# Click-Gaussian: Interactive Segmentation to Any 3D Gaussians

Seokhun Choi[1*], Hyeonseop Song[1*],
Jaechul Kim[1], Taehyeong Kim[2†], and Hoseok Do[1†]

[1] AI Lab, CTO Division, LG Electronics, Republic of Korea
{seokhun.choi, hyeonseop.song, jaechul1220.kim, hoseok.do}@lge.com
[2] Dept. of Biosystems Engineering, Seoul National University, Republic of Korea
taehyeong.kim@snu.ac.kr
Project page: https://seokhunchoi.github.io/Click-Gaussian

**Abstract.** Interactive segmentation of 3D Gaussians opens a great opportunity for real-time manipulation of 3D scenes thanks to the real-time rendering capability of 3D Gaussian Splatting. However, the current methods suffer from time-consuming post-processing to deal with noisy segmentation output. Also, they struggle to provide detailed segmentation, which is important for fine-grained manipulation of 3D scenes. In this study, we propose Click-Gaussian, which learns distinguishable feature fields of two-level granularity, facilitating segmentation without time-consuming post-processing. We delve into challenges stemming from inconsistently learned feature fields resulting from 2D segmentation obtained independently from a 3D scene. 3D segmentation accuracy deteriorates when 2D segmentation results across the views, primary cues for 3D segmentation, are in conflict. To overcome these issues, we propose Global Feature-guided Learning (GFL). GFL constructs the clusters of global feature candidates from noisy 2D segments across the views, which smooths out noises when training the features of 3D Gaussians. Our method runs in 10 ms per click, 15 to 130 times as fast as the previous methods, while also significantly improving segmentation accuracy.

**Keywords:** Interactive Segmentation · 3D Gaussian Splatting · 3D Feature Field · Contrastive Learning · View-consistency

## 1 Introduction

Recent progress in neural rendering technologies, such as Neural Radiance Fields (NeRF) [30], along with novel 3D scene representation methods like 3D Gaussian Splatting (3DGS) [17], have significantly impacted the field of photorealistic image synthesis within complex 3D environments. These innovations extend into practical applications, enabling advancements in diverse domains such as virtual

---

[*]Equal contribution.
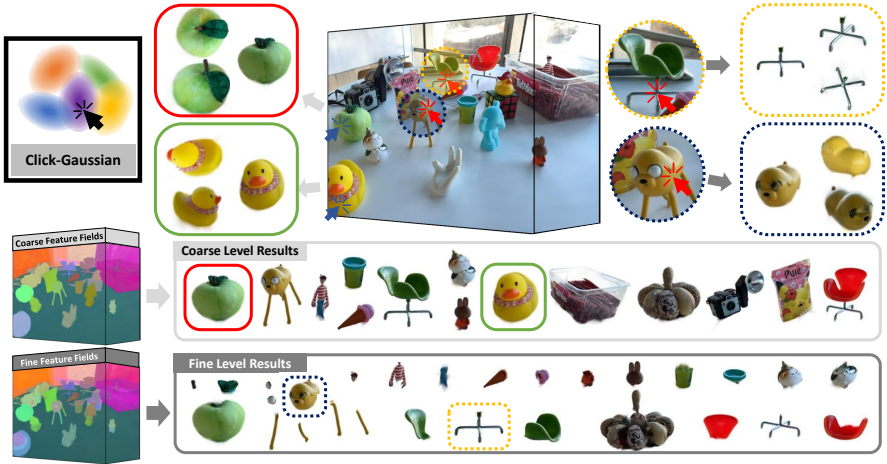
[†]Co-corresponding author.

**Fig. 1:** We present Click-Gaussian, a swift and precise method for interactive segmentation of 3D Gaussians using two-level granularity feature fields derived from 2D segmentation masks. Once trained, it enables users to select and segment desired objects at coarse and fine levels with a single click, completing the process within 10 ms.

and augmented reality [16, 45], digital content creation [25, 33, 54], and real-time rendering [8] for interactive systems that demand both high fidelity and efficiency. For these applications, accurately and efficiently segmenting objects within scenes is important [16], and presents ongoing challenges, particularly in distinguishing elements within diverse 3D environments.

Recently, various segmentation methods [5, 48, 53] based on 3DGS have been proposed, leveraging advantages of 3DGS such as enhanced rendering efficiency and superior reconstruction quality. For instance, some methods [5, 53] learn feature fields of 3D Gaussians that are aligned to the semantic representations from a foundational model like Segment Anything Model (SAM) [20]. This approach enables explicit segmentation of 3D scenes via 3D feature fields, which is crucial for supporting real-time applications and ensuring precise manipulation of intricate environments across diverse tasks. However, these methods face challenges in learning distinguishable feature fields in a scene, necessitating extensive post-processing to achieve clear segmentation. This reliance on time-consuming post-processing significantly impedes the efficiency benefits of 3DGS, creating a bottleneck for applications requiring rapid and direct manipulation of 3D scenes. An alternative approach [48] addresses 3D segmentation by utilizing object tracking mechanisms [10] to pre-assign SAM-based segment identities. However, this method's efficacy is contingent on successful object tracking, potentially excluding untracked objects from the segmentation process. This limitation suggests the potential benefit of developing more robust segmentation techniques capable of comprehensively handling diverse objects within complex scenes.

These considerations motivate the exploration of 3D segmentation methods that provide distinguishable feature fields without extensive post-processing. Progress in these techniques could significantly improve real-time interaction

with 3D scenes, thereby enabling more intuitive and responsive experiences in 3D object manipulation tasks. Such advancements have the potential to not only enhance 3D scene editing capabilities but also broaden the practical applications of 3D scene representation across diverse fields.

In this study, we propose Click-Gaussian, depicted in Fig. 1, as a practical and efficient method for interactive segmentation of 3D Gaussians pre-trained on real-world scenes. By elevating 2D segmentation masks extracted from the SAM into enriched 3D Gaussian's augmented features with two-level granularity, *i.e.*, coarse and fine levels, our approach facilitates fine-grained segmentation. The two-level granularity enables Click-Gaussian to capture scene elements at different scales in 3D environments, enhancing the precision and details of segmentation outcomes. This is achieved through the incorporation of a granularity prior for Gaussian's features, coupled with the employment of a contrastive feature learning method based on the SAM's 2D segmentation masks. Additionally, a significant hurdle in this process is the inconsistency of 2D masks across different views, which impedes the training of consistent and distinguishable semantic features. To address this issue, we introduce Global Feature-guided Learning (GFL), a novel strategy that systematically aggregates global feature candidates across the training views to coherently inform the development of 3D feature fields. The GFL technique enhances the robustness and reliability of feature learning, mitigating the impact of inherent ambiguities present in individual 2D segmentation masks. We demonstrate the effectiveness of Click-Gaussian through comprehensive experiments on complex real-world scenes, evaluating both segmentation accuracy and computational efficiency. Our results indicate that this approach offers a promising solution for precise and rapid 3D scene manipulation, potentially facilitating applications across various domains.

In summary, our key contributions are as follows:

– We propose Click-Gaussian, which enables interactive segmentation of any 3D Gaussians by utilizing two-level feature fields derived from 2D segmentation masks using contrastive learning methods and a granularity prior.
– To tackle the issue of 2D mask inconsistency across training views, we propose GFL, a novel approach that gathers global feature candidates from an entire scene to consistently guide Gaussian's feature learning.
– Through extensive experiments on complicated real-world scenes, we confirm the effectiveness of our approach, demonstrating its suitability for interactive segmentation by significantly enhancing accuracy and processing time.

## 2   Related Work

**3D Gaussian Representations.** 3D Gaussian Splatting [17] has emerged as a promising method for real-time scene rendering, offering superior visual quality. This has inspired research [11, 27, 46, 47] into dynamic scene reconstruction, leveraging its fast rendering capabilities through the design of deformation fields [11, 46, 47]. Moreover, the research has expanded into 3D [9, 42, 50] and 4D [25, 33] content generation by incorporating diffusion models [26, 35]. These

studies demonstrate the efficient rendering and high visual fidelity of 3D Gaussian representations in various applications. Our study further extends these capabilities by focusing on the segmentation of 3D Gaussians, while maintaining their inherent advantages.

**Feature Distillation for 3D Segmentation.** Recent approaches to 3D segmentation can be broadly categorized into two main strategies: feature distillation and mask-lifting techniques. Feature distillation approaches [7, 13, 18, 22, 43, 53] aim to transfer high-dimensional features from 2D vision foundation models [4, 32] into 3D representations. For instance, DFFs [22], N3F [43], and ISRF [13] utilize DINO [4], while LeRF [18] employs CLIP [32]. However, these foundational models, not specifically designed for segmentation tasks, make such approaches struggle to achieve fine-grained segmentation. More recent studies, like those by Chen *et al.* [7] and Feature3DGS [53], distill SAM's encoder features into 3D and use the SAM's decoder to interpret 2D rendered feature maps for segmentation. However, the computational demands of SAM's decoder limit real-time interactive 3D segmentation. In contrast, our approach achieves finer segmentation in real-time by lifting SAM-generated 2D masks to 3D space.

**2D Mask-lifting for 3D Segmentation.** In addition to feature distillation approaches, recent studies have explored lifting 2D segmentation masks into 3D space [5, 6, 19, 31, 34, 48, 52]. For instance, SA3D [6] and NVOS [34] utilize user prompts (*e.g.*, points or scribbles) to derive segmentation masks for a target object in reference views, subsequently training a neural field with these masks for object segmentation. MVseg [31] uses a video segmenter [4] to get multi-view masks. While the aforementioned approaches focus on single-object segmentation, other studies [5, 19, 48, 52] have developed methods for segmenting multiple objects simultaneously. OmniSeg3D [52] trains a feature field using a hierarchical contrastive learning method with 2D segmentation masks, achieving fine-grained segmentation by adjusting cosine similarity thresholds. GARField [19] addresses inconsistent SAM-generated masks across views by introducing a scale-conditioned feature field. However, both use NeRF-based structures, which face computational challenges during rendering, limiting real-time performance.

In the context of 3DGS, SAGA [5] also employs a contrastive learning method with SAM-generated masks. It projects SAM's features into a low-dimensional space via a trainable MLP, imitating these features to address inconsistency issues. However, the distilled SAM's feature is detrimental to the segmentation method that uses feature cosine similarity at the inference stage, requiring extensive post-processing for accuracy. Gau-Group [48] takes a different approach, applying a zero-shot tracker [10] to address mask inconsistencies under the assumption that training images form a video sequence. This assumption, though, limits generalizability, and the method struggles with untracked SAM masks. In contrast, our proposed Global Feature-guided Learning (GFL) method ensures view-consistent training signal by leveraging globally aggregated feature candidates throughout a scene without assuming sequential image inputs.
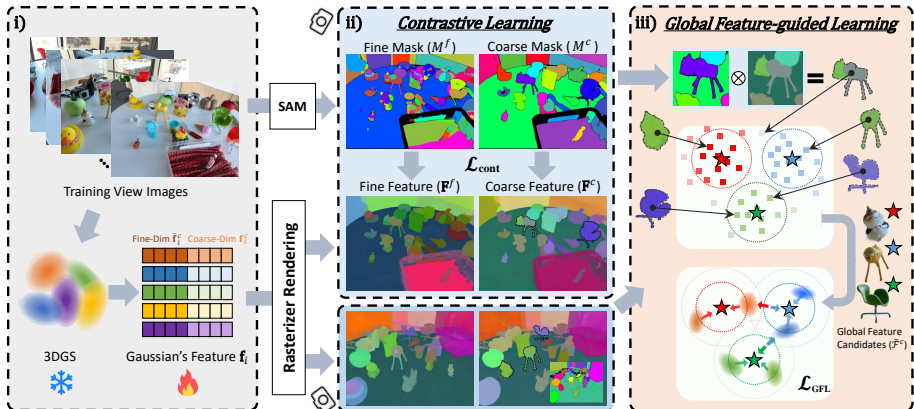
**Fig. 2:** Overview of the proposed method. i) Our approach augments pre-trained 3D Gaussians with two-level granularity features $\mathbf{f}_i$. ii) These features are trained through contrastive learning, utilizing 2D rendered feature maps $\mathbf{F}$ and their corresponding SAM-generated masks $M$. iii) To address inconsistencies in mask signals across views, we introduce a Global Feature-guided Learning approach. For clarity, Global Feature-guided Learning at the fine level is omitted from the illustration.

## 3  Methods

We propose Click-Gaussian, a 3D segmentation method that augments pre-trained 3D Gaussians with effective and distinct 3D feature fields, enabling real-time segmentation capabilities for 3D Gaussian representations. To achieve this, we initially utilize the *automatic mask generation module* of SAM [20] for all training views of a scene, then organize generated masks based on their segment areas to derive coarse and fine level masks for each image. The information from these two-level masks is then incorporated into 3D Gaussians by splitting each Gaussian's feature space using a granularity prior, facilitating the representation of both levels of detail (Sec. 3.2). We train these augmented features through contrastive learning, applied to 2D rendered feature maps in conjunction with the masks (Sec. 3.3). To enhance the consistency of feature learning across different viewpoints, we propose **G**lobal **F**eature-guided **L**earning (GFL), which aggregates global feature candidates across the scene during training (Sec. 3.4). Additionally, we employ several regularization methods in our training process to further stabilize and refine the training of Click-Gaussian's features (Sec. 3.5). The comprehensive methodology is illustrated in Fig. 2.

### 3.1  Preliminary: 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) represents a 3D scene with explicit 3D Gaussians and uses differentiable rasterizer [17] for rendering. Formally, given a training image set $\mathcal{I} = \{I^v\}_{v=1}^{V}$ with camera poses, it aims to learn a set of 3D Gaussians $G = \{g_i\}_{i=1}^{N}$, where $V$ is the number of training images, $N$ is the number of Gaussians, and $g_i = \{\mathbf{p}_i, \mathbf{s}_i, \mathbf{q}_i, o_i, \mathbf{c}_i\}$ is $i$-th Gaussian's trainable parameters.

Here, $\mathbf{p}_i \in \mathbb{R}^3$ is each Gaussian's center position. A scaling factor $\mathbf{s}_i \in \mathbb{R}^3$ and a quaternion $\mathbf{q}_i \in \mathbb{R}^4$ are used to represent each Gaussian's 3D covariance. $o_i \in \mathbb{R}$ is an opacity value, and $\mathbf{c}_i$ is a color represented with spherical harmonics coefficients [36]. After projecting 3D Gaussians onto 2D image space with a given camera pose, 3DGS uses the rasterizer to compute a color $\mathbf{C}$ on a pixel by performing $\alpha$-blending [23,24] on $\mathcal{N}$ depth-ordered points overlapping the pixel:

$$\mathbf{C} = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i T_i, \tag{1}$$

where $\alpha_i$ is calculated by evaluating the influences of each projected Gaussian using splatted 2D covariance [51], opacity $o_i$, and pixel distance, and $T_i = \prod_{j=1}^{i-1}(1 - \alpha_j)$ is the transmittance.

## 3.2   Feature Fields of Click-Gaussian

Click-Gaussian operates by equipping each 3D Gaussian in a scene with additional features for segmentation. Specifically, given a 3D Gaussian $g_i$, each Gaussian is augmented with a $D$-dimensional feature vector $\mathbf{f}_i \in \mathbb{R}^D$ for 3D segmentation, resulting in $\tilde{g}_i = g_i \cup \{\mathbf{f}_i\}$. We split $\mathbf{f}_i$ into $\mathbf{f}_i^c \in \mathbb{R}^{D^c}$ and $\bar{\mathbf{f}}_i^c \in \mathbb{R}^{D-D^c}$, enabling Click-Gaussian to learn features well on both the coarse and fine level masks. We use $\mathbf{f}_i^c$ as a coarse-level feature, and $\mathbf{f}_i^f = \mathbf{f}_i^c \oplus \bar{\mathbf{f}}_i^c$, not $\bar{\mathbf{f}}_i^c$, as a fine-level feature where $\oplus$ is a concatenate function. This is motivated by the intrinsic dependency between two levels in the real world, called granularity prior (*e.g.*, if two objects $A$ and $B$ are different at the coarse level, then each fine part $a \subset A$ and $b \subset B$ are naturally different), to make fine-level feature learning more effective. For our experimental setup, we set $D^c = 12$ and $D = 24$ and freeze other parameters of Gaussians except features. Using the rasterizer, we can compute two-level features $\mathbf{F}^l$ on a pixel, akin to the method outlined in Eq. (1):

$$\mathbf{F}^l = \sum_{i \in \mathcal{N}} \mathbf{f}_i^l \alpha_i T_i, \tag{2}$$

where $l = \{f, c\}$ is the granularity level. The computation of two-level features for each pixel is conducted in a single forward pass.

## 3.3   Contrastive Learning

We use cosine similarity based contrastive learning to train distinctive features with a set of two-level masks. To illustrate this concretely, consider a two-level mask $M^l \in \mathcal{M}$ for a training image $I \in \mathcal{I}$, where $l = \{f, c\}$ is the granularity level. For pixels $p_1$ and $p_2$, if their mask values are the same, *i.e.*, $M_{p_1}^l = M_{p_2}^l$, we aim to maximize the cosine similarity between their rendered features:

$$\mathcal{L}_{\text{pos}}^{\text{cont}} = -\frac{1}{|P_1||P_2|} \sum_{l}^{\{f,c\}} \sum_{p_1}^{P_1} \sum_{p_2}^{P_2} \mathbb{1}\left[M_{p_1}^l = M_{p_2}^l\right] \mathbf{S}^l(p_1, p_2), \tag{3}$$

where $\mathbb{1}$ is the indicator function, $P_1$ and $P_2$ are the set of sampled pixels, $|\cdot|$ is the number of elements in a set, and $\mathbf{S}^l(p_1, p_2) = \langle \mathbf{F}_{p_1}^l, \mathbf{F}_{p_2}^l \rangle$ is the cosine similarity between rendered features of two pixels. Conversely, for pixels with different mask values, *i.e.*, $M_{p_1}^l \neq M_{p_2}^l$, we constrain their rendered features' cosine similarity to not exceed a specified margin, $\tau^l$:

$$\mathcal{L}_{\text{neg}}^{\text{cont}} = \frac{1}{|P_1||P_2|} \sum_l^{\{f,c\}} \sum_{p_1}^{P_1} \sum_{p_2}^{P_2} \mathbb{1}\left[M_{p_1}^l \neq M_{p_2}^l\right] \mathbb{1}\left[\mathbf{S}^l(p_1, p_2) > \tau^l\right] \mathbf{S}^l(p_1, p_2). \quad (4)$$

Considering that two points may represent distinct parts at the fine level yet be classified as the same object at the coarse level, we apply stop gradient operations, *sg*, to the coarse-level components during optimization for negative contrastive loss on fine-level features: $\mathbf{F}^f = sg(\mathbf{F}^c) \oplus \bar{\mathbf{F}}^c$. This method effectively focuses the training process on elements critical for discerning fine-level distinction. We set the margins $\tau^f = 0.75$ and $\tau^c = 0.5$ for all experimental settings. The total contrastive learning loss is defined as:

$$\mathcal{L}_{\text{cont}} = \mathcal{L}_{\text{pos}}^{\text{cont}} + \lambda_{\text{neg}}^{\text{cont}} \mathcal{L}_{\text{neg}}^{\text{cont}}, \quad (5)$$

where $\lambda_{\text{neg}}^{\text{cont}}$ is a hyperparameter for balancing the two losses.

### 3.4 Global Feature-guided Learning

Click-Gaussian's features, despite being trained through contrastive learning, face challenges due to inconsistencies in SAM-generated masks across training viewpoints. This issue arises from the independent use of the masks in each view, potentially leading to unreliable training signals. To address this, we propose Global Feature-guided Learning (GFL), a method that continuously acquires global feature candidates to provide non-conflicting and reliable supervision.

**Global Feature Candidates.** After a specified number of training iterations, we calculate the average features for each two-level mask across all training views. This is accomplished by rendering 2D feature maps and applying average pooling to each mask for all training views. Formally, for two-level masks for all viewpoints $\mathcal{M}^l = \{M^{l,v}\}_{v=1}^V$, where $l = \{f, c\}$ denotes granularity level, $V$ is the number of training viewpoints, and $M^{l,v} \in \mathbb{Z}^{H \times W}$ represents a mask for viewpoint $v$ at level $l$, the average features are calculated as follows:

$$\mathcal{F}^l = \left\{ \bar{F}_s^{l,v} \in \mathbb{R}^{D^l} \,\middle|\, \bar{F}_s^{l,v} = \frac{1}{|\mathcal{P}_s^{l,v}|} \sum_{p \in \mathcal{P}_s^{l,v}} \mathbf{F}_p^{l,v}, 1 \leq s \leq \max_v M^{l,v} \right\}. \quad (6)$$

Here, $\mathcal{P}_s^{l,v} = \{p \mid M_p^{l,v} = s\}$ is a set of pixels with the same segment identiy (ID) in mask $M^{l,v}$, and $D^l$ is the feature dimension at level $l$. This average pooling procedure is done rapidly without gradient calculation, thanks to the

real-time rendering speed of 3DGS at inference time. We then obtain $C^l$ global feature candidates for each level, denoted as $\tilde{\mathcal{F}}^l$, across a scene by applying the HDBSCAN [28] clustering algorithm to each set $\mathcal{F}^l$. These global feature candidates are periodically updated to obtain the latest global features. Notably, as these global clusters are derived by grouping rendered features from noisy 2D segments across all views, they become the most representative features for the entire scene, effectively mitigating inconsistencies in the SAM-generated masks.

**Global Feature-guided Learning.** Global feature candidates enable supervision of Click-Gaussian in a view-consistent manner. For a Gaussian's feature $\mathbf{f}_i^l$, we guide the feature to belong to a specific global cluster and to be far from the others. This involves identifying $c_i^l$, the cluster ID where the $i$-th Gaussian's feature is most likely to belong at level $l$: $c_i^l = \arg\max_c \tilde{\mathbf{S}}^l(i, c)$, where $\tilde{\mathbf{S}}^l(i, c) = \langle \mathbf{f}_i^l, \tilde{\mathcal{F}}_c^l \rangle$ is the cosine similarity between the $i$-th Gaussian's feature and global cluster feature with ID $c$. The GFL loss function for supervising the Gaussian's feature to belong to the most likely global cluster is defined as:

$$\mathcal{L}_{\text{pos}}^{\text{GFL}} = -\frac{1}{N} \sum_l^{\{f,c\}} \sum_i^N \mathbb{1}\left[\tilde{\mathbf{S}}^l(i, c_i^l) > \tau^g\right] \tilde{\mathbf{S}}^l(i, c_i^l). \tag{7}$$

Here, $\tau^g$ is a threshold for determining whether to belong to the cluster, and we set $\tau^g = 0.9$ across all our experiments. Conversely, the GFL loss function guiding the Gaussian's feature away from other global clusters is defined as:

$$\mathcal{L}_{\text{neg}}^{\text{GFL}} = \frac{1}{N} \sum_l^{\{f,c\}} \sum_i^N \frac{1}{C^l} \sum_{c \neq c_i^l}^{C^l} \mathbb{1}\left[\tilde{\mathbf{S}}^l(i, c) > \tau^l\right] \tilde{\mathbf{S}}^l(i, c), \tag{8}$$

where $\tau^l$ is described in Eq. (4). The total GFL loss is thus formulated as:

$$\mathcal{L}_{\text{GFL}} = \mathcal{L}_{\text{pos}}^{\text{GFL}} + \mathcal{L}_{\text{neg}}^{\text{GFL}}. \tag{9}$$

Applying the GFL loss directly to Gaussian's features using global clusters enhances their distinctiveness and noise robustness through reliable supervision, which is vital for accurate 3D segmentation as shown in Sec. 4.3.

## 3.5 Regularization

**Hypersphere Regularization.** Features with excessively large norms underestimate the participation of other features in the rendering process in Eq. (2), impeding effective learning of all Gaussian's features. To prevent any single Gaussian's feature from dominating in the $\alpha$-blending process [23, 24] of the feature rendering, similar to [52], we constrain Gaussian's features to lie on the surface of the hypersphere:

$$\mathcal{L}_{\text{3D-norm}} = \frac{1}{N} \sum_{i=1}^N \left( \|\mathbf{f}_i^c\|_2 - 1 \right)^2 + \left( \|\bar{\mathbf{f}}_i^c\|_2 - 1 \right)^2. \tag{10}$$

**Rendered Feature Regularization.** Due to the hypersphere regularization in Eq. (10), each Gaussian's feature of level $l$ lie on the surface of a hypersphere of radius $r^l$, with $r^c = 1$ for coarse and $r^f = \sqrt{2}$ for fine levels. However, the norm of the rendered feature, $||\mathbf{F}_p^l||_2$, is less than $r^l$, as feature vectors $\mathbf{f}_i$ in different directions are integrated by Eq. (2). This implies that Gaussian's features $\mathbf{f}_i$ contributing to $\mathbf{F}_p^l$ for a single pixel (*i.e.*, the same object) can vary. To ensure all contributing features $\mathbf{f}_i$ for rendering $\mathbf{F}_p^l$ are aligned in the same direction, we apply the following regularization on the rendered feature:

$$\mathcal{L}_{\text{2D-norm}} = \frac{1}{HW} \sum_l^{\{f,c\}} \sum_p^{HW} \left( ||\mathbf{F}_p^l||_2 - r^l \right)^2 . \tag{11}$$

**Spatial Consistency Regularization.** Following the approach of [48], we leverage 3D spatial information to ensure that proximate Gaussians exhibit similar features. At the outset of training, we construct a KD-tree [1] using the 3D positions of Gaussians to facilitate efficient queries for spatially proximate neighbors. Throughout the training process, we sample $N_s$ Gaussians and adjust their features to align with those of their $K$-nearest neighbors in 3D space:

$$\mathcal{L}_{\text{spatial}} = -\frac{1}{N_s K} \sum_i^{N_s} \sum_k^K \langle \mathbf{f}_i, \mathbf{f}_k \rangle, \tag{12}$$

where $\langle \cdot, \cdot \rangle$ is the cosine similarity operation. For all experiments, we set $N_s = 100,000$ and $K = 5$. Finally, our total objective for training Click-Gaussian is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{cont}} + \lambda_1 \mathcal{L}_{\text{GFL}} + \lambda_2 \mathcal{L}_{\text{3D-norm}} + \lambda_3 \mathcal{L}_{\text{2D-norm}} + \lambda_4 \mathcal{L}_{\text{spatial}}, \tag{13}$$

where $\lambda_1$, $\lambda_2$, $\lambda_3$, and $\lambda_4$ are hyperparameters balancing the respective loss terms.

## 4 Experiments

### 4.1 Experimental Settings

We implemented Click-Gaussian using the 3DGS codebase [17], adopting its default settings for pre-trained Gaussians. The hyperparameters were set as follows: $\lambda_{\text{neg}}^{\text{cont}} = 0.1$, $\lambda_1 = 10.0$, $\lambda_2 = 0.2$, $\lambda_3 = 0.2$, and $\lambda_4 = 0.5$. We employed the Adam optimizer with a learning rate of 0.01 for Gaussian's features. For contrastive learning, we sampled 10k pixels for each training iteration using importance sampling based on mask pixel count. In HDBSCAN, we set the epsilons for coarse and fine features clustering to $1 \times 10^{-2}$ and $1 \times 10^{-3}$, respectively, with the minimum cluster size proportional to the number of training views. We trained Click-Gaussian for 3,000 iterations, incorporating Global Feature-Guided learning from the 2,000th iteration onward. The entire training process took approximately 13 minutes on an NVIDIA RTX A5000 GPU.

| Model | Figurines | | | Ramen | | | Teatime | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Coarse | Fine | All | Coarse | Fine | All | Coarse | Fine | All | Coarse | Fine | All |
| Gau-Group [48] | 69.7 | 17.1 | 33.8 | 77.0 | 75.1 | 75.8 | 71.7 | 7.78 | 43.3 | 72.8 | 31.5 | 49.9 |
| OmniSeg3D [52] | 87.1 | 28.0 | 46.8 | 77.3 | 72.1 | 74.2 | 73.8 | 46.9 | 61.8 | 79.4 | 46.9 | 60.9 |
| Feature3DGS [53] | 70.4 | 49.2 | 55.9 | 65.9 | 73.0 | 70.2 | 60.6 | 68.4 | 64.0 | 65.6 | 63.5 | 63.4 |
| GARField [19] | 89.2 | 59.4 | 70.1 | 75.7 | **86.9** | 82.4 | 77.8 | 67.8 | 73.2 | 80.9 | 71.4 | 75.2 |
| Ours | **93.2** | **75.9** | **81.4** | **90.9** | 86.7 | **88.4** | **83.2** | **90.4** | **86.4** | **89.1** | **84.3** | **85.4** |

**Table 1:** Quantitative comparison with baselines on LERF-Mask Dataset. Segmentation performance was evaluated in mIoU with coarse and fine level masks for each of three scenes, and 'All' reflects mean value across all objects and levels. 'Average' shows mean performance across scenes. Our approach outperforms all prior works, especially in fine-level segmentation.

We used the official code's *automatic mask generation module* for SAM mask creation, which extracts masks without distinguishing levels, allowing us to get only the highest-confidence segments in an image. The coarse and fine masks are then assigned per pixel by the largest and smallest segments, respectively (see supplementary materials for details). To evaluate our approach's 3D segmentation performance using these masks, we employ two public real-world datasets: LERF-Mask dataset [48] and SPIn-NeRF dataset [31]. The LERF-Mask dataset comprises three scenes [18] with manually annotated ground truth masks for large objects. We further annotated several masks for fine objects within each scene using Make-Sense [39] to evaluate fine-grained segmentation performance. The SPIn-NeRF dataset offers multi-view masks for single objects in the widely used NeRF datasets [12, 21, 29, 30, 49] that include both forward-facing [29] and 360-degree inward-facing [12, 21, 30, 49] setups.

**Segmentation Procedure.** After training for each dataset, we computed global feature candidates for subsequent segmentation tasks. We adopted a label propagation method [5, 6, 48] for 3D multi-view segmentation task by using a ground truth 2D mask from a reference view to identify target object IDs and evaluated generated 2D masks. Specifically, we calculated cosine similarities between the rendered 2D feature map and the pre-computed global clusters, applying a 0.9 threshold to determine cluster IDs of the target object. Using the IDs, we generated 2D masks for test views and evaluated performance by calculating mIoU between rendered and ground truth masks. This approach aligns with common scenarios where users segment a target object in the reference view to obtain results for other views. For 3D Gaussian extraction task, we used cosine similarity between rendered features corresponding to user-provided point prompts and two-level global feature candidates to retrieve clicked objects efficiently.

## 4.2    Comparisons

**Comparison on LERF-Mask Dataset.** To demonstrate Click-Gaussian's segmentation superiority, we compared it with various baselines using the LERF-Mask dataset. For Gau-Group, target object IDs in the reference view were
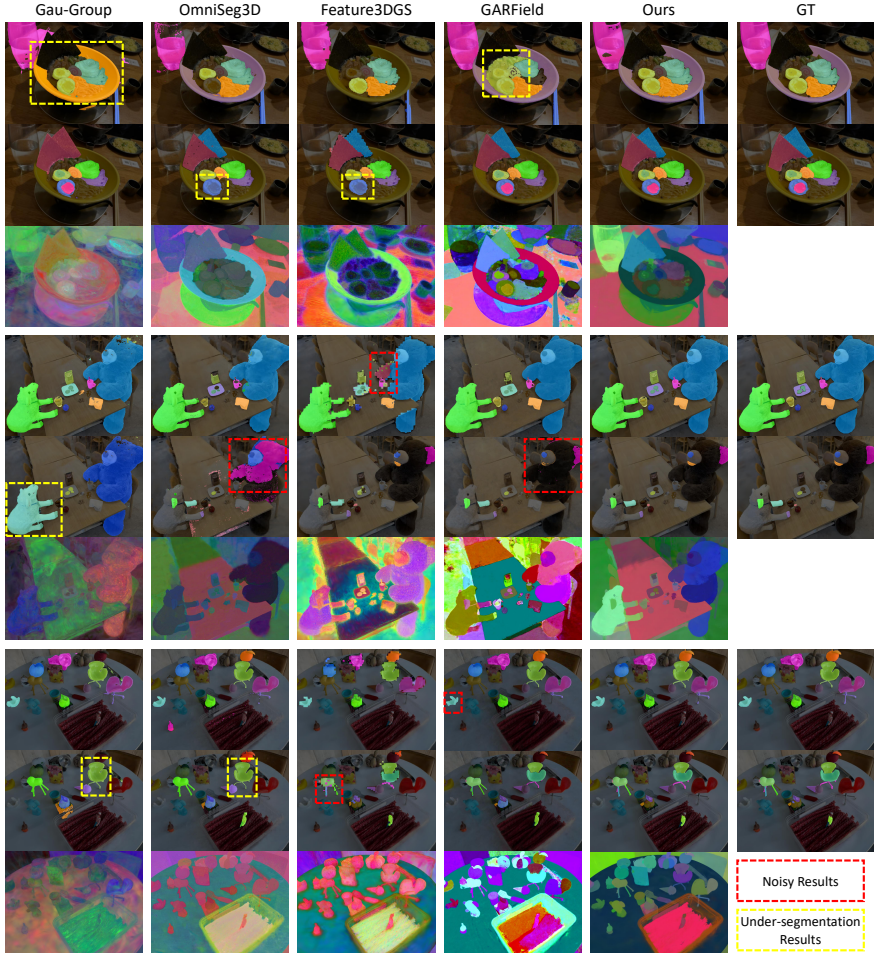
**Fig. 3:** Comparison with baselines on LERF-Mask Dataset. The results are displayed in three lines per scene (Teatime, Ramen, and Figurines in order). Each scene's first two rows show coarse and fine level segmentation results, respectively, and the third row shows the PCA visualizations of each model's finest-level feature field. Our approach demonstrates superior segmentation ability in both coarse and fine levels. Red and yellow boxes indicate noisy and under-segmentation results, respectively.

identified using a classifier and a ground truth mask [48]. OmniSeg3D's segmentation involved adjusting the cosine similarity threshold [52] from 0 to 1 in 0.01 increments, finding the optimal threshold for each target object. Feature3DGS utilized a rendered 2D feature map and the SAM's decoder for segmentation [53], selecting the best match for the target object. GARField applied a NeRF-based scale-conditioned field, selecting the best scale (0 to 1 in 0.05 steps) for each target object in the reference view [19].

As shown in Tab. 1, our method outperforms all baselines. Gau-Group underperforms in fine-level segmentation, due to its tracking methodology limitations.
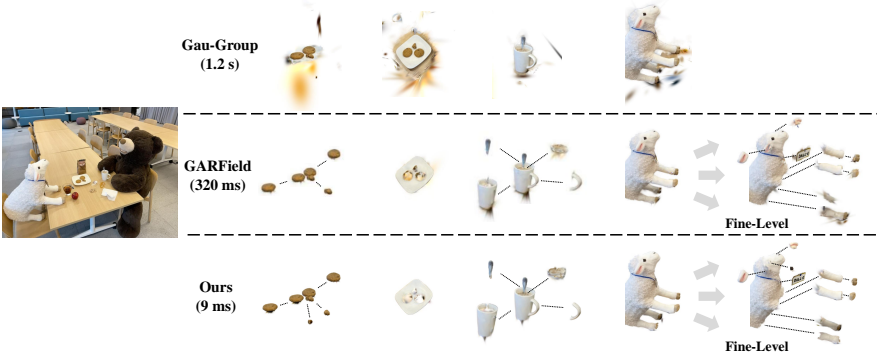
**Fig. 4:** Comparison with Gau-Group and GARField. Our approach performs more detailed and cleaner extractions of Gaussians, up to 130 times faster than other baselines.
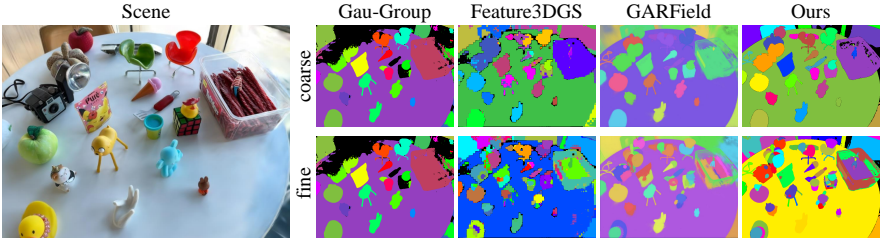


**Fig. 5:** Comparison for automatic segmentation of everything on novel views. Our method shows more exact and fine-grained results against baselines. For detailed experimental procedures, please refer to the supplementary materials. Gau-Group, unable to differentiate levels, presents identical coarse and fine segmentation results.

OmniSeg3D shows improved performance than Gau-Group using a hierarchical field separated by cosine similarity, but lacks detailed segmentation. Feature3DGS enhances performance via SAM's decoder but requires a longer processing time. GARField shows comparable results with scale-conditioned feature fields, yet struggles with finer segmentation accuracy. Fig. 3 qualitatively demonstrates our method's superiority, together with PCA visualizations of the finest level feature fields of all models.

We also qualitatively compared 3D Gaussian extraction performance using user-provided point prompts. As shown in Fig. 4, our method yields finer and cleaner extractions while operating up to 130 times faster than competing approaches. Additionally, as depicted in Fig. 5, we compared performance in automatic segmentation of everything on a novel view, where Click-Gaussian shows more exact and fine-grained results than all baselines.

**Comparison on SPIn-NeRF Dataset.** We extend our analysis to include additional baselines on the SPIn-NeRF dataset, evaluating our method's performance across various views per scene. The first two baseline models in Tab. 2, MVSeg and SA3D, can segment only a single foreground object at a time.

|  | MVSeg [31] | SA3D [6] | SAGA [5] | Ours |
|---|---|---|---|---|
| mIoU (%) | 90.9 | 92.4 | 88.0 | **94.0** |

**Table 2:** Quantitative comparison with baselines on SPIn-NeRF Dataset. We report average mIoU across ten real-world scenes. Our method outperforms not only MVSeg and SA3D, which segment only a single foreground object per training, but also SAGA.
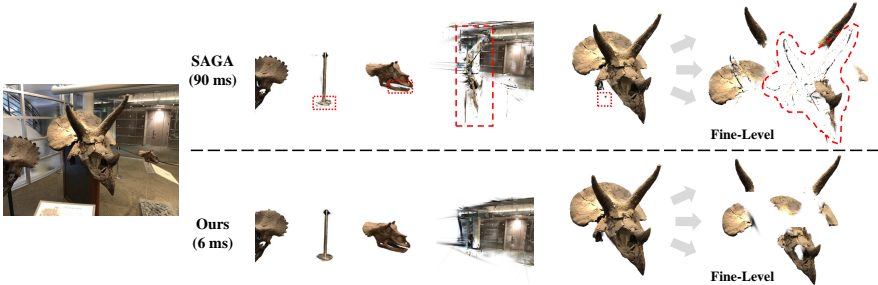


**Fig. 6:** Comparison with SAGA. Our method achieves more precise Gaussian extraction, highlighted by red dotted lines, and runs about 15 times faster than SAGA.

|  | w/o $\mathcal{L}_{\text{2D-norm}}$ | w/o $\mathcal{L}_{\text{3D-norm}}$ | w/o $\mathcal{L}_{\text{spatial}}$ | w/o $\mathcal{L}_{\text{GFL}}$ | w/o prior | Ours |
|---|---|---|---|---|---|---|
| Coarse | **89.3** | 88.9 | 88.5 | 83.4 | 88.8 | 89.1 |
| Fine | 80.8 | 74.1 | 78.5 | 42.3 | 78.3 | **84.3** |
| All | 83.2 $(-2.6\%)$ | 80.3 $(-6.0\%)$ | 82.0 $(-4.0\%)$ | 58.6 $(-31.3\%)$ | 82.1 $(-3.9\%)$ | **85.4** |

**Table 3:** Ablation study to evaluate the contribution of each component. We remove granularity prior (w/o prior), GFL loss (w/o $\mathcal{L}_{\text{GFL}}$), and regularization losses (w/o $\mathcal{L}_{\text{2D-norm}}$, w/o $\mathcal{L}_{\text{3D-norm}}$, and w/o $\mathcal{L}_{\text{spatial}}$) from our complete method to assess each component's impact. Average mIoU values on LERF-Mask Dataset reported.

Click-Gaussian, which efficiently segments any objects within a single training, even outperforms both methods. It also surpasses SAGA, which utilizes low-dimensional SAM's feature fields. In Fig. 6, we further experimented on 3D Gaussian extraction using user-provided point prompts. Our method performed detailed and exact Gaussian extraction about 15 times faster than SAGA which uses time-consuming multiple post-processing steps [5].

## 4.3   Ablation Study

We evaluated the impact of removing granularity prior, GFL loss, and regularization terms from Click-Gaussian. Without the granularity prior (w/o prior), the model learns two-level features independently, lacking collaborative enhancement. As Tab. 3 shows, our complete model outperformed the one without prior in fine-level segmentation. This highlights how the intrinsic dependency between coarse and fine levels aids fine-level feature learning. GFL loss significantly improves fine-level segmentation by effectively addressing inconsistency and ambiguities of SAM-generated fine-level masks. Removing each regularization loss confirmed their collective importance in enhancing segmentation performance.

**Fig. 7:** Application of Click-Gaussian. Click-Gaussian enables interactive object selection within about 10ms for various edits. For text-based editing, we utilized CLIP-based editing methods with the source text *'flower'* and target text *'stained glass flower'*.

### 4.4   Versatile Applications

After training Click-Gaussian's two-level feature fields, various scene manipulation tasks become feasible, including object removal, resizing, repositioning, duplication, and text-based editing. As Fig. 7 shows, the two-level global clusters enable rapid object selection within a scene (approximately 10 ms), facilitating interactive local adjustments to selected Gaussians. Additionally, for text-based editing, we leveraged CLIP-based methods [32,40,44] to get faster results (within about 10 s) than diffusion-based methods [2, 14, 35].

## 5   Conclusion

We present Click-Gaussian, a swift and precise method enabling interactive fine-grained segmentation of pre-trained 3D Gaussians by lifting 2D segmentation masks into 3D feature fields of two-level granularity. Noticing from the intrinsic dependency between coarse and fine levels in the real world, we employ a granularity prior for feature division in the representation of feature fields. To address feature learning hindered by the cross-view inconsistency masks, an inherent issue in lifting 2D masks to 3D, we propose the Global Feature-guided Learning method for more consistent feature field training. Once Click-Gaussian is trained, users can select desired objects at coarse and fine levels more swiftly than previous methods. This enhanced capability has the potential to improve efficient and precise 3D environment modification across various applications.

**Limitations.** Our approach faces limitations due to its reliance on pre-trained 3DGS and the two-level granularity assumption. Feature learning may be hindered if a single Gaussian represents multiple objects, particularly when they are semantically distinct but chromatically similar. The two-level granularity assumption, lacking intermediate levels, could limit efficiency for varying granular levels and complex structures, potentially requiring multiple interactions to select desired segmentation regions.

# References

1. Bentley, J.L.: Multidimensional binary search trees used for associative searching. Commun. ACM **18**(9), 509–517 (sep 1975). https://doi.org/10.1145/361002.361007, https://doi.org/10.1145/361002.361007

2. Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18392–18402 (2023)

3. Brooks, T., Peebles, B., Holmes, C., DePue, W., Guo, Y., Jing, L., Schnurr, D., Taylor, J., Luhman, T., Luhman, E., Ng, C., Wang, R., Ramesh, A.: Video generation models as world simulators (2024), https://openai.com/research/video-generation-models-as-world-simulators

4. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9650–9660 (2021)

5. Cen, J., Fang, J., Yang, C., Xie, L., Zhang, X., Shen, W., Tian, Q.: Segment any 3d gaussians (2024), https://arxiv.org/abs/2312.00860v1

6. Cen, J., Zhou, Z., Fang, J., Shen, W., Xie, L., Jiang, D., Zhang, X., Tian, Q., et al.: Segment anything in 3d with nerfs. Advances in Neural Information Processing Systems **36** (2024)

7. Chen, X., Tang, J., Wan, D., Wang, J., Zeng, G.: Interactive segment anything nerf with feature imitation. arXiv preprint arXiv:2305.16233 (2023)

8. Chen, Z., Funkhouser, T., Hedman, P., Tagliasacchi, A.: Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In: The Conference on Computer Vision and Pattern Recognition (CVPR) (2023)

9. Chen, Z., Wang, F., Liu, H.: Text-to-3d using gaussian splatting. arXiv preprint arXiv:2309.16585 (2023)

10. Cheng, H.K., Oh, S.W., Price, B., Schwing, A., Lee, J.Y.: Tracking anything with decoupled video segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1316–1326 (2023)

11. Cotton, R.J., Peyton, C.: Dynamic gaussian splatting from markerless motion capture reconstruct infants movements. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops. pp. 60–68 (January 2024)

12. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022)

13. Goel, R., Sirikonda, D., Saini, S., Narayanan, P.: Interactive Segmentation of Radiance Fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)

14. Haque, A., Tancik, M., Efros, A., Holynski, A., Kanazawa, A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2023)

15. Hoffstadt, J., Cothren, P., Contributors: Dearpygui. https://github.com/hoffstadt/DearPyGui

16. Jiang, Y., Yu, C., Xie, T., Li, X., Feng, Y., Wang, H., Li, M., Lau, H., Gao, F., Yang, Y., Jiang, C.: Vr-gs: A physical dynamics-aware interactive gaussian splatting system in virtual reality. arXiv preprint arXiv:2401.16663 (2024)

17. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4) (July 2023), https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/

18. Kerr, J., Kim, C.M., Goldberg, K., Kanazawa, A., Tancik, M.: Lerf: Language embedded radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19729–19739 (2023)
19. Kim, C.M., Wu, M., Kerr, J., Goldberg, K., Tancik, M., Kanazawa, A.: Garfield: Group anything with radiance fields (2024)
20. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollar, P., Girshick, R.: Segment anything. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 4015–4026 (October 2023)
21. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: benchmarking large-scale scene reconstruction. ACM Trans. Graph. **36**(4) (jul 2017). https://doi.org/10.1145/3072959.3073599, https://doi.org/10.1145/3072959.3073599
22. Kobayashi, S., Matsumoto, E., Sitzmann, V.: Decomposing nerf for editing via feature field distillation. In: Advances in Neural Information Processing Systems. vol. 35 (2022), https://arxiv.org/pdf/2205.15585.pdf
23. Kopanas, G., Leimkühler, T., Rainer, G., Jambon, C., Drettakis, G.: Neural point catacaustics for novel-view synthesis of reflections. ACM Transactions on Graphics (TOG) **41**(6), 1–15 (2022)
24. Kopanas, G., Philip, J., Leimkühler, T., Drettakis, G.: Point-based neural rendering with per-view optimization. In: Computer Graphics Forum. vol. 40, pp. 29–43. Wiley Online Library (2021)
25. Ling, H., Kim, S.W., Torralba, A., Fidler, S., Kreis, K.: Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. arXiv preprint arXiv:2312.13763 (2023)
26. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023)
27. Luiten, J., Kopanas, G., Leibe, B., Ramanan, D.: Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. arXiv preprint arXiv:2308.09713 (2023)
28. McInnes, L., Healy, J., Astels, S.: hdbscan: Hierarchical density based clustering. Journal of Open Source Software **2**(11),  205 (2017). https://doi.org/10.21105/joss.00205, https://doi.org/10.21105/joss.00205
29. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. ACM Transactions on Graphics (TOG) **38**(4), 1–14 (2019)
30. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021)
31. Mirzaei, A., Aumentado-Armstrong, T., Derpanis, K.G., Kelly, J., Brubaker, M.A., Gilitschenski, I., Levinshtein, A.: Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20669–20679 (2023)
32. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
33. Ren, J., Pan, L., Tang, J., Zhang, C., Cao, A., Zeng, G., Liu, Z.: Dreamgaussian4d: Generative 4d gaussian splatting. arXiv preprint arXiv:2312.17142 (2023)

34. Ren, Z., Agarwala†, A., Russell†, B., Schwing†, A.G., Wang†, O.: Neural volumetric object selection. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022), († alphabetic ordering)
35. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022)
36. Sara Fridovich-Keil and Alex Yu, Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022)
37. Schönberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
38. Schönberger, J.L., Zheng, E., Pollefeys, M., Frahm, J.M.: Pixelwise view selection for unstructured multi-view stereo. In: European Conference on Computer Vision (ECCV) (2016)
39. Skalski, P.: Make Sense. https://github.com/SkalskiP/make-sense/ (2019)
40. Song, H., Choi, S., Do, H., Lee, C., Kim, T.: Blending-nerf: Text-driven localized editing in neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 14383–14393 (October 2023)
41. Tang, J., Chen, X., Wan, D., Wang, J., Zeng, G.: Segment-anything nerf. https://github.com/ashawkey/Segment-Anything-NeRF (2023)
42. Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653 (2023)
43. Tschernezki, V., Laina, I., Larlus, D., Vedaldi, A.: Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In: Proceedings of the International Conference on 3D Vision (3DV) (2022)
44. Wang, C., Chai, M., He, M., Chen, D., Liao, J.: Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3835–3844 (2022)
45. Xu, L., Agrawal, V., Laney, W., Garcia, T., Bansal, A., Kim, C., Rota Bulò, S., Porzi, L., Kontschieder, P., Božič, A., Lin, D., Zollhöfer, M., Richardt, C.: VR-NeRF: High-fidelity virtualized walkable spaces. In: SIGGRAPH Asia Conference Proceedings (2023). https://doi.org/10.1145/3610548.3618139, https://vr-nerf.github.io
46. Yang, Z., Yang, H., Pan, Z., Zhu, X., Zhang, L.: Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. arXiv preprint arXiv:2310.10642 (2023)
47. Yang, Z., Gao, X., Zhou, W., Jiao, S., Zhang, Y., Jin, X.: Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. arXiv preprint arXiv:2309.13101 (2023)
48. Ye, M., Danelljan, M., Yu, F., Ke, L.: Gaussian grouping: Segment and edit anything in 3d scenes. arXiv preprint arXiv:2312.00732 (2023)
49. Yen-Chen, L., Florence, P., Barron, J.T., Lin, T.Y., Rodriguez, A., Isola, P.: NeRF-Supervision: Learning dense object descriptors from neural radiance fields. In: IEEE Conference on Robotics and Automation (ICRA) (2022)
50. Yi, T., Fang, J., Wu, G., Xie, L., Zhang, X., Liu, W., Tian, Q., Wang, X.: Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. arXiv preprint arXiv:2310.08529 (2023)
51. Yifan, W., Serena, F., Wu, S., Öztireli, C., Sorkine-Hornung, O.: Differentiable surface splatting for point-based geometry processing. ACM Transactions on Graphics (TOG) 38(6), 1–14 (2019)
52. Ying, H., Yin, Y., Zhang, J., Wang, F., Yu, T., Huang, R., Fang, L.: Omniseg3d: Omniversal 3d segmentation via hierarchical contrastive learning (2023)

53. Zhou, S., Chang, H., Jiang, S., Fan, Z., Zhu, Z., Xu, D., Chari, P., You, S., Wang, Z., Kadambi, A.: Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. arXiv preprint arXiv:2312.03203 (2023)
54. Zielonka, W., Bagautdinov, T., Saito, S., Zollhöfer, M., Thies, J., Romero, J.: Drivable 3d gaussian avatars. arXiv preprint arXiv:2311.08581 (2023)

# Supplementary Material of
# Click-Gaussian: Interactive Segmentation
# to Any 3D Gaussians

## A    GUI-based Implementation for Click-Gaussian

To showcase interactive segmentation and manipulation using Click-Gaussian, we design a Graphical User Interface (GUI) tool based on DearPyGui [15, 41], a fast and powerful GUI toolkit for Python. As shown in Fig. 8, our GUI is designed to allow users to easily click and segment objects at coarse and fine levels, and provides tools for real-time manipulation tasks such as resizing, translation, removal, and text-based editing for intuitive interaction with the segmented objects. The supplementary video demonstrates the effectiveness of our method in enabling real-time interactive scene manipulation, showcasing its fast and precise 3D segmentation performance. We encourage readers to view this video for a comprehensive understanding of the proposed approach's capabilities.



**Fig. 8:** Graphical User Interface (GUI) for Click-Gaussian.

# B  SAM-based Multi-level Mask Generation

We utilized the official code's *automatic mask generation module* for SAM mask creation, which extracts masks without distinguishing levels, allowing us to get only the highest-confidence segments in an image. These segments are then assigned to two masks by area: if multiple segments are assigned to a single pixel, the coarse-level mask prioritizes the identity of the larger segment, while the fine-level mask favors the identity of the smaller segments. This approach enables us to assign a single mask identity per pixel at each level, facilitating stable contrastive learning.

**Comparative Analysis of Multi-level Mask Strategies.** Our method can adopt SAM's three-level masks (whole, part, and subpart) in two ways: three-level-score and three-level-area. Each approach prioritizes the highest score segment and smallest segment, respectively, for each level. In these cases, we split $\mathbf{f}_i \in \mathbb{R}^{24}$ into three levels of granularity. As shown in Fig. 9, the three-level-area outperforms the three-level-score in fine-level mIoU due to finer-grained mask supervision (*e.g.*, egg white and yolk), demonstrating the efficacy of the area-based prioritization. Additionally, our method using two-level masks surpasses the three-level-area thanks to the mask completeness and training efficiency: It has fewer unassigned identities than the three-level-area and learns feature fields more efficiently with the same feature dimension of 24. For these advantages, we adopt the two-level granularity assumption.
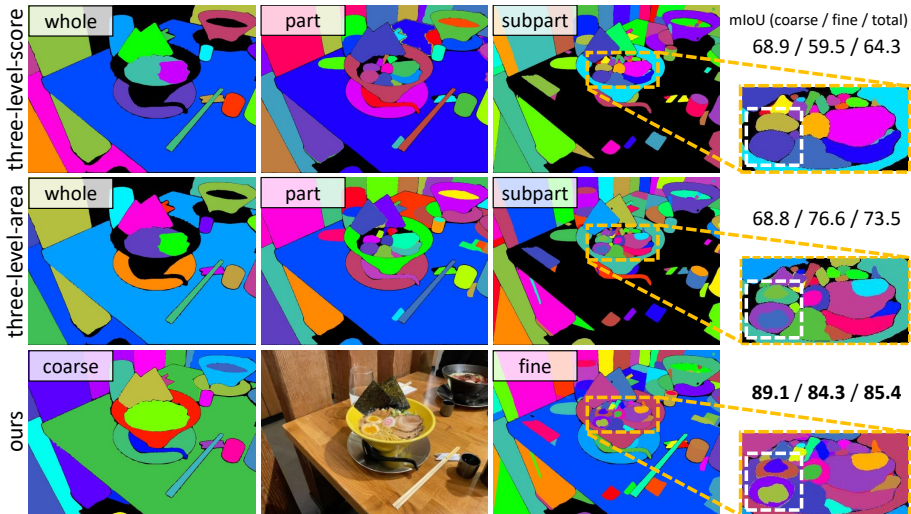


**Fig. 9:** Performance comparison of different mask generation strategies. Black areas indicate pixels with unassigned identities.

# C Annotations for Evaluating Fine-grained Segmentation

We evaluate our approach's segmentation performance using the LERF-Mask dataset [48], a public real-world dataset for 3D segmentation tasks. This dataset comprises three scenes (Figurines, Ramen, and Teatime) [18] with manually annotated ground truth masks for semantically large objects, as shown in the first two rows of Fig. 10. To assess fine-grained segmentation performance, we additionally annotated masks for smaller objects within each scene using Make-Sense [39], a free online image labeling tool, as shown in the last two rows of Fig. 10. This additional annotation is necessary due to the lack of datasets suitable for fine-level comparison. Note that the annotation process was conducted independently from our segmentation experiments.



**Fig. 10:** Annotations for evaluating fine-grained segmentation. The first two rows of each scene show the ground truth annotations for evaluating coarse-level segmentation with two sampled test views. On the other hand, each scene's last two rows show the ground truth annotations for evaluating fine-level segmentation.

# D    Additional Experiments and Results

## D.1    3D Editing in AI-generated Videos

OpenAI recently announced Sora [3], a groundbreaking text-to-video generation model, showing a promising path towards building general-purpose world simulators. These simulators can be further improved by enabling interactive modification of generated realistic environments through accurate and fast 3D segmentation methods like Click-Gaussian, enhancing their functionality and user interaction capabilities. To demonstrate Click-Gaussian's versatility in scene segmentation and manipulation on these generated scenes, we applied our method to videos (Santorini[1] and Snow-village[2]) generated by Sora. As shown in Fig. 11, after pre-training 3DGS on each generated video using COLMAP [37,38], users can flexibly make desired modifications, resulting in more creative and diverse 3D environments with Click-Gaussian.
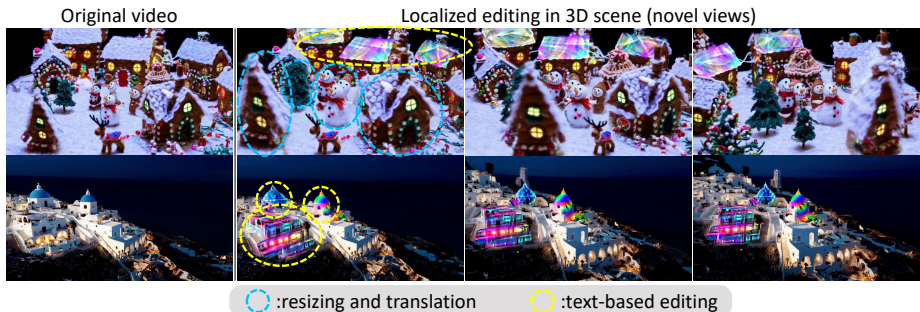


**Fig. 11:** Versatile applications of Click-Gaussian on synthetic videos generated by Sora. After pre-training 3DGS on Sora-generated videos, users can flexibly modify the reconstructed 3D scenes in real-time, including resizing, translation (sky blue circle), and text-based editing (yellow circle). In the Snow-village scene (first row), we manipulated the scene by enlarging and translating three snowmen, two houses, and a tree, while stylizing other house roofs to crystal. In the Santorini scene (second row), we applied text-based editing to buildings, transforming them into cyberpunk neon, crystal, and rainbow styles from the bottom left, respectively.

---

[1] https://cdn.openai.com/sora/videos/santorini.mp4
[2] https://cdn.openai.com/tmp/s/interp/b2.mp4. This video has no official name, so we refer to it as Snow-village.

## D.2 Open-vocabulary 3D Object Localization

Once trained, our method can perform open-vocabulary 3D object localization as shown in Fig. 12, using the obtained global feature candidates, which we call global clusters. Specifically, for all two-level global clusters, we render only the Gaussians corresponding to each cluster in multiple views (10 randomly sampled views) as shown in Fig. 13. We then input these rendered images into the CLIP image encoder [32] to obtain the CLIP embeddings of each cluster. Thanks to the real-time rendering speed of 3DGS, this process of obtaining CLIP embeddings for all global clusters completes in 20–40 seconds, depending on the number of global clusters in the scene. Note that this process only needs to be performed once before any text query. Subsequently, given text queries, open-vocabulary 3D object localization is performed by returning the global cluster with the highest cosine similarity between the obtained image embeddings of all global clusters and the text query embedding. Fig. 12 qualitatively demonstrates that our approach precisely localizes 3D objects for given text queries using globally obtained clusters.
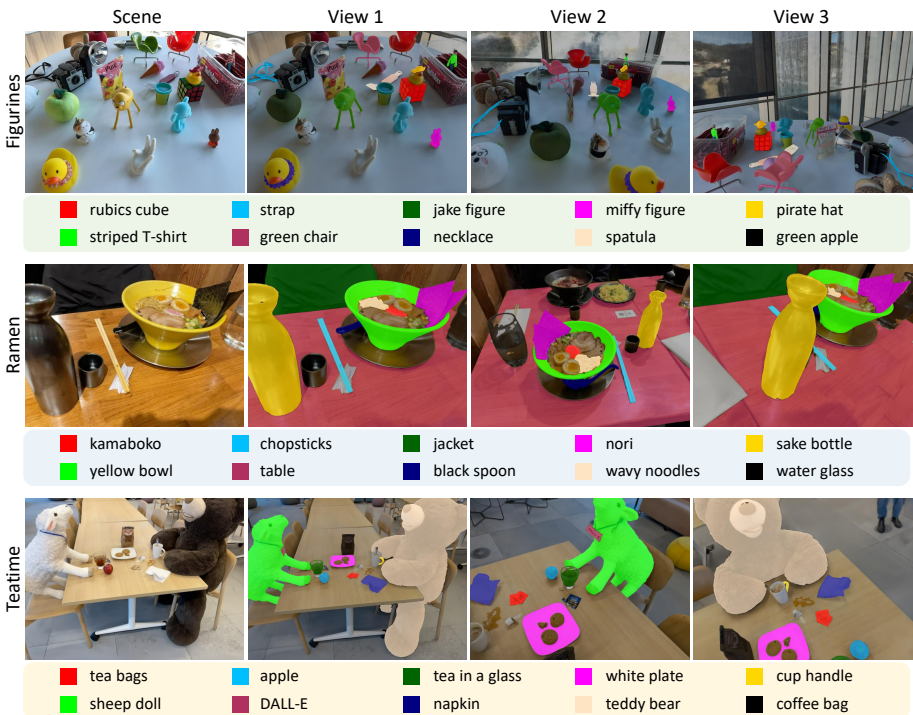


**Fig. 12:** Open-vocabulary 3D object localization results on the LERF-Mask Dataset. Segmentation results are color-overlaid for visualization in three different scenes.
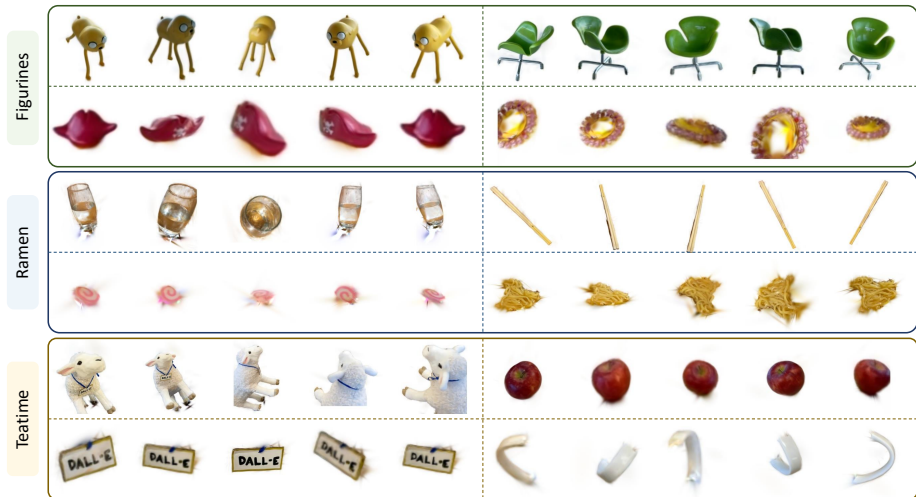
**Fig. 13:** Examples of rendered images using the Gaussians corresponding to each global cluster. Five representative images are shown per cluster for simplicity. These images are used to obtain CLIP embeddings for each cluster via the CLIP image encoder.

### D.3   Additional Results for 3D Segmentation

**Experiments on LeRF Dataset.** In addition to user-guided segmentation, our approach can also automatically segment everything by calculating the cosine similarity between the rendered 2D feature map and global clusters' features, assigning a global cluster ID with the maximum similarity value to each pixel. By performing this process for each of the two granularity levels, we obtain automatic segmentation results at both coarse and fine levels. Figs. 14, 15, and 16 show the results of automatic segmentation for several complicated real-world scenes from the LeRF dataset [18], along with PCA visualizations of rendered feature maps at two levels. These results qualitatively demonstrate that Click-Gaussian achieves high-fidelity, fine-grained segmentation of everything in complex real-world scenes.

**Experiments on SPIn-NeRF Dataset.** We further showcase the 3D multi-view segmentation results on the SPIn-NeRF Dataset using the label propagation method, as illustrated in Fig. 17. These results offer additional examples demonstrating the effectiveness of Click-Gaussian across various real-world scenes.
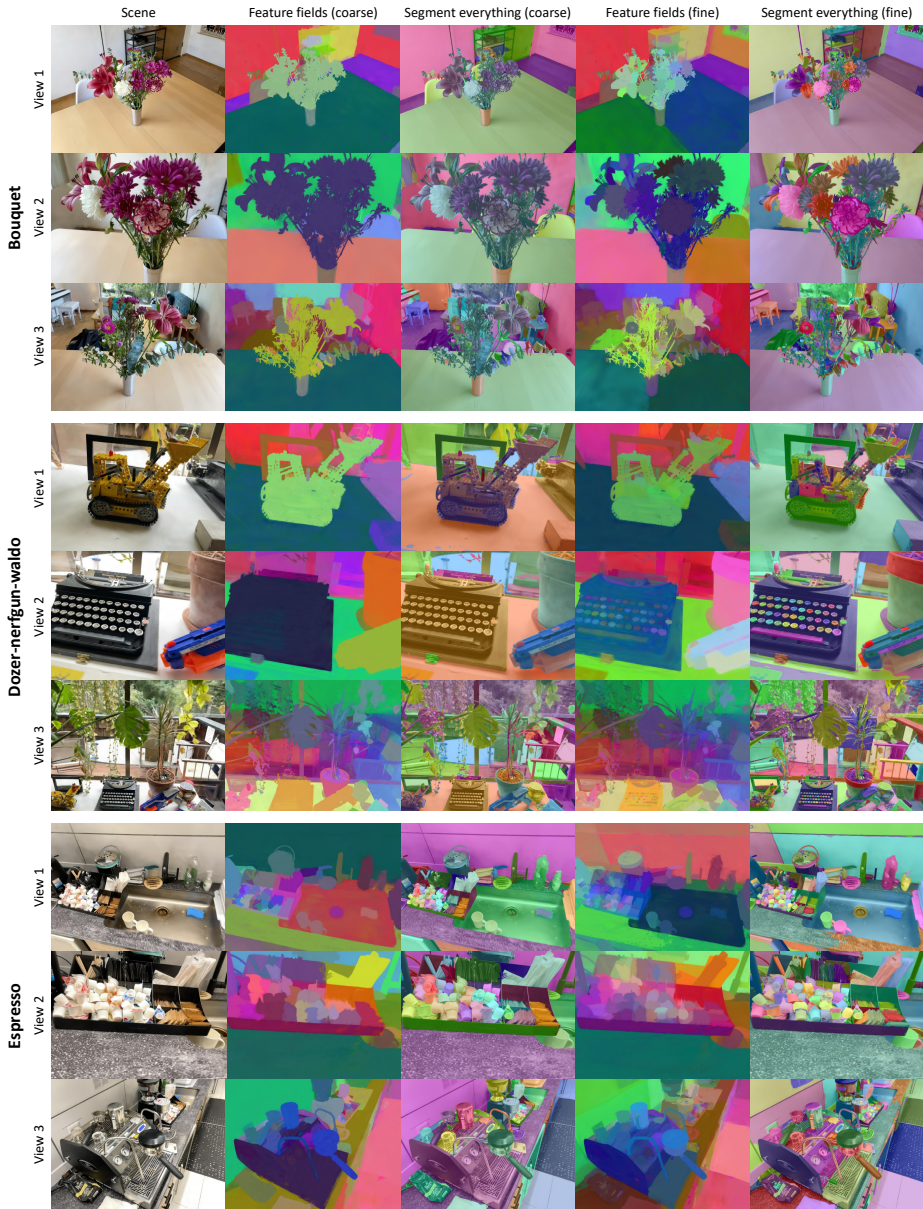
**Fig. 14:** Segmentation of everything results on the LeRF Dataset. We present automatic segmentation results (third and fifth columns) along with PCA visualizations of rendered feature maps (second and fourth columns) at two granularity levels for Bouquet, Dozer-nerfgun-waldo, and Espresso scenes (first column) from the LeRF Dataset. Objects classified with the same ID in the segmentation results share the same overlaid color across the three given views, as each global cluster ID remains consistent throughout a scene.

**Fig. 15:** Segmentation of everything results on the LeRF Dataset. We present automatic segmentation results (third and fifth columns) along with PCA visualizations of rendered feature maps (second and fourth columns) at two granularity levels for Figurines, Fruit-aisle, and Ramen scenes (first column) from the LeRF Dataset. Objects classified with the same ID in the segmentation results share the same overlaid color across the three given views, as each global cluster ID remains consistent throughout a scene.
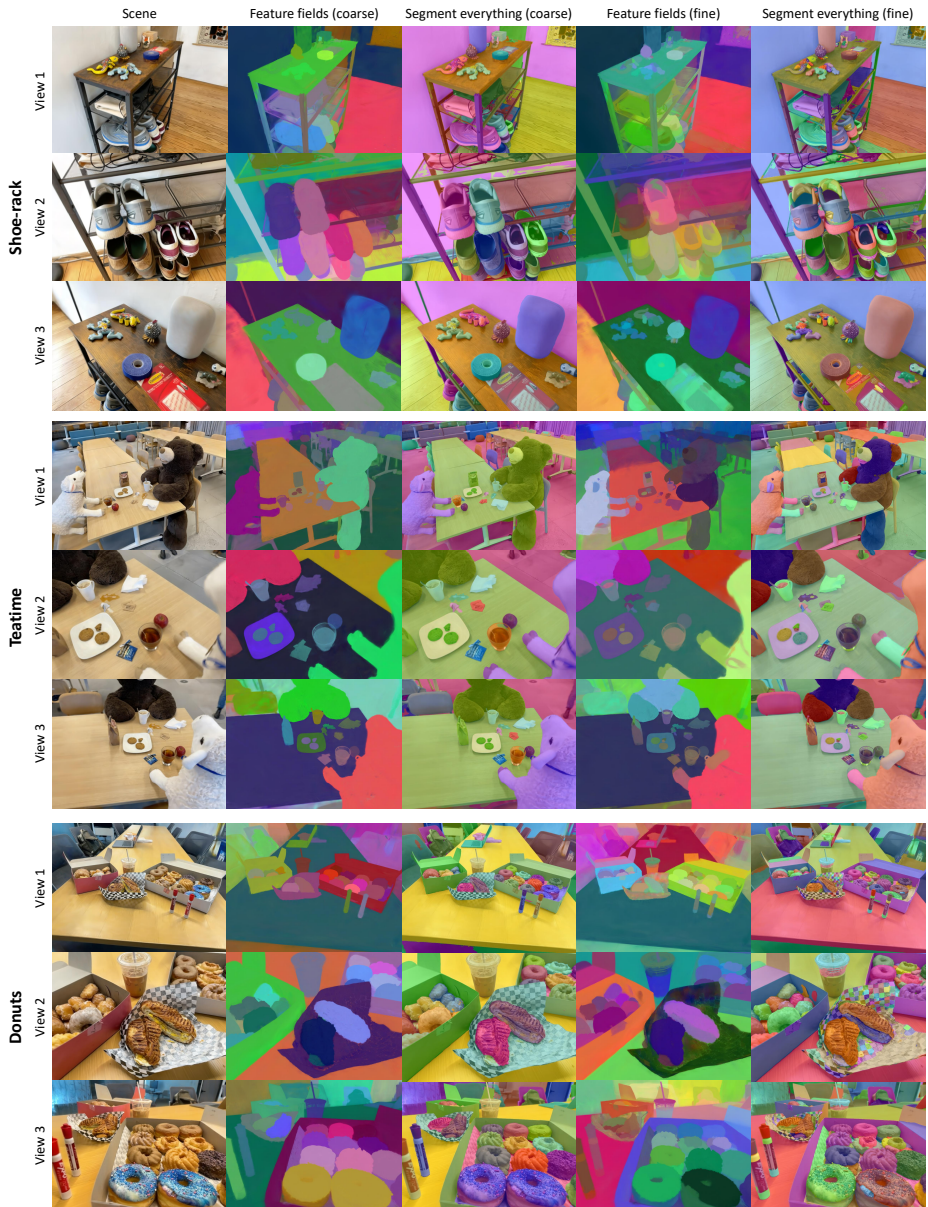
**Fig. 16:** Segmentation of everything results on the LeRF Dataset. We present automatic segmentation results (third and fifth columns) along with PCA visualizations of rendered feature maps (second and fourth columns) at two granularity levels for Shoe-rack, Teatime, and Donuts scenes (first column) from the LeRF Dataset. Objects classified with the same ID in the segmentation results share the same overlaid color across the three given views, as each global cluster ID remains consistent throughout a scene.
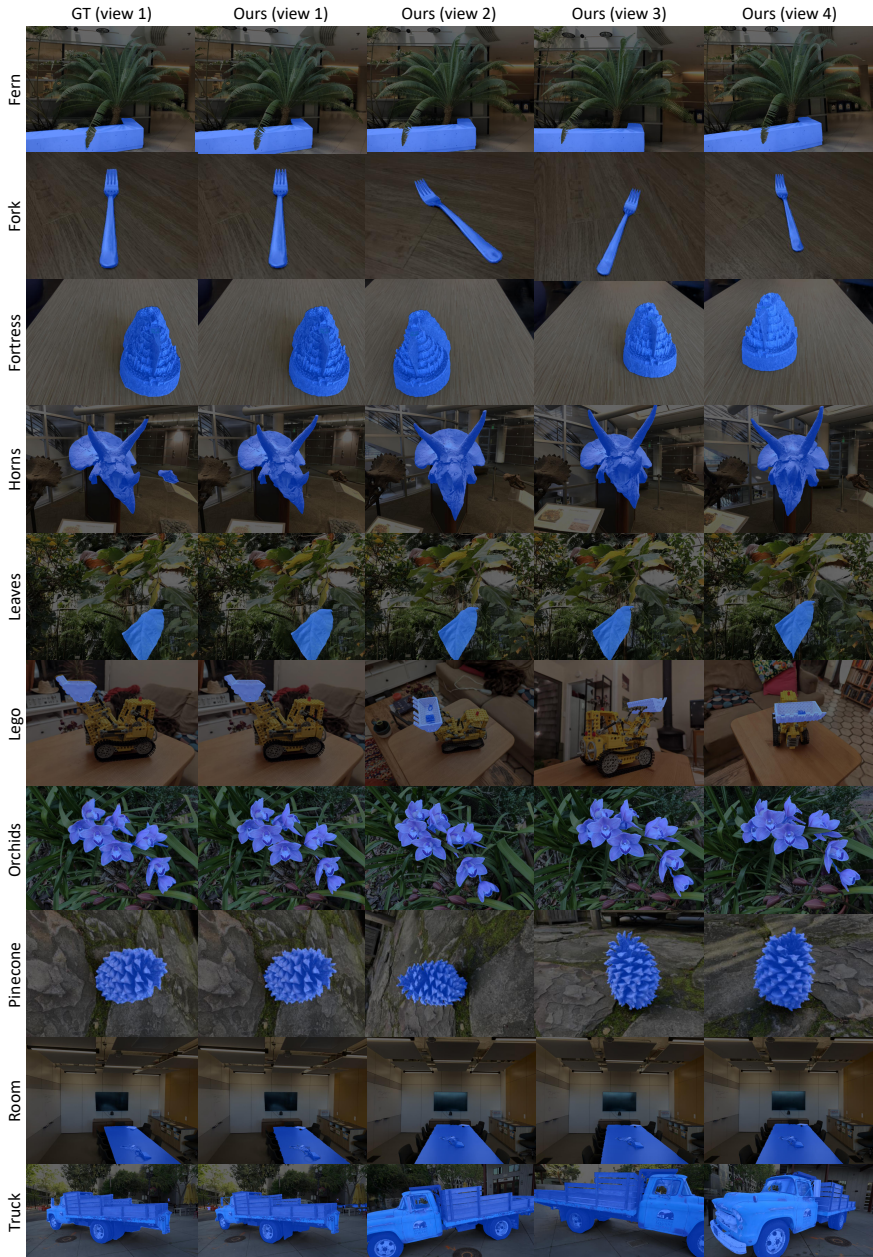
**Fig. 17:** 3D segmentation results on the SPIn-NeRF Dataset. We use the label propagation method based on the ground truth mask of a reference view (first column) to identify the cluster IDs belonging to the target object. These IDs are then used to generate 2D masks for test views (subsequent columns).