# Segment Any 3D Gaussians

**Jiazhong Cen[1], Jiemin Fang[2], Chen Yang[1], Lingxi Xie[2], Xiaopeng Zhang[2], Wei Shen[1*], Qi Tian[2]**

[1]MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
[2]Huawei Technologies Co., Ltd.
{jiazhongcen, ycyangchen, wei.shen}@sjtu.edu.cn,
{jaminfong, 198808xc, zxphistory}@gmail.com, tian.qi1@huawei.com

## Abstract

This paper presents SAGA (Segment Any 3D GAussians), a highly efficient 3D promptable segmentation method based on 3D Gaussian Splatting (3D-GS). Given 2D visual prompts as input, SAGA can segment the corresponding 3D target represented by 3D Gaussians within **4 ms**. This is achieved by attaching an scale-gated affinity feature to each 3D Gaussian to endow it a new property towards multi-granularity segmentation. Specifically, a scale-aware contrastive training strategy is proposed for the scale-gated affinity feature learning. It 1) distills the segmentation capability of the Segment Anything Model (SAM) from 2D masks into the affinity features and 2) employs a soft scale-gate mechanism to deal with multi-granularity ambiguity in 3D segmentation through adjusting the magnitude of each feature channel according to a specified 3D physical scale. Evaluations demonstrate that SAGA achieves real-time multi-granularity segmentation with quality comparable to state-of-the-art methods. As one of the first methods addressing promptable segmentation in 3D-GS, the simplicity and effectiveness of SAGA pave the way for future advancements in this field.

**Code** — https://github.com/Jumpat/SegAnyGAussians

## 1   Introduction

Promptable segmentation has attracted increasing attention and has seen significant advancements, particularly with the development of 2D segmentation foundation models such as the Segment Anything Model (SAM) (Kirillov et al. 2023). However, 3D promptable segmentation remains relatively unexplored due to the scarcity of 3D data and the high cost of annotation. To address these challenges, many studies (Cen et al. 2023; Chen et al. 2023; Ying et al. 2024; Kim et al. 2024; Fan et al. 2023; Liu et al. 2024b) have proposed to extend SAM's 2D segmentation capabilities to 3D using radiance fields, achieving notable success.

In this paper, we focus on promptable segmentation in 3D Gaussian Splatting (3D-GS) (Kerbl et al. 2023), which represents a significant milestone in radiance fields research due to its superior rendering quality and efficiency compared to its predecessors. We highlight that, in contrast to previous

radiance fields, the explicit 3D Gaussian structure is an ideal carrier for 3D segmentation, as segmentation capabilities can be integrated into 3D-GS as an intrinsic attribute, without necessitating an additional bulky segmentation module.

Accordingly, we propose SAGA (Segment Any 3D GAussians), a 3D promptable segmentation method that integrates the segmentation capabilities of SAM into 3D-GS seamlessly. SAGA takes 2D visual prompts as input and outputs the corresponding 3D target represented by 3D Gaussians. To achieve this purpose, two primary challenges are faced. First, SAGA should figure out an efficient way to endow each 3D Gaussian with the ability of 3D segmentation, so that the high efficiency of 3D-GS can be preserved. Second, as a robust promptable segmentation method, SAGA must effectively address multi-granularity ambiguity, where a single 3D Gaussian may belong to different parts or objects at varying levels of granularity.

To address the two challenges, SAGA respectively introduces two solutions. First, SAGA attaches an affinity feature to each 3D Gaussian in a scene to endow it with a new property towards segmentation. The similarity between two affinity features indicates whether the corresponding 3D Gaussians belong to the same 3D target. Second, inspired by GARField (Kim et al. 2024), SAGA employs a soft scale-gate mechanism to handle multi-granularity ambiguity. Depending on a specified 3D physical scale, the scale gate adjusts the magnitude of each feature channel. This mechanism maps the Gaussian affinity features into different sub-spaces for various scales, thereby preserving the multi-granularity information and meanwhile mitigating the distraction in feature learning brought by multi-granularity ambiguity. To realize the two solutions, SAGA proposes a scale-aware contrastive training strategy, which distills the segmentation capability of SAM from 2D masks into scale-gated affinity features. This strategy determines the correlations between a pair of pixels within an image based on 3D scales. These correlations are then used to supervise the rendered affinity features through a correspondence distillation loss. The correlation information is transmitted to the Gaussian affinity features via backpropagation facilitated by the differentiable rasterization algorithm. After training, SAGA achieves real-time multi-granularity segmentation precisely.

---

Figure 1: SAGA performs promptable multi-granularity segmentation within **milliseconds**. Prompts are marked by points.

## 2 Related Work

**2D Promptable Segmentation** The task of 2D promptable segmentation is proposed by Kirillov et al. (2023), which aims to return segmentation masks given input prompts that specify the segmentation target in an image. To address this problem, they introduce the Segment Anything Model (SAM), a groundbreaking segmentation foundation model. A similar model to SAM is SEEM (Zou et al. 2023), which also achieves competitive performance. Prior to these models, the most closely related task to promptable 2D segmentation is interactive image segmentation (Boykov and Jolly 2001; Grady 2006; Gulshan et al. 2010; Rother, Kolmogorov, and Blake 2004; Chen et al. 2022b; Sofiiuk, Petrov, and Konushin 2022; Liu et al. 2023b). Inspired by the success of SAM, many studies (Yang et al. 2023; Xu et al. 2023; Guo et al. 2024; Yin et al. 2024) proposed to use SAM for 3D segmentation. Different from them, we focus on lifting the ability of SAM to 3D via 3D-GS.

**3D Segmentation in Radiance Fields** With the success of radiance fields (Mildenhall et al. 2020; Sun, Sun, and Chen 2022; Chen et al. 2022a; Barron et al. 2022; Müller et al. 2022; Hedman et al. 2024; Fridovich-Keil et al. 2022; Wizadwongsa et al. 2021; Lindell, Martel, and Wetzstein 2021; Fang et al. 2022), numerous studies have explored 3D segmentation within them. Zhi et al. (2021) proposed Semantic-NeRF, demonstrating the potential of Neural Radiance Field (NeRF) in semantic propagation and refinement. NVOS (Ren et al. 2022) introduced an interactive approach to select 3D objects from NeRF by training a lightweight MLP using custom-designed 3D features. By using 2D self-supervised models, approaches like N3F (Tschernezki et al. 2022), DFF (Kobayashi, Matsumoto, and Sitzmann 2022), and ISRF (Goel et al. 2023) aim to elevate 2D

visual features to 3D by training additional feature fields that can output 2D feature maps, imitating the original 2D features from different views. NeRF-SOS (Fan et al. 2023) and ContrastiveLift (Bhalgat et al. 2023) distill the 2D feature similarities into 3D features. There are also some other instance segmentation and semantic segmentation approaches (Stelzner, Kersting, and Kosiorek 2021; Niemeyer and Geiger 2021; Yu, Guibas, and Wu 2022; Liu et al. 2022, 2023c; Bing, Chen, and Yang 2023; Fu et al. 2022; Vora et al. 2022; Siddiqui et al. 2023) for radiance fields. Combined with CLIP (Radford et al. 2021), some approaches (Kerr et al. 2023; Liu et al. 2023a; Bhalgat et al. 2024; Qin et al. 2024) proposed to conduct open-vocabulary 3D segmentation in radiance fields. With the popularity of SAM, a stream of studies (Kim et al. 2024; Ying et al. 2024; Ye et al. 2024; Cen et al. 2023; Lyu et al. 2024) proposed lifting the segmentation ability of SAM to 3D with radiance fields. SA3D (Cen et al. 2023) adopts an iterative pipeline to refine the 3D mask grids with SAM. GaussianGrouping (Ye et al. 2024) uses video tracking technology to align the inconsistent 2D masks extracted by SAM across different views and assigns labels to 3D Gaussians in a 3D-GS model with the aligned masks. OmniSeg3D (Ying et al. 2024) employs a hierarchical contrastive learning method to automatically learn segmentation from multi-view 2D masks extracted by SAM.

The approach most closely related to SAGA is GARField (Kim et al. 2024), which addresses multigranularity ambiguity in 3D segmentation using 3D physical scale, inspiring SAGA's scale gate mechanism. However, GARField's reliance on implicit feature fields for outputting 3D features requires repeated queries for segmentation at different scales, reducing efficiency. In contrast, the scalegate mechanism of SAGA enhances efficiency by integrating directly with 3D-GS without additional computation.

# 3 Method

In this section, we first give a brief review of 3D Gaussian Splatting (3D-GS) (Kerbl et al. 2023) and the scale-conditioned 3D features (Kerr et al. 2023; Kim et al. 2024). Then we introduce the overall pipeline of SAGA, followed by explanation of the scale-gated Gaussian affinity features and the scale-aware contrastive learning.

## 3.1 Preliminary

**3D Gaussian Splatting (3D-GS)**  Given a training dataset $\mathcal{I}$ of multi-view 2D images with camera poses, 3D-GS learns a set of 3D colored Gaussians $\mathcal{G} = \{\mathbf{g}_1, \mathbf{g}_2, ..., \mathbf{g}_N\}$, where $N$ denotes the number of 3D Gaussians in the scene. The mean of a Gaussian represents its position and the covariance indicates its scale. Accordingly, 3D-GS proposes a novel differentiable rasterization technology for efficient training and rendering. Given a specific camera pose, 3D-GS projects the 3D Gaussians to 2D and computes the color $\mathbf{C}(\mathbf{p})$ of a pixel $\mathbf{p}$ by blending a set of ordered Gaussians $\mathcal{G}_{\mathbf{p}}$ overlapping the pixel. Let $\mathbf{g}_i^{\mathbf{P}}$ denote the i-th Gaussian in $\mathcal{G}_{\mathbf{p}}$, this process is formulated as:

$$\mathbf{C}(\mathbf{p}) = \sum_{i=1}^{|\mathcal{G}_{\mathbf{P}}|} \mathbf{c}_{\mathbf{g}_i^{\mathbf{P}}} \alpha_{\mathbf{g}_i^{\mathbf{P}}} \prod_{j=1}^{i-1}(1 - \alpha_{\mathbf{g}_j^{\mathbf{P}}}), \qquad (1)$$

where $\mathbf{c}_{\mathbf{g}_i^{\mathbf{P}}}$ is the color of $\mathbf{g}_i^{\mathbf{P}}$ and $\alpha_{\mathbf{g}_i^{\mathbf{P}}}$ is given by evaluating the corresponding 2D Gaussian with covariance $\Sigma$ multiplied with a learned per-Gaussian opacity.

**Scale-Conditioned 3D Feature**  LERF (Kerr et al. 2023) first proposes the concept of a scale-conditioned feature field for learning from global image embeddings obtained from CLIP. GARField (Kim et al. 2024) then introduces it into the area of radiance field segmentation to tackle the multi-granularity ambiguity. To compute the 3D mask scale $s_{\mathbf{M}}$ of a 2D mask $\mathbf{M}$, GARField projects $\mathbf{M}$ into 3D space with the camera intrinsic parameters and depth information predicted by a pre-trained radiance field. Let $\mathcal{P}$ denote the obtained point cloud, $\mathcal{X}(\mathcal{P}), \mathcal{Y}(\mathcal{P}), \mathcal{Z}(\mathcal{P})$ denote the set of 3D coordinate components of $\mathcal{P}$, the mask scale $s_{\mathbf{M}}$ is:

$$s_{\mathbf{M}} = 2\sqrt{\texttt{std}(\mathcal{X}(\mathcal{P}))^2 + \texttt{std}(\mathcal{Y}(\mathcal{P}))^2 + \texttt{std}(\mathcal{Z}(\mathcal{P}))^2}, \qquad (2)$$

where $\texttt{std}(\cdot)$ denotes the standard variation of a set of scalars. Since these scales are computed in 3D space, they are generally consistent across different views. SAGA uses the 3D scales for multi-granularity segmentation but realizes it in a more efficient way.

## 3.2 Overall Pipeline

The main components of SAGA are shown in Figure 2. Given a pre-trained 3D-GS model $\mathcal{G}$, SAGA attaches a Gaussian affinity feature $\mathbf{f}_{\mathbf{g}} \in \mathbb{R}^D$ for each 3D Gaussian $\mathbf{g}$ in $\mathcal{G}$. $D$ denotes the feature dimension. To handle the inherent multi-granularity ambiguity of 3D promptable segmentation, SAGA employs a soft scale gate mechanism to project these features into different scale-gated feature subspaces for various scales $s$.

To train the affinity features, SAGA extracts a set of multi-granularity masks $\mathcal{M}_{\mathbf{I}} = \{\mathbf{M}_{\mathbf{I}}^i \in \{0, 1\}^{HW} | i = 1, ..., N_{\mathbf{I}}\}$ for each image $\mathbf{I}$ in the training set $\mathcal{I}$ with SAM. $H, W$ denotes the height and width of $\mathbf{I}$ respectively. $N_{\mathbf{I}}$ is the number of extracted masks. For each mask $\mathbf{M}_{\mathbf{I}}^i$, its 3D physical scale $s_{\mathbf{M}_{\mathbf{I}}^i}$ is calculated using the depth predicted by $\mathcal{G}$ with the camera pose, as shown in Equation (2). Subsequently, SAGA employs a scale-aware contrastive learning strategy (Section 3.4) to distill the multi-granularity segmentation ability embedded in multi-view 2D masks into the scale-gated affinity features. After training, at given scales, the affinity feature similarities between two Gaussians indicate whether they belong to the same 3D target.

During inference (Section 3.5), given a specific viewpoint, SAGA converts 2D visual prompts (points with scales) into corresponding 3D scale-gated query features to segment the 3D target by evaluating feature similarities with 3D affinity features. Additionally, with well-trained affinity features, 3D scene decomposition is achievable through simple clustering. Moreover, by integrating with CLIP, SAGA can perform open-vocabulary segmentation (see Section A.3 of the supplement) without requiring language fields.

## 3.3 Gaussian Affinity Feature

At the core of SAGA is the Gaussian affinity features $\mathcal{F} = \{\mathbf{f}_{\mathbf{g}} \mid \mathbf{g} \in \mathcal{G}\}$, which are learned from the multi-view 2D masks extracted by SAM. To tackle the inherent multi-granularity ambiguity in promptable segmentation, inspired by GARfield (Kim et al. 2024), we introduce a scale-gate mechanism to split the feature space into different subspaces for various 3D physical scales. Then, a 3D Gaussian can belong to different segmentation targets at different granularities without conflict.

**Scale-Gated Affinity Features**  Given a Gaussian affinity feature $\mathbf{f}_{\mathbf{g}}$ and a specific scale $s$, SAGA adopts a scale gate to adapt the magnitude of different feature channels accordingly. The scale gate is defined as a mapping $\mathcal{S} : [0, 1] \rightarrow [0, 1]^D$, which projects a scale scalar $s \in [0, 1]$ to its corresponding soft gate vector $\mathcal{S}(s)$. To maximize segmentation efficiency, the scale gate adopts a extremely streamlined design, which is a single linear layer followed by a sigmoid function. At the scale of $s$ the scale-gated affinity feature is:

$$\mathbf{f}_{\mathbf{g}}^s = \mathcal{S}(s) \odot \mathbf{f}_{\mathbf{g}}, \qquad (3)$$

where $\odot$ denotes the Hadamard product. Thanks to the simplicity of the scale-gate mechanism, the time overhead caused by scale changing is negligible.

Since all Gaussian affinity features share a common scale gate at scale $s$, during training, we can first render the affinity features to 2D and then apply the scale gate to the 2D rendered features, i.e.,

$$\mathbf{F}(\mathbf{p}) = \sum_{i=1}^{|\mathcal{G}_{\mathbf{P}}|} \mathbf{f}_{\mathbf{g}_i^{\mathbf{P}}} \alpha_{\mathbf{g}_i^{\mathbf{P}}} \prod_{j=1}^{i-1}(1 - \alpha_{\mathbf{g}_j^{\mathbf{P}}}), \qquad (4)$$

$$\mathbf{F}^s(\mathbf{p}) = \mathcal{S}(s) \odot \mathbf{F}(\mathbf{p}). \qquad (5)$$

During inference, the scale gate is directly applied to the 3D Gaussian affinity features for conducting 3D segmentation.
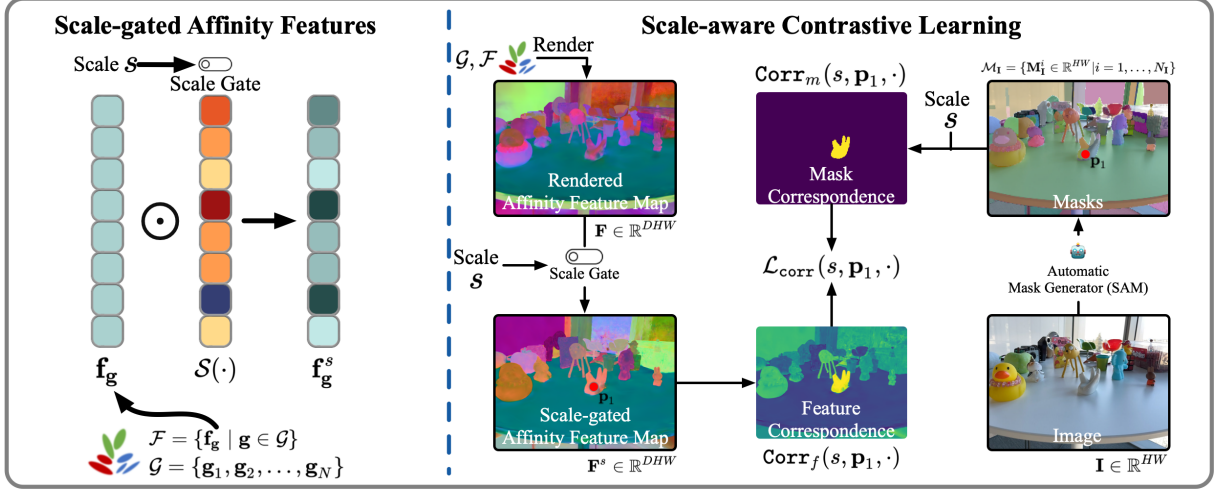
Figure 2: The architecture of SAGA. *Left*: SAGA attaches a Gaussian affinity feature to each 3D Gaussian. The magnitude of different affinity feature channels are adjusted by a soft scale gate to handle multi-granularity ambiguity. *Right*: SAGA distills segmentation ability of SAM into affinity features attached to 3D Gaussians in the 3D-GS model through scale-aware contrastive learning.

**Local Feature Smoothing**    In practice, we find that there are many noisy Gaussians in the 3D space that exhibit unexpectedly high feature similarities with the segmentation target. This may occur for various reasons, such as insufficient training due to small weights in rasterization or incorrect geometry structure learned by 3D-GS. To tackle this problem, we adopt the spatial locality prior of 3D Gaussians. During training, SAGA uses the smoothed affinity feature of a Gaussian $\mathbf{g}$ to replace its original feature $\mathbf{f_g}$, *i.e.*, $\mathbf{f_g} \leftarrow \frac{1}{K} \sum_{\mathbf{g}' \in \text{KNN}(\mathbf{g})} \mathbf{f_{g'}}$. $\text{KNN}(\mathbf{g})$ denotes K-nearest neighbors of $\mathbf{g}$. After training, the affinity feature for each 3D Gaussian is saved as its smoothed feature.

### 3.4    Scale-Aware Contrastive Learning

As introduced in Section 3.3, each 3D Gaussian $\mathbf{g}$ is assigned with an affinity feature $\mathbf{f_g}$. To train these features, we employ a scale-aware contrastive learning strategy to distill the pixel-wise correlation information from 2D masks into 3D Gaussians via the differentiable rasterization.

**Scale-Aware Pixel Identity Vector**    To conduct scale-aware contrastive learning, for an image $\mathbf{I}$, we first convert the automatically extracted 2D masks $\mathcal{M}_\mathbf{I}$ to scale-aware supervision signal. For this purpose, we assign a scale-aware pixel identity vector $\mathbf{V}(s, \mathbf{p}) \in \{0, 1\}^{N_\mathbf{I}}$ to each pixel $\mathbf{p}$ in $\mathbf{I}$. The identity vectors reflect the 2D masks that a pixel belong to at specific scales. If two pixels $\mathbf{p}_1, \mathbf{p}_2$ share at least a same mask at a given scale (*i.e.*, $\mathbf{V}(s, \mathbf{p}_1) \cdot \mathbf{V}(s, \mathbf{p}_2) > 0$), they should have similar features at scale $s$.

To obtain $\mathbf{V}(s, \mathbf{p})$, we first sort the mask set $\mathcal{M}_\mathbf{I}$ in descending order according to their mask scales and get an ordered mask list $\mathcal{O}_\mathbf{I} = (\mathbf{M}_\mathbf{I}^{(1)}, ..., \mathbf{M}_\mathbf{I}^{(N_\mathbf{I})})$, where $s_{\mathbf{M}_\mathbf{I}^{(1)}} > ... > s_{\mathbf{M}_\mathbf{I}^{(N_\mathbf{I})}}$. Then, for a pixel $\mathbf{p}$, when $s_{\mathbf{M}_\mathbf{I}^{(i)}} < s$, the i-th entry of $\mathbf{V}(s, \mathbf{p})$ is set to $\mathbf{M}_\mathbf{I}^{(i)}(\mathbf{p})$. When $s_{\mathbf{M}_\mathbf{I}^{(i)}} \geq s$, the i-th

entry of $\mathbf{V}(s, \mathbf{p})$ equals to 1 only if $\mathbf{M}_\mathbf{I}^{(i)}(\mathbf{p}) = 1$ and all smaller masks in $\{\mathbf{M}_\mathbf{I}^{(j)} \mid s \leq s_{\mathbf{M}_\mathbf{I}^{(j)}} < s_{\mathbf{M}_\mathbf{I}^{(i)}}\}$ equals to 0 at pixel $\mathbf{p}$. Formally, we have:

$$\mathbf{V}^i(s, \mathbf{p}) = \begin{cases} \mathbf{M}_\mathbf{I}^{(i)}(\mathbf{p}) & \text{if } s_{\mathbf{M}_\mathbf{I}^{(i)}} < s \text{ or } \mathcal{C}(\mathbf{p}) \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{C}(\mathbf{p}) \triangleq \left( \forall \mathbf{M} \in \{\mathbf{M}_\mathbf{I}^{(j)} \mid s \leq s_{\mathbf{M}_\mathbf{I}^{(j)}} < s_{\mathbf{M}_\mathbf{I}^{(i)}}\}, \mathbf{M}(\mathbf{p}) = 0 \right)$$
(6)

This assignment of pixel identity vectors is based on the fact that if a pixel belongs to a specific mask at a given scale, it will continue to belong to that mask at larger scales.

**Loss Function**    We adapt the correspondence distillation loss (Hamilton et al. 2022) for training the scale-gated Gaussian affinity features. Concretely, for two pixels $\mathbf{p}_1, \mathbf{p}_2$ at a given scale $s$, their mask correspondence is given by:

$$\texttt{Corr}_m(s, \mathbf{p}_1, \mathbf{p}_2) = \mathbb{1}(\mathbf{V}(s, \mathbf{p}_1) \cdot \mathbf{V}(s, \mathbf{p}_2)), \quad (7)$$

where $\mathbb{1}(\cdot)$ is the indicator function, which is equal to 1 when the input greater than or equal to 0. The feature correspondence between two pixels is defined as the cosine similarity between their scale-gated features:

$$\texttt{Corr}_f(s, \mathbf{p}_1, \mathbf{p}_2) = \langle \mathbf{F}^s(\mathbf{p}_1), \mathbf{F}^s(\mathbf{p}_2) \rangle. \quad (8)$$

The correspondence distillation loss $\mathcal{L}_{\text{corr}}(\mathbf{p}_1, \mathbf{p}_2)$ between two pixels is given by[1]:

$$\begin{aligned} \mathcal{L}_{\text{corr}}(s, \mathbf{p}_1, \mathbf{p}_2) = &(1 - 2 \cdot \texttt{Corr}_m(s, \mathbf{p}_1, \mathbf{p}_2)) \\ &\cdot \max(\texttt{Corr}_f(s, \mathbf{p}_1, \mathbf{p}_2), 0) \end{aligned} \quad (9)$$

---

[1]The feature correspondence is clipped at 0 to stabilize training. Please refer to Hamilton et al. (2022) for more details.

**Feature Norm Regularization** During training, the 2D features are obtained by rendering with 3D affinity features. This indicates a misalignment between 2D and 3D features. As revealed in Equation (4), a 2D feature is a linear combination of multiple 3D features, each with distinct directions. In such situations, SAGA may show good segmentation ability on the rendered feature map but perform poorly in 3D space. This motivates us to introduce a feature norm regularization. Concretely, during rendering the 2D feature map, the 3D features are first normalized as unit vectors, *i.e.*,

$$\mathbf{F}(\mathbf{p}) = \sum_{i=1}^{|\mathcal{G}_{\mathbf{p}}|} \frac{\mathbf{f}_{\mathbf{g}_i^{\mathbf{p}}}}{||\mathbf{f}_{\mathbf{g}_i^{\mathbf{p}}}||_2} \alpha_{\mathbf{g}_i^{\mathbf{p}}} \prod_{j=1}^{i-1} (1 - \alpha_{\mathbf{g}_j^{\mathbf{p}}}). \quad (10)$$

Accordingly, $||\mathbf{F}(\mathbf{p})||_2$ ranges in $[0, 1]$. When 3D features along a ray are perfectly aligned, $||\mathbf{F}(\mathbf{p})||_2 = 1$. Thus, we impose a regularization on the rendered feature norm:

$$\mathcal{L}_{\text{norm}}(\mathbf{p}) = 1 - ||\mathbf{F}(\mathbf{p})||_2 \quad (11)$$

With the feature norm regularization term, for an iteration of training, the loss of SAGA is defined as:

$$\mathcal{L} = \sum_{(\mathbf{p}_1, \mathbf{p}_2) \in \delta(\mathbf{I}) \times \delta(\mathbf{I})} \mathcal{L}_{\text{corr}}(\mathbf{p}_1, \mathbf{p}_2) + \sum_{\mathbf{p} \in \delta(\mathbf{I})} \mathcal{L}_{\text{norm}}(\mathbf{p}), \quad (12)$$

where $\delta(\mathbf{I})$ denotes the set of pixels within the image $\mathbf{I}$.

**Additional Training Strategy** During the training, an unavoidable issue is the data imbalance, reflected by: 1) Most pixel pairs keep to be positive or negative regardless of scale variations, making the learned feature insensitive to scales; 2) The majority of pixel pairs shows negative correspondence, resulting in feature collapse; 3) Large targets that occupy more pixels in images have more effect on the optimization, leading to bad performance of segmenting small targets. We tackle this problem by resampling the pixel-pairs and re-weighting the loss function for different samples. Please refer to Section A.1 of the appendix.

## 3.5 Inference

With well-trained Gaussian affinity features, SAGA can conduct various segmentation tasks in the 3D space. For promptable segmentation, SAGA takes **2D point prompts** at specific view and the scale as input. Then, SAGA segments the 3D target by matching scale-gated 3D Gaussian affinity features with the 2D query features selected from the rendered feature map according to prompt points. For automatic scene decomposing, SAGA employs HDBSCAN to cluster the affinity features directly in the 3D space[2]. Additionally, we design a vote-based segmentation mechanism to integrate SAGA with CLIP for conducting open vocabulary segmentation (see Section A.3 of the supplement).

## 4 Experiments

In this section, we demonstrate the effectiveness of SAGA both quantitatively and qualitatively. For implementation details and experiment settings, please refer to Section A.2.

[2]For efficiency, SAGA uniformly selects 1% of the Gaussians from the 3D-GS model for clustering.

| Method | mIoU (%) | mAcc (%) |
|---|---|---|
| NVOS (Ren et al. 2022) | 70.1 | 92.0 |
| ISRF (Goel et al. 2023) | 83.8 | 96.4 |
| SA3D (Cen et al. 2023) | 90.3 | 98.2 |
| OmniSeg3D (Ying et al. 2024) | 91.7 | 98.4 |
| GauGroup (Ye et al. 2024) | 85.6 | 97.3 |
| SA3D-GS (Cen et al. 2024) | 92.2 | 98.5 |
| SAGA (ours) | **92.6** | **98.6** |

Table 1: Results on NVOS dataset.

## 4.1 Datasets

For promptable segmentation experiments, we utilize two datasets: NVOS (Ren et al. 2022) and SPIn-NeRF(Mirzaei et al. 2023). The former is derived from the LLFF dataset (Mildenhall et al. 2019) and the latter is a combination of subsets of data from established NeRF-related datasets (Mildenhall et al. 2019, 2020; Lin et al. 2022; Knapitsch et al. 2017; Fridovich-Keil et al. 2022). For open-vocabulary segmentation experiments, we adopt the 3D-OVS dataset (Liu et al. 2023a). For qualitative analysis (Figure 3 and 4), we employ various datasets including LLFF (Mildenhall et al. 2019), MIP-360 (Barron et al. 2022), Tanks&Temple (Knapitsch et al. 2017), and Replica (Straub et al. 2019). These datasets encompass indoor and outdoor scenes, forward-facing and 360-degree scenes, as well as synthetic and real scenes.

## 4.2 Quantitative Results

**NVOS Dataset** As shown in Table 1, SAGA outperforms previous segmentation approaches for both 3D-GS and other radiance fields, *i.e.*, +0.4 mIoU over the previous SOTA SA3D-GS and +0.9 mIoU over the OmniSeg3D.

**SPIn-NeRF Dataset** Results on the SPIn-NeRF dataset can be found in Table 2. SA3D performs on par with the previous SOTA OmniSeg3D. The minor performance degradation is attributed to sub-optimal geometry learned by 3D-GS. For example, 3D-GS models the reflection effects with numerous outlier Gaussians that are not aligned with the exact geometry of the object. Excluding these Gaussians results in empty holes in the segmentation mask for certain views, while including them introduces noise in other views. Nevertheless, we believe the segmentation accuracy of SAGA can meet most requirements.

**Open-Vocabulary Semantic Segmentation** As shown in Table 3, SAGA demonstrates superior results across all scenes in the 3D-OVS dataset. For more details about open-vocabulary segmentation, please refer to Section A.3.

**Time Consumption Analysis** In Table 4, we reveal the time consumption of SAGA and compare it with existing promptable segmentation methods in radiance fields. ISRF trains a feature field by mimicking the 2D multi-view visual features extracted by DINO (Caron et al. 2021), thus enjoys faster convergence speed. However, this leads to less accurate segmentation, necessitating extensive post-processing.
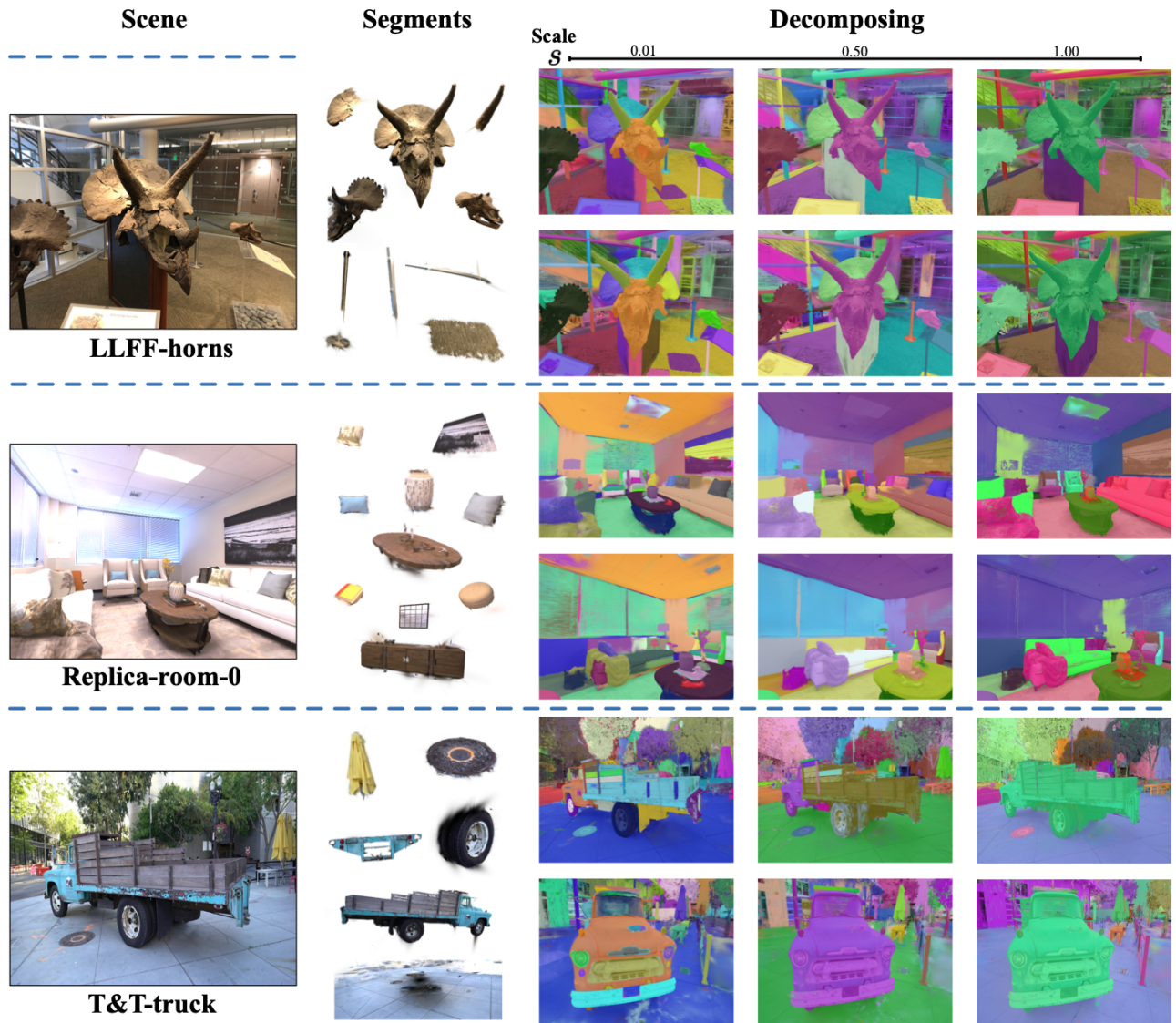
Figure 3: Qualitative results of SAGA across different scenes. We provide both the targets segmented via 2D point prompts and the "segment everything" results.

| Method | mIoU (%) | mAcc (%) |
|---|---|---|
| MVSeg (Mirzaei et al. 2023) | 90.9 | 98.9 |
| SA3D (Cen et al. 2023) | 92.4 | 98.9 |
| OmniSeg3D (Ying et al. 2024) | **94.3** | **99.3** |
| GauGroup (Ye et al. 2024) | 86.5 | 98.9 |
| SA-GS (Hu et al. 2024) | 89.9 | 98.7 |
| SA3D-GS (Cen et al. 2024) | 93.2 | 99.1 |
| SAGA (ours) | 93.4 | 99.2 |

Table 2: Results on SPIn-NeRF dataset.

| Method | bed | bench | room | sofa | lawn | mean |
|---|---|---|---|---|---|---|
| LERF | 73.5 | 53.2 | 46.6 | 27.0 | 73.7 | 54.8 |
| 3D-OVS | 89.5 | 89.3 | 92.8 | 74.0 | 88.2 | 86.8 |
| LangSplat | 92.5 | 94.2 | 94.1 | 90.0 | 96.1 | 93.4 |
| N2F2 | 93.8 | 92.6 | 93.5 | 92.1 | 96.3 | 93.9 |
| SAGA | **97.4** | **95.4** | **96.8** | **93.5** | **96.6** | **96.0** |

Table 3: Results on 3D-OVS dataset (mIoU).

| Method | Training | Inference |
|---|---|---|
| SA3D (Cen et al. 2023) | - | 45 s |
| SA-GS (Hu et al. 2024) | - | 15 s |
| ISRF (Goel et al. 2023) | 2.5 mins | 3.3 s |
| OmniSeg3D (Ying et al. 2024) | 15~40 mins | 50~100 ms |
| GARField (Kim et al. 2024) | 20~60 mins | 30~70 ms |
| SAGA (ours) | 10~40 mins | **2~5 ms** |

Table 4: Time consumption comparison.



Figure 4: SAGA can maintain the high frequency texture details captured by 3D-GS. We reveal the inherent structure of these details by shrinking the Gaussians by 60%.

SA3D and SA-GS employ an iterative mask refinement pipeline, which eliminates the need for training but incurs significant inference time consumption. Compared to methods that distill segmentation capabilities from SAM masks, such as OmniSeg3D and GARField, SAGA demonstrates much faster inference speed and comparable training speed.

### 4.3 Qualitative Results

Figure 3 shows that SAGA achieves fine-grained segmentation at various scales across different scenes. The compact Gaussian affinity features enable scene decomposition through simple clustering in the 3D-GS model. The needle-like artifacts on the edges of segmented targets are due to 3D-GS overfitting multi-view RGB without object awareness. Building upon 3D-GS, which can capture high-frequency texture details, SAGA can effectively segment thin, fine-grained structures, as shown in Figure 4. By shrinking the 3D Gaussians, we reveal the underlying structural modeling capabilities of 3D-GS and demonstrate the completeness of the segmentation results. Since GARField lacks quantitative results on the NVOS and SPIn-NeRF datasets, we provide qualitative comparisons in Section A.6 to highlight the effectiveness of SAGA's learned features.

### 4.4 Ablation Study

**Local Feature Smoothing (LFS) & Feature Norm Regularization (FNR)** Both LFS and FNR impose constraints on the Gaussian affinity features. We present visualization results to illustrate their roles.

In Figure 5, when segmenting a 3D object with a cosine similarity threshold of 0.75, without the feature smoothing, the result shows many false positives, which reveals that outliers are primarily eliminated by the feature smoothing op-



Figure 5: Ablation study on effects of local feature smoothing (Smooth) and feature norm regularization (Feature Norm). Outliers are primarily eliminated through local feature smoothing. Feature norm regularization helps features of inner Gaussians align better with those of the surface.



Figure 6: Failure cases of SAGA. The targets of interest are labeled by red border.

eration. Unlike local feature smoothing, feature norm regularization primarily impacts the Gaussians within objects. When raising the similarity score threshold to 0.95, the apple segmented by SAGA with feature norm regularization remains intact, while the one without it quickly becomes translucent. This phenomenon supports our assumption that 3D features are not perfectly aligned with 2D features, as introduced in Section 3.4. Imposing feature norm regularization helps align the affinity features of 3D Gaussians by pulling the features along a ray in the same direction.

## 5 Limitation

SAGA learns the affinity features from multi-view 2D masks extracted by SAM. This makes SAGA hardly segment objects that are not appeared in these masks. As shown in Figure 6, this limitation is particularly evident when the target of interest is small. Enhancing the generalization ability of SAGA to unrecognized targets during the automatic extraction stage is a promising direction.

## 6 Conclusion

In this paper, we propose SAGA, a 3D promptable segmentation method for 3D Gaussian Splatting (3D-GS). SAGA injects the segmentation capability of SAM into Gaussian affinity features for all 3D Gaussians in a 3D-GS model, endowing them with a new property towards segmentation. To preserve the multi-granularity segmentation ability of SAM and the efficiency of 3D-GS, SAGA introduces a lightweight scale-gate mechanism, which adapts the affinity features according to different 3D physical scales with minimal computation overhead. After training, SAGA can achieve real-time fine-grained 3D segmentation. Comprehensive experiments are conducted to demonstrate the effectiveness of SAGA. As one of the first methods addressing promptable segmentation in 3D-GS, the simplicity and effectiveness of SAGA pave the way for future advancements in this field.

## Acknowledgments

## References

Barron, J. T.; Mildenhall, B.; Verbin, D.; Srinivasan, P. P.; and Hedman, P. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In *CVPR*.

Bhalgat, Y.; Laina, I.; Henriques, J. a. F.; Vedaldi, A.; and Zisserman, A. 2023. Contrastive Lift: 3D Object Instance Segmentation by Slow-Fast Contrastive Fusion. In *NeurIPS*.

Bhalgat, Y.; Laina, I.; Henriques, J. F.; Zisserman, A.; and Vedaldi, A. 2024. N2F2: Hierarchical Scene Understanding with Nested Neural Feature Fields. arXiv:2403.10997.

Bing, W.; Chen, L.; and Yang, B. 2023. DM-NeRF: 3D Scene Geometry Decomposition and Manipulation from 2D Images. In *ICLR*.

Boykov, Y. Y.; and Jolly, M.-P. 2001. Interactive graph cuts for optimal boundary & region segmentation of objects in ND images. In *ICCV*.

Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; and Joulin, A. 2021. Emerging Properties in Self-Supervised Vision Transformers. In *ICCV*.

Cen, J.; Fang, J.; Zhou, Z.; Yang, C.; Xie, L.; Zhang, X.; Shen, W.; and Tian, Q. 2024. Segment Anything in 3D with Radiance Fields. arXiv:2304.12308.

Cen, J.; Zhou, Z.; Fang, J.; Yang, C.; Shen, W.; Xie, L.; Jiang, D.; Zhang, X.; and Tian, Q. 2023. Segment Anything in 3D with NeRFs. In *NeurIPS*.

Chen, A.; Xu, Z.; Geiger, A.; Yu, J.; and Su, H. 2022a. TensoRF: Tensorial Radiance Fields. In *ECCV*.

Chen, X.; Tang, J.; Wan, D.; Wang, J.; and Zeng, G. 2023. Interactive Segment Anything NeRF with Feature Imitation. arXiv:2305.16233.

Chen, X.; Zhao, Z.; Zhang, Y.; Duan, M.; Qi, D.; and Zhao, H. 2022b. Focalclick: Towards practical interactive image segmentation. In *CVPR*.

Fan, Z.; Wang, P.; Jiang, Y.; Gong, X.; Xu, D.; and Wang, Z. 2023. NeRF-SOS: Any-View Self-supervised Object Segmentation on Complex Scenes. In *ICLR*.

Fang, J.; Yi, T.; Wang, X.; Xie, L.; Zhang, X.; Liu, W.; Nießner, M.; and Tian, Q. 2022. Fast Dynamic Radiance Fields with Time-Aware Neural Voxels. In *SIGGRAPH Asia 2022 Conference Papers*.

Fridovich-Keil, S.; Yu, A.; Tancik, M.; Chen, Q.; Recht, B.; and Kanazawa, A. 2022. Plenoxels: Radiance Fields without Neural Networks. In *CVPR*.

Fu, X.; Zhang, S.; Chen, T.; Lu, Y.; Zhu, L.; Zhou, X.; Geiger, A.; and Liao, Y. 2022. Panoptic NeRF: 3D-to-2D Label Transfer for Panoptic Urban Scene Segmentation. In *3DV*.

Goel, R.; Sirikonda, D.; Saini, S.; and Narayanan, P. 2023. Interactive segmentation of radiance fields. In *CVPR*.

Grady, L. 2006. Random walks for image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*

Gulshan, V.; Rother, C.; Criminisi, A.; Blake, A.; and Zisserman, A. 2010. Geodesic star convexity for interactive image segmentation. In *CVPR*.

Guo, H.; Zhu, H.; Peng, S.; Wang, Y.; Shen, Y.; Hu, R.; and Zhou, X. 2024. SAM-guided Graph Cut for 3D Instance Segmentation. arXiv:2312.08372.

Hamilton, M.; Zhang, Z.; Hariharan, B.; Snavely, N.; and Freeman, W. T. 2022. Unsupervised Semantic Segmentation by Distilling Feature Correspondences. In *ICLR*.

Hedman, P.; Srinivasan, P. P.; Mildenhall, B.; Reiser, C.; Barron, J. T.; and Debevec, P. 2024. Baking Neural Radiance Fields for Real-Time View Synthesis. *IEEE TPAMI*.

Hu, X.; Wang, Y.; Fan, L.; Fan, J.; Peng, J.; Lei, Z.; Li, Q.; and Zhang, Z. 2024. SAGD: Boundary-Enhanced Segment Anything in 3D Gaussian via Gaussian Decomposition. arXiv:2401.17857.

Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM TOG*.

Kerr, J.; Kim, C. M.; Goldberg, K.; Kanazawa, A.; and Tancik, M. 2023. Lerf: Language embedded radiance fields. In *ICCV*.

Kim, C. M.; Wu, M.; Kerr, J.; Tancik, M.; Goldberg, K.; and Kanazawa, A. 2024. GARField: Group Anything with Radiance Fields. In *CVPR*.

Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; et al. 2023. Segment anything. In *ICCV*.

Knapitsch, A.; Park, J.; Zhou, Q.-Y.; and Koltun, V. 2017. Tanks and Temples: Benchmarking Large-Scale Scene Reconstruction. *ACM TOG*.

Kobayashi, S.; Matsumoto, E.; and Sitzmann, V. 2022. Decomposing NeRF for Editing via Feature Field Distillation. In *NeurIPS*.

Liao, G.; Li, J.; Bao, Z.; Ye, X.; Wang, J.; Li, Q.; and Liu, K. 2024. CLIP-GS: CLIP-Informed Gaussian Splatting for Real-time and View-consistent 3D Semantic Understanding. arXiv:2404.14249.

Lin, Y.; Florence, P.; Barron, J. T.; Lin, T.; Rodriguez, A.; and Isola, P. 2022. NeRF-Supervision: Learning Dense Object Descriptors from Neural Radiance Fields. In *ICRA*.

Lindell, D. B.; Martel, J. N. P.; and Wetzstein, G. 2021. AutoInt: Automatic Integration for Fast Neural Volume Rendering. In *CVPR*.

Liu, K.; Zhan, F.; Zhang, J.; XU, M.; Yu, Y.; Saddik, A. E.; Theobalt, C.; Xing, E.; and Lu, S. 2023a. Weakly Supervised 3D Open-vocabulary Segmentation. In *NeurIPS*.

Liu, Q.; Xu, Z.; Bertasius, G.; and Niethammer, M. 2023b. Simpleclick: Interactive image segmentation with simple vision transformers. In *ICCV*.

Liu, S.; Zeng, Z.; Ren, T.; Li, F.; Zhang, H.; Yang, J.; Jiang, Q.; Li, C.; Yang, J.; Su, H.; Zhu, J.; and Zhang, L. 2024a. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. arXiv:2303.05499.

Liu, X.; Chen, J.; Yu, H.; Tai, Y.; and Tang, C. 2022. Unsupervised Multi-View Object Segmentation Using Radiance Field Propagation. In *NeurIPS*.

Liu, Y.; Hu, B.; Huang, J.; Tai, Y.-W.; and Tang, C.-K. 2023c. Instance neural radiance field. In *ICCV*.

Liu, Y.; Hu, B.; Tang, C.-K.; and Tai, Y.-W. 2024b. SANeRF-HQ: Segment Anything for NeRF in High Quality. In *CVPR*.

Lyu, W.; Li, X.; Kundu, A.; Tsai, Y.-H.; and Yang, M.-H. 2024. Gaga: Group Any Gaussians via 3D-aware Memory Bank. arXiv:2404.07977.

Mildenhall, B.; Srinivasan, P. P.; Cayon, R. O.; Kalantari, N. K.; Ramamoorthi, R.; Ng, R.; and Kar, A. 2019. Local light field fusion: practical view synthesis with prescriptive sampling guidelines. *ACM TOG*.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.

Mirzaei, A.; Aumentado-Armstrong, T.; Derpanis, K. G.; Kelly, J.; Brubaker, M. A.; Gilitschenski, I.; and Levinshtein, A. 2023. SPIn-NeRF: Multiview Segmentation and Perceptual Inpainting with Neural Radiance Fields. In *CVPR*.

Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*.

Niemeyer, M.; and Geiger, A. 2021. GIRAFFE: Representing Scenes As Compositional Generative Neural Feature Fields. In *CVPR*.

Qin, M.; Li, W.; Zhou, J.; Wang, H.; and Pfister, H. 2024. LangSplat: 3D Language Gaussian Splatting. In *CVPR*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; Krueger, G.; and Sutskever, I. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*.

Ren, Z.; Agarwala, A.; Russell, B. C.; Schwing, A. G.; and Wang, O. 2022. Neural Volumetric Object Selection. In *CVPR*.

Rother, C.; Kolmogorov, V.; and Blake, A. 2004. "GrabCut": interactive foreground extraction using iterated graph cuts. *ACM TOG*.

Siddiqui, Y.; Porzi, L.; Bulò, S. R.; Müller, N.; Nießner, M.; Dai, A.; and Kontschieder, P. 2023. Panoptic lifting for 3d scene understanding with neural fields. In *CVPR*.

Sofiiuk, K.; Petrov, I. A.; and Konushin, A. 2022. Reviving iterative training with mask guidance for interactive segmentation. In *ICIP*.

Stelzner, K.; Kersting, K.; and Kosiorek, A. R. 2021. Decomposing 3D Scenes into Objects via Unsupervised Volume Segmentation. arXiv:2104.01148.

Straub, J.; Whelan, T.; Ma, L.; Chen, Y.; Wijmans, E.; Green, S.; Engel, J. J.; Mur-Artal, R.; Ren, C.; Verma, S.; Clarkson, A.; Yan, M.; Budge, B.; Yan, Y.; Pan, X.; Yon, J.; Zou, Y.; Leon, K.; Carter, N.; Briales, J.; Gillingham, T.;

Mueggler, E.; Pesqueira, L.; Savva, M.; Batra, D.; Strasdat, H. M.; Nardi, R. D.; Goesele, M.; Lovegrove, S.; and Newcombe, R. 2019. The Replica Dataset: A Digital Replica of Indoor Spaces. arXiv:1906.05797.

Sun, C.; Sun, M.; and Chen, H. 2022. Direct Voxel Grid Optimization: Super-fast Convergence for Radiance Fields Reconstruction. In *CVPR*.

Tschernezki, V.; Laina, I.; Larlus, D.; and Vedaldi, A. 2022. Neural Feature Fusion Fields: 3D Distillation of Self-Supervised 2D Image Representations. In *3DV*.

Vora, S.; Radwan, N.; Greff, K.; Meyer, H.; Genova, K.; Sajjadi, M. S.; Pot, E.; Tagliasacchi, A.; and Duckworth, D. 2022. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *TMLR*.

Wizadwongsa, S.; Phongthawee, P.; Yenphraphai, J.; and Suwajanakorn, S. 2021. NeX: Real-Time View Synthesis With Neural Basis Expansion. In *CVPR*.

Xu, M.; Yin, X.; Qiu, L.; Liu, Y.; Tong, X.; and Han, X. 2023. SAMPro3D: Locating SAM Prompts in 3D for Zero-Shot Scene Segmentation. arXiv:2311.17707.

Yang, Y.; Wu, X.; He, T.; Zhao, H.; and Liu, X. 2023. SAM3D: Segment Anything in 3D Scenes. arXiv:2306.03908.

Ye, M.; Danelljan, M.; Yu, F.; and Ke, L. 2024. Gaussian Grouping: Segment and Edit Anything in 3D Scenes. arXiv:2312.00732.

Yin, Y.; Liu, Y.; Xiao, Y.; Cohen-Or, D.; Huang, J.; and Chen, B. 2024. Sai3d: Segment any instance in 3d scenes. In *CVPR*.

Ying, H.; Yin, Y.; Zhang, J.; Wang, F.; Yu, T.; Huang, R.; and Fang, L. 2024. OmniSeg3D: Omniversal 3D Segmentation via Hierarchical Contrastive Learning. In *CVPR*.

Yu, H.; Guibas, L. J.; and Wu, J. 2022. Unsupervised Discovery of Object Radiance Fields. In *ICLR*.

Zhi, S.; Laidlow, T.; Leutenegger, S.; and Davison, A. J. 2021. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *ICCV*.

Zou, X.; Yang, J.; Zhang, H.; Li, F.; Li, L.; Wang, J.; Wang, L.; Gao, J.; and Lee, Y. J. 2023. Segment everything everywhere all at once. In *NeurIPS*.

# A  Appendix

In this appendix we provide the concrete training strategy (Section A.1) of SAGA and implementation details (Section A.2). Then, we provide details about the open-vocabulary segmentation ability of SAGA and analyze its limitation (Section A.3). We also provide an interpretability analysis about the scale gate mechanism (Section A.4) to reveal the underlying principle of SAGA. Section A.5 evaluates the robustness and generalizability of SAGA by applying it to more kinds of radiance fields, and Section A.6 presents additional visualization results to demonstrate its effectiveness.

## A.1  Detailed Additional Training Stategy

During the contrastive-based learning, an unavoidable problem is the data imbalance. In SAGA, the data imbalance is reflected in the following three aspects: 1) Scale-sensitivity imbalance. The majority of pixel pairs exhibit insensitivity to changes in scale. In other words, during a training iteration, most pixel pairs maintain their positive or negative classification regardless of scale variations. This makes the scale gate collapse to constant output; 2) Positive-negative samples imbalance. The majority of pixel pairs shows negative correspondence, resulting in segmentation features degradation; 3) Target-size imbalance. Large targets that occupy more pixels in images have more effect on the optimization, which leads to bad performance of segmenting small targets.

We use a resampling strategy to tackle the scale-sensitivity imbalance and positive-negative samples imbalance. Then, we adopt a pixel-wise re-weighting strategy to tackle the target-size imbalance.

**Resampling**  In each iteration of training, we randomly select $N_s$ scales and $N_p$ pixels within an image and form $N_p \times N_p$ pixel pairs. Calculating the mask correspondence for these pixel pairs at every sampled scale results in a scale-conditioned correspondence matrix $\mathbf{R} \in \{0,1\}^{N_s N_p N_p}$. We split the pixel pairs into three sets according to $\mathbf{R}$: 1) Inconsistent set: $\mathcal{Q}^{\text{in}} = \{(\mathbf{p}_1, \mathbf{p}_2) \mid \exists s_1, s_2, \ \mathbf{R}(s_1, \mathbf{p}_1, \mathbf{p}_2) \neq \mathbf{R}(s_2, \mathbf{p}_1, \mathbf{p}_2)\}$; 2) Consistent positive set: $\mathcal{Q}^{\text{pos}} = \{(\mathbf{p}_1, \mathbf{p}_2) \mid \forall s, \ \mathbf{R}(s, \mathbf{p}_1, \mathbf{p}_2) = 1\}$; 3) Consistent negative set: $\mathcal{Q}^{\text{neg}} = \{(\mathbf{p}_1, \mathbf{p}_2) \mid \forall s, \ \mathbf{R}(s, \mathbf{p}_1, \mathbf{p}_2) = 0\}$. During the loss calculation, all pixel pairs in $\mathcal{Q}^{\text{in}}$ are involved. Then, we randomly select $\frac{|\mathcal{Q}^{\text{in}}|}{2}$ pixel pairs in both $\mathcal{Q}^{\text{pos}}$ and $\mathcal{Q}^{\text{neg}}$ respectively. To tackle the hard samples in training, we also add pixel pairs in $\mathcal{Q}^{\text{neg}}$ which have feature correspondences larger than 0.5 and pixel pairs in $\mathcal{Q}^{\text{pos}}$ which have feature correspondence smaller than 0.75 into loss calculation. This design not only ensures the sensitivity of the loss to scale changing but also keeps the balance of positive pairs and negative pairs. Let $\phi(\cdot)$ denote the sampling operation, after the resampling, we get three sets of pixel pairs, *i.e.*, $\mathcal{Q}^{\text{in}}, \phi(\mathcal{Q}^{\text{pos}}), \phi(\mathcal{Q}^{\text{neg}})$.

**Re-weighting**  Considering two masks $\mathbf{M}_\mathbf{I}^1, \mathbf{M}_\mathbf{I}^2$ in $\mathcal{M}_\mathbf{I}$, when uniformly sampling a pair of pixels, the probability that the pair is from $\mathbf{M}_\mathbf{I}^1, \mathbf{M}_\mathbf{I}^2$ is proportional to the product of the number of positive pixels of the two mask. This indicates that the optimization process is dominated by large

masks. To re-weight the loss, we first calculate the mean mask size $m_\mathbf{p} = \frac{1}{|\mathcal{K}_\mathbf{p}|} \sum_{\mathbf{M} \in \mathcal{K}_\mathbf{p}} \sum_{i=1,j=1}^{H,W} \mathbf{M}(i,j)$ for each pixel, where $\mathcal{K}_\mathbf{p} = \{\mathbf{M} \mid \mathbf{M} \in \mathcal{M}_\mathbf{I}, \mathbf{M}(\mathbf{p}) = 1\}$. For a pixel pair $(\mathbf{p}_1, \mathbf{p}_2)$, the loss weight is defined as $\omega(\mathbf{p}_1, \mathbf{p}_2) = \frac{1}{m_{\mathbf{p}_1} m_{\mathbf{p}_2}}$. All weights in a training iteration are then min-max normalized to the range $[1, 10]$ to ensure stable training.

With the resampling, re-weighting strategies, the overall loss function of SAGA is:

$$
\begin{aligned}
\mathcal{L} = & \frac{\sum_{(\mathbf{p}_1, \mathbf{p}_2) \in \mathcal{Q}^{\text{in}} \cup \phi(\mathcal{Q}^{\text{pos}})} \omega(\mathbf{p}_1, \mathbf{p}_2) \mathcal{L}_{\text{corr}}(\mathbf{p}_1, \mathbf{p}_2)}{|\mathcal{Q}^{\text{in}} \cup \phi(\mathcal{Q}^{\text{pos}})|} + \\
& \frac{\sum_{(\mathbf{p}_1, \mathbf{p}_2) \in \mathcal{Q}^{\text{in}} \cup \phi(\mathcal{Q}^{\text{neg}})} \omega(\mathbf{p}_1, \mathbf{p}_2) \mathcal{L}_{\text{corr}}(\mathbf{p}_1, \mathbf{p}_2)}{|\mathcal{Q}^{\text{in}} \cup \phi(\mathcal{Q}^{\text{neg}})|} + \\
& \frac{1}{HW} \sum_{\mathbf{p} \in \delta(\mathbf{I})} \mathcal{L}_{\text{norm}}(\mathbf{p}),
\end{aligned}
\tag{13}
$$

where $\delta(\mathbf{I})$ denotes the set of pixels within the image $\mathbf{I}$.

## A.2  Implementation Details

Across different scenes, SAGA maintains consistent hyper-parameters. The feature dimension $D$ is set to 32. The $K$ of KNN used in local feature smoothing is set to 16. Training of Gaussian affinity features lasts for 10,000 iterations. In each iteration, we randomly sample eight different scales and 1,000 pixels ($1000^2$ pixel pairs) from the view for training. For different loss terms, we do not adjust any loss balance coefficients in experiments.

We extract the multi-view 2D masks with the SAM ViT-H model. For open-vocabulary segmentation, we use the Open-CLIP ViT-B/16 model. All training and inference is conducted on a single Nvidia RTX 3090 GPU.

**Experiment Setting**  For the NVOS dataset, we randomly select positive and negative points from the scribbles on the reference view (provided by the NVOS dataset) to conduct promptable 3D segmentation. We then render the 3D segmentation result on the target view and evaluate the Intersection over Union (IoU) and pixel-wise accuracy against the provided ground truth. For each scene in the SPIn-NeRF dataset, we randomly select a subset of points within and outside the mask of the reference view as positive and negative prompts.

## A.3  Vote-based Open-vocabulary Segmentation

To enable open-vocabulary 3D segmentation in radiance fields, previous studies (Kerr et al. 2023; Liu et al. 2023a; Qin et al. 2024; Bhalgat et al. 2024; Liao et al. 2024) focus on aligning 3D language feature fields with the visual features extracted by CLIP (Radford et al. 2021) image encoder. Then 3D segmentation can be achieved by querying the language fields with the textual features. The process of training language feature field can be regarded as multi-view feature fusing, which, in essence, is a kind of vote mechanism. This motivate us to see whether SAGA can handle open-vocabulary segmentation with well-trained Gaussian affinity features at minimum modification.

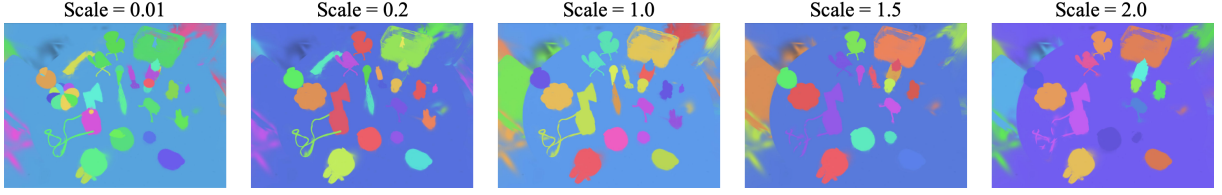| Scale = 0.01 | Scale = 0.2 | Scale = 1.0 | Scale = 1.5 | Scale = 2.0 |

Figure A1: Adapting the scale gate mechanism with GARField achieves competitive results, demonstrating the potential of SAGA across different radiance fields.

Readers may wonder why we do not use grounding methods like Grounding-DINO (Liu et al. 2024a) to locate the object and convert the location into point prompts, which can then be fed to SAGA. To answer this question, we emphasize the setting of open-vocabulary segmentation. We follow a stricter approach than previous studies, such as SA3D (Cen et al. 2023) and GaussianGrouping (Ye et al. 2024), which rely on Grounded-SAM to convert language prompts to visual prompts in a specific view. Specifying this view introduces additional prior knowledge. In SAGA, users do not need to query a specific view; all that is required for open-vocabulary segmentation is a text prompt. With this idea in mind, we introduce a vote-based open-vocabulary segmentation strategy.

**Constructing Vote Graph by Clustering** To cluster the multi-view masks, an intuitive way is to cluster their corresponding segmentation features. However, this is infeasible since the segmentation features in SAGA are scale-conditioned. Segmentation features for different masks are probably in different feature subspace. This drives us to find a kind of **global features** that is consistent across different scales for multi-view masks to enable global clustering.

We propose to use the segmented Gaussians as the global feature. First, we uniformly sample a set of anchor Gaussians $\mathcal{A}$ from $\mathcal{G}$. Then, for a 2D mask $\mathbf{M} \in \mathcal{M}_{\mathbf{I}}$ with scale $s_{\mathbf{M}}$, we calculate its scale-conditioned feature as[3]:

$$\mathbf{f}_{\mathbf{M}} = \frac{1}{\delta(\mathbf{M})} \sum_{\mathbf{p} \in \delta(\mathbf{M})} \mathbf{F}^{s_{\mathbf{M}}}(\mathbf{p}). \qquad (14)$$

Then we compute the similarities between $\mathbf{f}_{\mathbf{M}}$ and all anchor Gaussians in $\mathcal{A}$ to get a segmentation result $\mathcal{A}_{\mathbf{M}} = \{\mathbf{g} \mid \langle \mathbf{f}_{\mathbf{g}}^{s_{\mathbf{M}}}, \mathbf{f}_{\mathbf{M}} \rangle > \tau, \mathbf{g} \in \mathcal{A}\}$. The distance between two masks $\mathbf{M}_1, \mathbf{M}_2$ is defined as their intersection over union of this 3D segmentation result, *i.e.*, $D(\mathbf{M}_1, \mathbf{M}_2) = \frac{|\mathcal{A}_{\mathbf{M}_1} \cap \mathcal{A}_{\mathbf{M}_2}|}{|\mathcal{A}_{\mathbf{M}_1} \cup \mathcal{A}_{\mathbf{M}_2}|}$. Then we perform HDBSCAN based on this distance map to cluster the 2D masks.

**Vote-based Segmentation** After clustering, we obtain a vote graph $\mathcal{V}$, where masks of the same 3D target (instance or part) are grouped together. In other words, each cluster centroid in $\mathcal{V}$ represents a potential segmentation target in the 3D space.

For each 2D mask $\mathbf{M}_{\mathbf{I}}$ of image $\mathbf{I}$, we extract its visual feature by feeding the masked image $\mathbf{I} \odot \mathbf{M}_{\mathbf{I}}$ to the CLIP visual encoder. During inference, given any text prompt, a relevancy score $r_{\mathbf{M}}$ is assigned to each mask $\mathbf{M}$ by comparing the textual feature with the visual feature of the mask[4]. For a cluster centroid $\mathbf{T}$ in $\mathcal{V}$, the relevancy scores of its corresponding 2D masks are aggregated to form its final relevancy score, *i.e.*,

$$r_{\mathbf{T}} = \frac{1}{|\mathcal{V}_{\mathbf{T}}|} \sum_{\mathbf{M} \in \mathcal{V}_{\mathbf{T}}} r_{\mathbf{M}}, \qquad (15)$$

where $\mathcal{V}_{\mathbf{T}}$ is the set containing all 2D masks corresponding to the cluster centroid $\mathbf{T}$. For semantic segmentation in the 3D-OVS dataset, which provides a category list for each scene, the label of each cluster is assigned as the category with the highest relevancy score.

**Limitation of SAGA in Open-vocabulary Segmentation** The vote-based open-vocabulary segmentation strategy encounters difficulties in certain scenarios. For instance, when considering a bowl of noodles with an egg in it, SAGA ideally should distinguish between categories such as the egg, the egg with noodles (contents of the bowl), and the bowl with noodles and egg. However, because "egg" is included in masks at larger scales, CLIP often misclassifies larger objects as "egg." This issue is inherently rooted in the multi-granularity ambiguity in semantics. Addressing this problem is a promising research direction.

Another limitation is common among current CLIP-SAM based methods, such as LangSplat (Qin et al. 2024) and N2F2 (Bhalgat et al. 2024). Both SAGA and these methods first use SAM to segment the images and then use CLIP to attach semantics to these segments by segmenting the corresponding regions of the image and feeding them to CLIP. However, the effectiveness of the CLIP visual encoder also depends on context. For example, SAM sometimes segments the texture on a wall. Without seeing the entire wall, CLIP struggles to recognize it. This lack of context hinders the ability to ground the segments accurately. This is an important yet currently unexplored issue.

### A.4 Interpretability Analysis

To better understand the scale gate mechanism, we analyze the weights of the learned scale gates through a statistical analysis across 47 scenes, each containing 32 scale gates, resulting in a total of $47 \times 32$ entries. Within these entries, 36.1% (543) of the gates are positive, while 63.9% (961) are

---

[3]For brevity, we continue to use $\delta(\cdot)$ to denote the set of positive pixels in 2D masks.

[4]We adopt the relevancy score introduced in LERF (Kerr et al. 2023), which has been proven robust.

negative, meaning that a typical scene has about 12 positive gates and 20 negative gates. In other words, when larger-scale features are used as input, more gates tend to close, and fewer gates remain open, which aligns with the intuitive understanding that finer-grained segmentation requires more features to effectively capture detailed information.

## A.5 SAGA with Hybrid Radiance Fields

Although SAGA is designed for segmentation in 3D-GS, its scale gate mechanism is not confined to a particular type of radiance field. We demonstrate this versatility by adapting SAGA to GARField, replacing the scale-conditioned affinity field with a scale-gated affinity field. As illustrated in Figure A1, the performance remains competitive.

Furthermore, we apply SAGA to InstantNGP (Müller et al. 2022), which can be viewed as an explicit version of GARField without the MLP. As shown in Figure A2, SAGA continues to produce competitive results. It is worth noting that direct clustering on the hash grid of InstantNGP is infeasible. To address this, we follow GARField's approach, using 3D-GS to query the hash grids and extract 3D features for clustering.

## A.6 More Qualitative Results

We present additional qualitative results in Figures A5 and A4. In Figure A3, it is evident that the affinity features in SAGA exhibit better stability compared to GARField. We perform clustering on the affinity features across the entire scene. At a smaller scale (0.01), both SAGA and GARField achieve fine-grained segmentation. However, when conducting coarse segmentation at a larger scale, GARField tends to merge smaller objects into larger ones (*e.g.*, the table). This behavior can be attributed to two factors: first, GARField assumes that the entire scene belongs to the same "object" when the scale exceeds that of any existing objects within the scene. Second, GARField employs an implicit feature field to fit the scale-conditioned affinity features. As noted by Mildenhall et al. (2020), radiance fields often produce over-smoothed predictions, which can result in the loss of small objects at larger scales. In contrast, SAGA assigns an explicit affinity feature to each 3D Gaussian in the 3D-GS model and uses a simple one-layer linear layer (*i.e.*, the scale gate) to model the scale-conditioned effect. This approach allows SAGA to better preserve smaller objects. This is also evidenced by Figure A1: by replacing the multi-layer perceptron of GARField with our proposed scale-gate mechanism, GARField can preserve the small objects on the table.
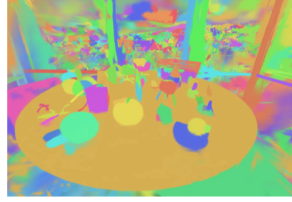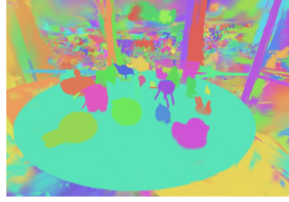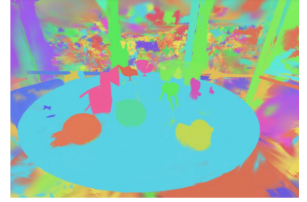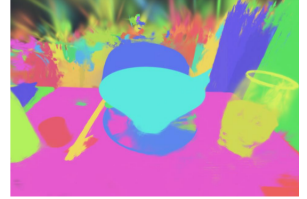
Figure A2: Applying SAGA to Instant-NGP achieves competitive segmentation performance, further demonstrating the generalizability and robustness of SAGA across different radiance field representations.
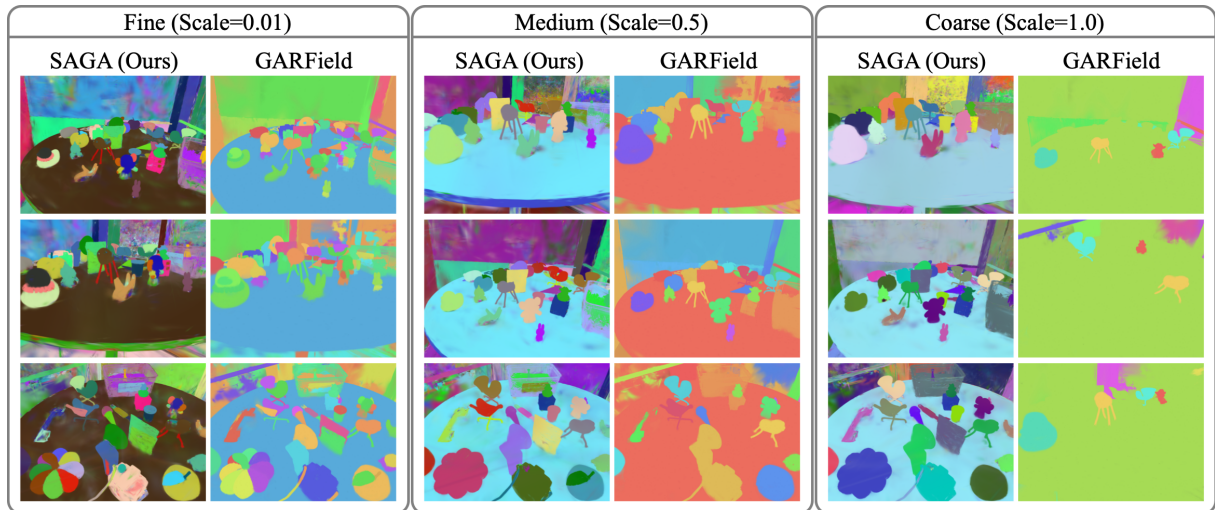


Figure A3: Qualitative comparison with GARField. We conduct feature clustering across the whole scene (LERF-figurines). Compared to GARField, which employs an additional feature field to model affinity features, SAGA demonstrates greater stability by utilizing explicit affinity features. At larger scales, SAGA effectively preserves the perception of small objects without merging them with other targets.
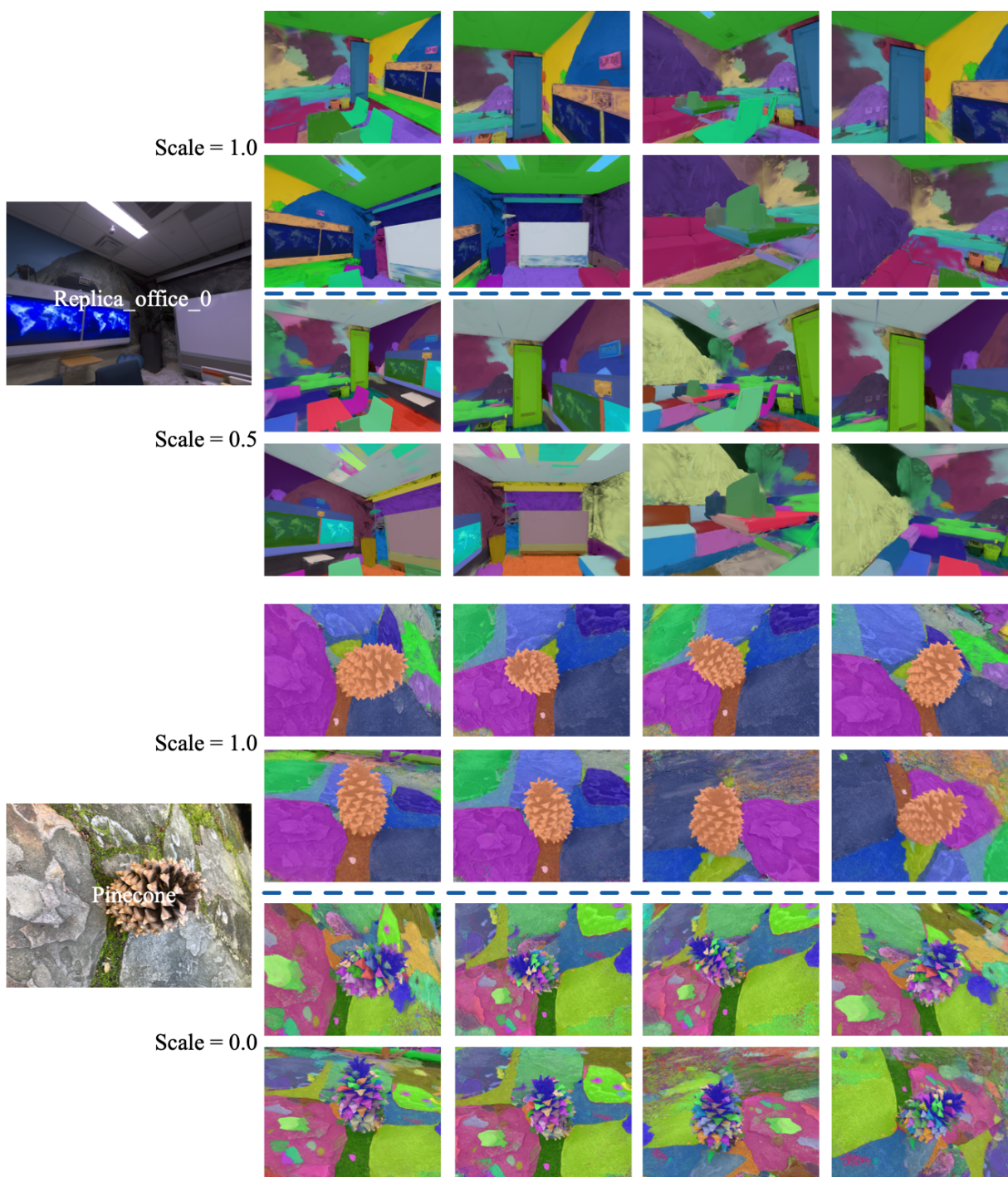
Figure A4: More qualitative results of SAGA.

MIP360-garden

LLFF-flower

LERF-Ramen

lego_real_night_radial

Replica_office_3

Figure A5: More qualitative results of SAGA.