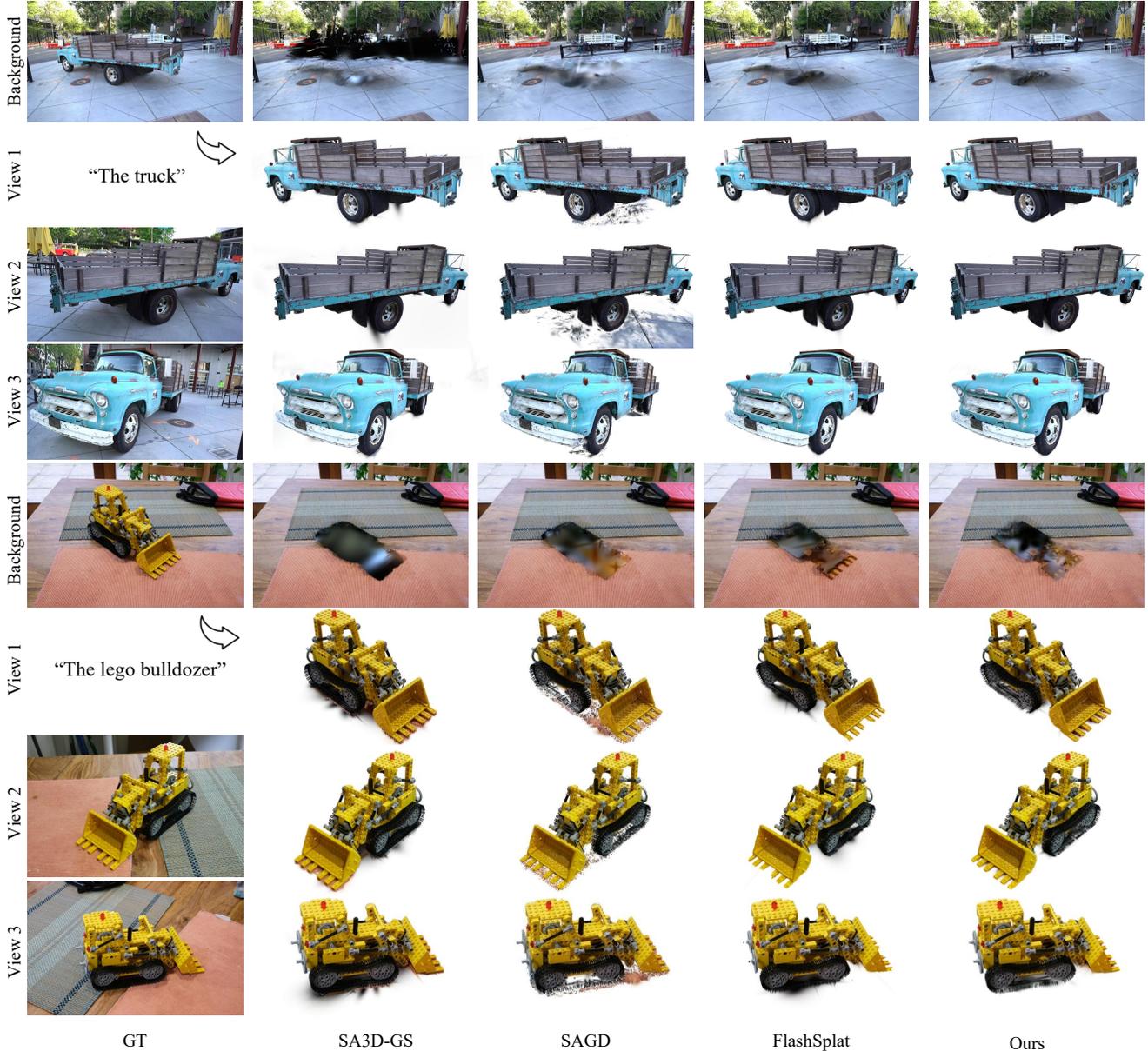


Figure 5. Qualitative result of single-object segmentation. The results show that our method segments the boundaries of the object more clearly, without blurred Gaussians, and the background is cleaner after the object removal.

4.3. Ablation Study

4.3.1. Joint Optimization of Semantics and Textures

An intuitive understanding is that relying solely on Gaussian splitting is insufficient. Although this method enables the Gaussian to adapt to object boundaries, it compromises visual quality. We conducted ablation experiments on the NVOS dataset to explore the relationship between semantics and texture, and the results are shown in Table 3.

When only the boundary-adaptive Gaussian splitting is applied, scene texture quality significantly degrades. In con-

Table 3. Texture quality results on NVOS dataset (PSNR). “Vanilla” indicates the original scene; “M.O” indicates mask-optimized; “T.O” indicates texture-optimized.

Method	Fern	Flower	Fortress	Horns_C	Horns_L	Leaves	Orchids	Trex	Mean
Vanilla	24.26	26.75	29.43	22.24	22.24	15.07	19.82	24.68	23.06
M.O	23.66	26.49	29.05	21.31	22.13	15.05	19.80	23.35	22.61
T.O	24.26	26.69	29.45	22.27	22.27	15.06	19.83	24.71	23.07
M.O+T.O	24.29	26.82	29.48	22.24	22.25	15.09	19.91	24.97	23.13

trast, joint optimization of mask labels and textures improves scene segmentation accuracy while preserving scene quality. For comparison, we performed the same number of iterations

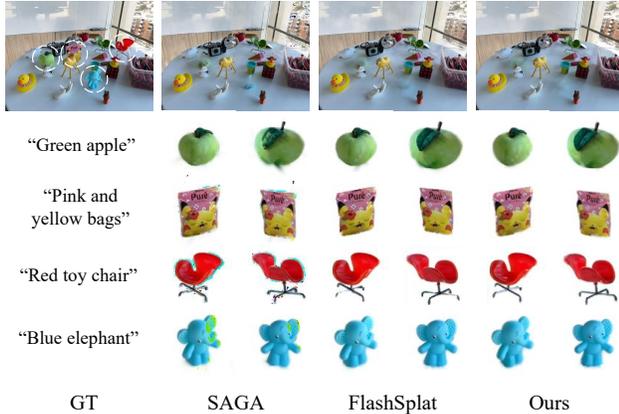


Figure 6. Qualitative result of multi-object segmentation. Our method gets more accurate segmentation results and clearer background quality after object removal compared to contrast method.

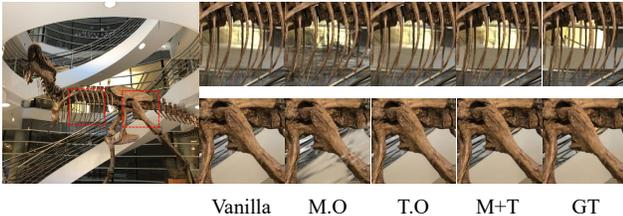


Figure 7. Visualization results in the T rex scene. “Vanilla” indicates the original result; “M.O” indicates mask-optimized; “T.O” indicates texture-optimized; “M+T” indicates M.O and T.O.

Table 4. Ablation results on NVOS dataset. BAGS indicates the boundary-adaptive Gaussian splitting; BGTR indicates the boundary-guided texture restoration; RAEM indicates robustness against erroneous masks.

Component			Performance	
BAGS	BGTR	RAEM	mIoU (%)	mAcc (%)
			91.2	98.3
✓			91.9	98.5
✓	✓		91.9	98.4
✓	✓	✓	92.1	98.6

using texture optimization alone, which was less effective. Visualization results in Figure 7 demonstrate that boundary Gaussian splitting based on mask label statistics can improve texture quality at object boundaries while maintaining independence between the foreground and background.

4.3.2. Robustness Against Erroneous Masks

After joint optimization, tiny ambiguous Gaussians remain along the segmentation boundaries due to predicted erroneous masks, as illustrated in Figure 4. To investigate this, we conducted ablation experiments on the NVOS dataset, as shown in Table 4. The results indicate that without Gaussian splitting, large ambiguous Gaussians yield the lowest

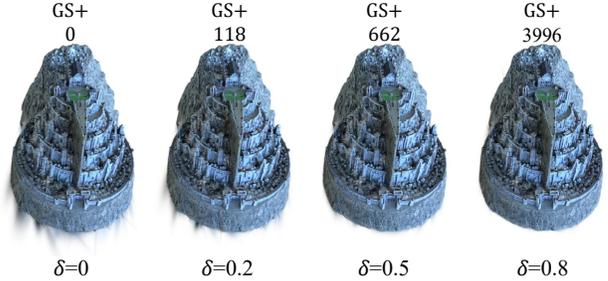


Figure 8. Ablation results in the Fortress scene. As δ increases, the number of ambiguous boundary Gaussians decreases; however, the total number of Gaussians increases rapidly.

metrics. Applying only Gaussian splitting improves segmentation metrics but compromises texture quality. Joint optimization improves texture quality at the expense of slightly decreasing segmentation accuracy. Finally, incorporating robust handling of tiny Gaussians gets optimal performance.

4.3.3. Ambiguity Supervision Threshold

In our method, the parameter δ functions as a threshold to control the number of ambiguous Gaussians. As δ increases, the discrimination against semantically ambiguous Gaussians strengthens, leading to more splits. We conducted an ablation study on δ in the Fortress scene, as shown in Figure 8. We observed that increasing δ reduces ambiguous boundary Gaussians and improves segmentation clarity. However, it also introduces a substantial number of additional Gaussians at the object edge. Therefore, selecting an optimal δ is crucial, with $\delta = 0.5$ proving effective for most scenarios.

5. Conclusion

In this paper, we propose COB-GS, a 3DGS refinement and segmentation approach that clearly segments scene boundaries. COB-GS is innovatively designed to jointly optimize semantics and textures, ensuring they complement each other. This process is supported by a boundary-adaptive Gaussian splitting method. Specifically, in the semantic optimization phase, we utilize semantic gradient statistics to identify and split the ambiguous Gaussians, aligning them with object boundaries. Then we enhance the scene texture along the precisely refined boundary structures. Experimental results demonstrate significant improvements in 3D segmentation performance, particularly in terms of clear object boundaries, accurate textures and robustness against inaccurate masks. In summary, COB-GS is the first 3DGS segmentation method explicitly designed to jointly optimize semantics and texture, ensuring they enhance one another and offering new insights into segmentation with learnable scene representations. Currently, 3D scene reconstruction encounters the challenge of floating artifacts, which are magnified after the segmentation. Future work should focus on effectively eliminating these floating artifacts using semantic information.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. In *ICCV*, pages 5855–5864, 2021. 6, 14
- [2] Yash Bhalgat, Iro Laina, João F Henriques, Andrew Zisserman, and Andrea Vedaldi. N2F2: Hierarchical Scene Understanding with Nested Neural Feature Fields. In *ECCV*, pages 197–214. Springer, 2025. 14
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. In *ICCV*, pages 9650–9660, 2021. 2, 6, 11, 12
- [4] Jiazhong Cen, Jiemin Fang, Chen Yang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Segment Any 3D Gaussians. *arXiv preprint arXiv:2312.00860*, 2023. 2, 6, 14
- [5] Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Chen Yang, Wei Shen, Lingxi Xie, Dongsheng Jiang, Xiaopeng Zhang, and Qi Tian. Segment Anything in 3D with NeRFs. *NeurIPS*, 36: 25971–25990, 2023. 2, 3, 4, 6, 11, 12
- [6] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelSplat: 3D Gaussian Splats from Image Pairs for Scalable Generalizable 3D Reconstruction. In *CVPR*, pages 19457–19467, 2024. 2
- [7] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. GaussianEditor: Swift and Controllable 3D Editing with Gaussian Splatting. In *CVPR*, pages 21476–21485, 2024. 2
- [8] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking Anything with Decoupled Video Segmentation. In *ICCV*, pages 1316–1326, 2023. 12
- [9] Kai Cheng, Xiaoxiao Long, Kaizhi Yang, Yao Yao, Wei Yin, Yuexin Ma, Wenping Wang, and Xuejin Chen. GaussianPro: 3D Gaussian splatting with progressive propagation. In *ICML*, pages 8123–8140. PMLR, 2024. 2
- [10] Seokhun Choi, Hyeonseop Song, Jaechul Kim, Taehyeong Kim, and Hoseok Do. Click-gaussian: Interactive segmentation to any 3d gaussians. In *ECCV*, pages 289–305. Springer, 2025. 2
- [11] Zhiwen Fan, Kevin Wang, Kairun Wen, Zehao Zhu, DeJia Xu, and Zhangyang Wang. LightGaussian: Unbounded 3D Gaussian Compression with 15x Reduction and 200+ FPS. *NeurIPS*, 2024. 2
- [12] Guangchi Fang and Bing Wang. Mini-Splatting: Representing Scenes with a Constrained Number of Gaussians. In *ECCV*, pages 165–181. Springer, 2024. 2
- [13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance Fields Without Neural Networks. In *CVPR*, pages 5501–5510, 2022. 2
- [14] Rahul Goel, Dhawal Sirikonda, Saurabh Saini, and P.J. Narayanan. Interactive Segmentation of Radiance Fields. In *CVPR*, 2023. 2, 6
- [15] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-NeRF2NeRF: Editing 3D Scenes with Instructions. In *ICCV*, pages 19740–19750, 2023. 2, 14
- [16] Xu Hu, Yuxi Wang, Lue Fan, Junsong Fan, Junran Peng, Zhen Lei, Qing Li, and Zhaoxiang Zhang. Semantic Anything in 3D. *arXiv preprint arXiv:2401.17857*, 2024. 2, 4, 6
- [17] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *TOG*, 42(4):1–14, 2023. 1, 2, 3, 11
- [18] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. LERF: Language Embedded Radiance Fields. In *ICCV*, pages 19729–19739, 2023. 2, 6, 12
- [19] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment Anything. *arXiv preprint arXiv:2304.02643*, 2023. 2, 11
- [20] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: benchmarking large-scale scene reconstruction. *TOG*, 36(4), 2017. 6
- [21] Joo Chan Lee, Daniel Rho, Xiangyu Sun, Jong Hwan Ko, and Eunbyung Park. Compact 3D Gaussian Representation for Radiance Field. In *CVPR*, pages 21719–21728, 2024. 2
- [22] Yiqing Liang, Numair Khan, Zhengqin Li, Thu Nguyen-Phuoc, Douglas Lanman, James Tompkin, and Lei Xiao. Gafre: Gaussian deformation fields for real-time dynamic novel view synthesis. *arXiv preprint arXiv:2312.11458*, 2023. 2
- [23] 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 *CVPR*, pages 6517–6526, 2024. 2
- [24] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. Weakly Supervised 3D Open-vocabulary Segmentation. *NeurIPS*, 36:53433–53456, 2023. 2
- [25] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *CVPR*, pages 20654–20664, 2024. 2
- [26] Mallick and Goel, Bernhard Kerbl, Francisco Vicente Carasco, Markus Steinberger, and Fernando De La Torre. Taming 3DGS: High-Quality Radiance Fields with Limited Resources. In *SIGGRAPH Asia 2024 Conference Papers*, 2024. 2
- [27] 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. *TOG*, 38(4), 2019. 6, 14
- [28] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*, 2020. 1, 2, 14

- [29] 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 *CVPR*, pages 20669–20679, 2023. [2](#)
- [30] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. LangSplat: 3D Language Gaussian Splatting. In *CVPR*, pages 20051–20060, 2024. [2](#), [12](#)
- [31] Ri-Zhao Qiu, Ge Yang, Weijia Zeng, and Xiaolong Wang. Language-Driven Physics-Based Scene Synthesis and Editing via Feature Splatting. In *ECCV*, 2024. [2](#)
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*, pages 8748–8763. PMLR, 2021. [2](#)
- [33] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. GARField: Group Anything with Radiance Fields. *arXiv preprint arXiv:2401.09419*, 2024. [2](#)
- [34] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. SAM 2: Segment Anything in Images and Videos. *arXiv preprint arXiv:2408.00714*, 2024. [2](#), [5](#), [11](#)
- [35] Zhongzheng Ren, Aseem Agarwala[†], Bryan Russell[†], Alexander G. Schwing[†], and Oliver Wang[†]. Neural Volumetric Object Selection. In *CVPR*, pages 4201–4211, 2022. [6](#)
- [36] Qihong Shen, Xingyi Yang, and Xinchao Wang. FlashSplat: 2D to 3D Gaussian Splatting Segmentation Solved Optimally. In *ECCV*, pages 456–472. Springer, 2024. [2](#), [4](#), [5](#), [6](#), [14](#)
- [37] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. In *CVPR*, pages 5333–5343, 2024. [2](#)
- [38] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct Voxel Grid Optimization: Super-Fast Convergence for Radiance Fields Reconstruction. In *CVPR*, pages 5459–5469, 2022. [2](#)
- [39] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. [2](#)
- [40] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring CLIP for Assessing the Look and Feel of Images. In *AAAI*, 2023. [6](#)
- [41] 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 *CVPR*, pages 20310–20320, 2024. [2](#)
- [42] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. PhysGaussian: Physics-Integrated 3D Gaussians for Generative Dynamics. In *CVPR*, pages 4389–4398, 2024. [2](#)
- [43] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth Anything: Unleashing the Power of Large-Scale Unlabeled Data. In *CVPR*, pages 10371–10381, 2024. [2](#)
- [44] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian Grouping: Segment and Edit Anything in 3D Scenes. In *ECCV*, 2024. [2](#), [12](#), [14](#)
- [45] Haiyang Ying, Yixuan Yin, Jinzhi Zhang, Fan Wang, Tao Yu, Ruqi Huang, and Lu Fang. OmniSeg3D: Omniversal 3D Segmentation via Hierarchical Contrastive Learning. In *CVPR*, pages 20612–20622, 2024. [6](#)
- [46] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-Splatting: Alias-free 3D Gaussian Splatting. In *CVPR*, pages 19447–19456, 2024. [2](#)
- [47] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J. Davison. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *ICCV*, pages 15838–15847, 2021. [2](#)
- [48] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejie Xu, Pradyumna Chari, Suyu You, Zhangyang Wang, and Achuta Kadambi. Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields. In *CVPR*, pages 21676–21685, 2024. [2](#)
- [49] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. FSGS: Real-Time Few-Shot View Synthesis using Gaussian Splatting. In *ECCV*, pages 145–163. Springer, 2024. [2](#)

COB-GS: Clear Object Boundaries in 3DGS Segmentation Based on Boundary-Adaptive Gaussian Splitting

Supplementary Material

In the supplementary material, we first introduce in detail our proposed two-stage mask generation based on text prompts in Sec. 6. Next, we present the concrete training strategy and implementation details of COB-GS in Sec. 7. In Sec. 8 and Sec. 9, we evaluate the open-vocabulary segmentation capability and the computational cost of COB-GS, respectively. In Sec. 10 and Sec. 11, we discussed the differences between mask-based and feature-based methods in terms of application scenarios and multi-granularity segmentation. Finally, additional visualizations of the segmentation results are presented in Sec. 12.

6. Two-Stage Mask Generation Based on Text Prompts

Mask-based 3D segmentation requires generating a set of masks for regions of interest from a collection of input images. Thus, the supervision data consists of V input views $\{I^v\}$ corresponding to 2D binary masks $\{M^v\}$. Each mask $M \in \mathbb{R}^{H \times W}$ contains discrete values of 0 and 1. The related work SA3D [5] improves optimization efficiency and mask view consistency by using Segment Anything Model (SAM) [19] to iteratively generate the mask for each frame. With the emergence of foundational models like SAM2 [34], mask prediction across video sequences has become feasible.

SAM2 retains the encoder-decoder structure of SAM, where the encoder S_e takes an image I as input. Unlike SAM, SAM2 employs memory attention S_m to utilize past frame features f_m as conditions for generating the current frame embedding e_I :

$$e_I = S_m(S_e(I), f_m) \quad (1)$$

The past frame features f_m are maintained in a FIFO memory queue. The decoder takes the current frame embedding e_I and the prompts P as input, outputting the corresponding 2D binary mask M :

$$M = S_d(e_I, P) \quad (2)$$

The prompts P include masks, boxes, points, or texts. The memory capability of SAM2 allows it to handle mask prediction for video sequences, which aligns with the input view conditions $\{I^v\}$ for the 3DGS task. However, when SAM2 performs mask prediction across video sequences, it encounters challenges with object continuity; specifically, it may fail to recognize severely occluded objects due to information discontinuity. To address this issue, we propose a two-stage mask generation method based on text prompts.

Algorithm 1 Two-stage mask generation

Input: Frame index idx , text prompt $text$, image set I , high confidence C_{high} , low confidence C_{low}
Result: Updated dictionary $video_segments$
Initialize dictionary $valid_idxs \leftarrow \{\}$
Initialize dictionary $video_segments \leftarrow \{\}$
SAM2.init_state(I)
 $image \leftarrow I[idx]$
 $boxes \leftarrow \text{Grounding DINO}(text, image, C_{low})$
SAM2.add_new_box($idx, boxes$)
for each frame $i, mask$ in SAM2(idx) **do**
 $video_segments[i] \leftarrow mask$
 $valid_idxs[i] \leftarrow$ if $mask$ is empty then 0 else 1
end for
for each key in $valid_idxs$ **do**
 if $valid_idxs[key] = 0$ **then**
 $boxes \leftarrow \text{Grounding DINO}(text, I[key], C_{high})$
 if $boxes$ is empty **then**
 continue
 end if
 SAM2.add_new_box($idx, boxes$)
 $max_sk \leftarrow \text{FindMaxSub}(valid_idxs, key)$
 for each frame $j, mask$ in SAM2(key, max_sk) **do**
 $video_segments[j] \leftarrow mask$
 $valid_idxs[j] \leftarrow$ if $mask$ is empty then 0 else 1
 end for
 end if
end for

In the coarse mask generation stage, we utilize Grounding DINO [3] to extract box prompts from the given prompt frame with lower text confidence, which are then used for full-sequence mask prediction to obtain preliminary results. In the fine-grained stage, we leverage Grounding DINO with higher text confidence to extract box prompts for subsequences within the original sequence that lack mask prediction results, which are then used for subsequence mask prediction. See Algorithm 1 for details.

7. Implementation Details

Our method is a post-processing method based on the original 3D Gaussian Splatting [17]. For each scene, we perform 30,000 iterations of training according to the parameters set by the original 3DGS to obtain the original 3DGS scene. COB-GS mainly consists of two components: optimization

Table 5. Results on LERF-mask dataset.

Method	Figurines		Ramen		Teatime	
	mIoU (%)	mBIoU (%)	mIoU (%)	mBIoU (%)	mIoU (%)	mBIoU (%)
DEVA [8]	46.2	45.1	56.8	51.1	54.3	52.2
LERF [18]	33.5	30.6	28.3	14.7	49.7	42.6
SA3D [5]	24.9	23.9	7.4	7.0	42.5	39.2
LangSplat [30]	52.8	50.5	50.4	44.7	69.5	65.6
Gaussian Grouping [44]	69.7	67.9	77.0	68.7	71.7	66.1
COB-GS (ours)	76.3	73.9	78.1	69.2	77.2	72.8



Figure 9. Visualization of the LERF-mask dataset [44]. The result of the segmentation is obtained under the specified text prompt.

process and robustness process. The optimization process involves alternating between mask optimization and texture optimization. For the mask optimization stage, we optimize the mask labels and perform Gaussian splitting. The learning rate of the mask labels is set to 0.1. For the texture optimization stage, we optimize the geometry and texture, and the learning rate of appearance follows the original 3DGS setting. Each stage is trained for $2 \times V$ iterations, where V is the number of input images. Two sets of hyperparameters are used for different scene types: for forward scenes, we set $\delta = 0.5$ and perform a total of $22 \times V$ iterations of alternating optimization; for surrounding scenes, we set $\delta = 0.8$ and conduct $14 \times V$ iterations of alternating optimization. The robustness process follows scene optimization and involves

extracting and refining ambiguous boundary Gaussians at scales smaller than the pixel scale. In our two-stage mask generation method, we utilize the SAM2 hiera.l model and the Grounding DINO swinb model. All experiments were conducted on a single NVIDIA RTX 3090 GPU.

8. Open-Vocabulary 3D Segmentation

To achieve open-vocabulary semantic segmentation, we follow the setup of existing methods [5, 44] and utilize Grounding DINO [3] to generate boxes for input images, similar to the approach in Sec. 6. We compare our method with the current state-of-the-art methods for open-vocabulary 3D segmentation using the LERF-mask dataset, which is annotated from test views of three 3D scenes in the LERF dataset.



Figure 10. Visualization of 3DGS segmentation. We utilize text prompts to obtain object masks and perform 3D segmentation across multiple scenes, including Horns, Orchids and Fortress from the LLFF dataset, Garden from MIP-360, Bear from the IN2N dataset, and Pinecone from NeRF.

The scenes contain severe object occlusions, and the mask boundaries of the test views are more complex. As shown in Table 5, our method demonstrates a clear advantage over current SOTA methods. Visual segmentation comparisons in Figure 9 reveal that our method provides more accurate segmentation predictions with clear boundaries, while Gaussian Grouping [44] exhibits blurriness in segmentation results.

9. Computation Cost

We evaluate the computational efficiency of COB-GS in comparison to state-of-the-art 3DGS segmentation methods, namely the feature-based SAGA [4] and the mask-based FlashSplat [36]. This evaluation is conducted on the Fortress scene ($V = 42$) from the LLFF dataset [27] using a single NVIDIA RTX 3090 GPU, with results presented in Table 6. We provide the total time cost (prep time+opt time+seg time) and the maximum VRAM of the entire reconstruction and segmentation pipeline. SAGA [4] requires 10,000 iterations of gradient descent to distill 2D masks into object features associated with each 3D Gaussian, resulting in substantial additional training time for scene optimization. Moreover, object segmentation remains time-consuming due to the need for network inference. FlashSplat [36] does not offer a mask extraction method, and assigning labels to each Gaussian through forward rendering is relatively time-consuming. In contrast, our extraction process relies entirely on inverse rendering, which ensures that texture optimization simultaneously optimizes scene labels. The optimization time is comparable to the speed of FlashSplat, and segmentation requires only filtering the labels.

Method	Prep Time	Opt Time	Seg Time	Total Time	Mem
SAGA [4]	145 s	20 min	200 ms	22.42 min	7.6 G
FlashSplat [36]	N/A	24 s	10 ms	N/A	2.4 G
COB-GS	4 s	24 s	8 ms	0.46 min	2.7 G

Table 6. Computation cost comparisons over the Fortress scene.

10. Application Scenarios

Tab. 6 shows that feature-based methods like SAGA [4] consume more time and memory than mask-based methods. This is because feature-based methods optimize high-dimensional features for the entire scene, while our approach focuses solely on optimizing labels. This difference is evident in the optimization time. However, SAGA [4] has the advantage of allowing multiple segmentations with a single training session, making it suitable for offline fixed scenes that require frequent segmentations, despite its high equipment demands. In contrast, mask-based methods assign labels directly on the reconstructed 3DGS, offering faster single segmentation. COB-GS achieves better visual quality

in a negligible time and suits edge applications that require fast and detailed single object segmentation. Furthermore, progress in foundational models and faster 3DGS training will further promote COB-GS’s real-time applicability.

11. Multi-granularity Segmentation

Unlike the feature-based method N2F2 [2], which requires a time-consuming process to integrate multi-granular high-dimensional features for scene reconstruction, the mask-based COB-GS decouples the processes of reconstruction and segmentation. This separation allows for greater flexibility, enabling arbitrary granularity to be finely achieved during the mask generation phase. As a result, COB-GS can obtain the region of interest efficiently within the preparation time of a single fast segmentation, streamlining the workflow and enhancing performance without the heavy computational burden associated with integrating complex features.

12. More Qualitative Results

To demonstrate the effectiveness of our proposed 3D segmentation method in producing clear object boundaries, we provide visualizations of 3D segmentation across multiple scenes, including the Horns, Orchids and Fortress from the LLFF dataset [27], the Garden from MIP-360 [1], the Bear from the IN2N dataset [15], and the Pinecone from NeRF [28], encompassing both forward and surrounding scenes. We obtain masks using text prompts, as described in Sec. 6. The results shown in Figure 10 clearly demonstrate that the object edges in our 3D segmentation results are very clear, while also maintaining high-quality textures for both the foreground and background.