

LangSurf: Language-Embedded Surface Gaussians for 3D Scene Understanding

Hao Li^{1*} Roy Qin^{2*†} Zhengyu Zou^{1*} Diqi He¹ Bohan Li³ Bingquan Dai² Dingwen Zhang^{1†}
Junwei Han¹

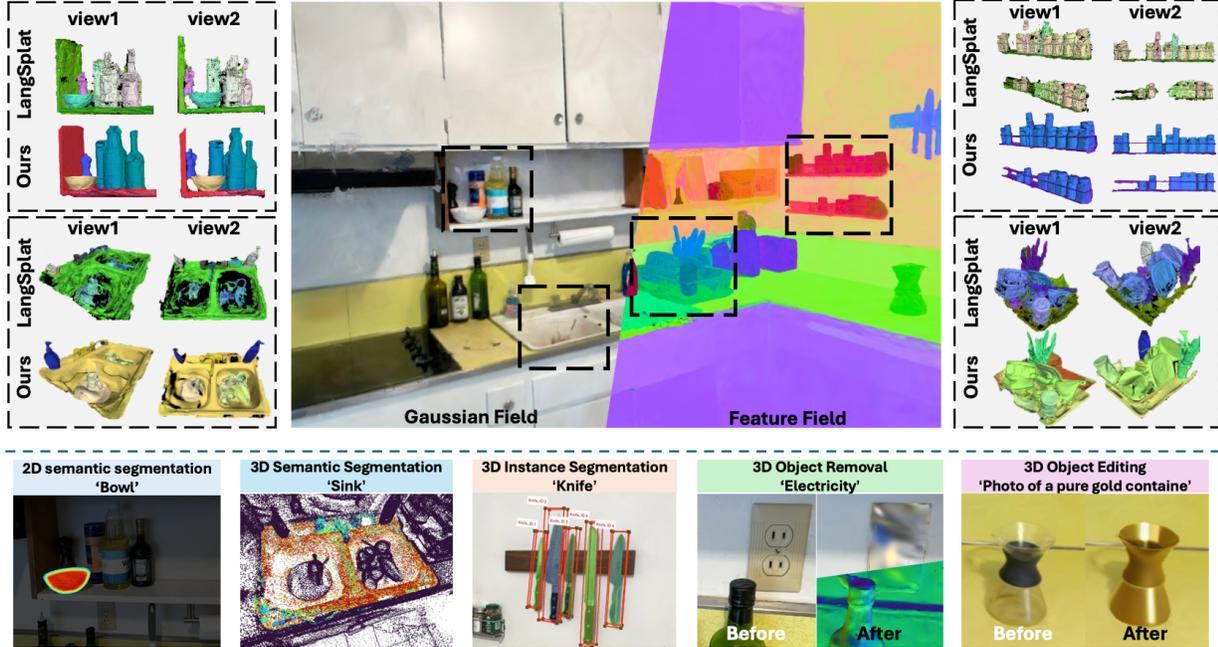


Figure 1: We proposed LangSurf, a model that aligns language features with object surfaces to enhance 3D scene understanding. The top left and right panels illustrate the qualitative differences between LangSurf and LangSplat (Qin et al., 2024a) in reconstructing the 3D language field from multiple viewpoints, demonstrating LangSurf’s superior alignment of semantic features with object surfaces. The bottom row shows a variety of downstream applications enabled by LangSurf.

Abstract

Applying Gaussian Splatting to perception tasks for 3D scene understanding is becoming increasingly popular. Most existing works primarily focus on rendering 2D feature maps from novel viewpoints, which leads to an imprecise 3D language field with outlier languages, ultimately failing to align objects in 3D space. To this end, we propose a Language-Embedded Surface Field (LangSurf), which accurately aligns the 3D language fields with the surface of objects, facilitating precise 2D and 3D segmentation with text

query, widely expanding the downstream tasks such as removal and editing. The core of LangSurf is a joint training strategy that flattens the language Gaussian on the object surfaces using geometry supervision and contrastive losses to assign accurate language features to the Gaussians of objects. In addition, we also introduce the Hierarchical-Context Awareness Module to extract features at the image level for contextual information then perform hierarchical mask pooling using masks segmented by SAM to obtain fine-grained language features in different hierarchies. Extensive experiments on open-vocabulary 2D and 3D semantic segmentation demonstrate that LangSurf outperforms the previous state-of-the-art method LangSplat by a large margin. As shown in Fig. 1, our method is capable of segmenting objects in 3D space, thus boosting the effectiveness of our approach in instance recogni-

*Equal contribution ¹Northwestern Polytechnical University
²Tsinghua University ³Shanghai Jiao Tong University. Correspondence to: Roy Qin <royqin.research@gmail.com>, Dingwen Zhang <zhangdingwen2006yy@gmail.com>.

tion, removal, and editing, which is also supported by comprehensive experiments. [Project Page](#).

1. Introduction

Recently, 3D scene understanding has emerged as a critical research focus. By integrating natural language with 3D scenes, systems can enable more intuitive human-computer interactions in applications such as virtual reality (Li et al., 2022b; Jiang et al., 2024; Jaritz et al., 2019; Li et al., 2024), autonomous driving (Caesar et al., 2020; Zou et al., 2025), and robotics (Awais et al., 2023; Kwon et al., 2023; Kerr et al., 2024; Zhang et al., 2025; Firoozi et al., 2023; Lu et al., 2025). However, accurately embedding semantic information within 3D space remains a significant challenge.

Current methods use NeRF (Neural Radiance Fields) (Mildenhall et al., 2021) or 3DGS (3D Gaussian splatting) (Kerbl et al., 2023) as the 3D representation, combined with language features from the CLIP (Radford et al., 2021) model, to enable open-vocabulary 3D querying. However, while the semantic maps generated in these methods are critical for supervising the 3D semantic field, they lack sufficient contextual information (Zhao et al., 2017). More specifically, these methods typically rely on sliding windows (as in LERF (Kerr et al., 2023)) or Segment Anything Model (SAM (Kirillov et al., 2023)) masks (as in LangSplat (Qin et al., 2024a)) to divide images into parts, which are then processed by CLIP to extract the corresponding semantic features. In such a process, the obtained semantic features only contain information from the local image regions, which can hardly be used to represent the semantics of low-texture regions (Vaswani, 2017), such as walls and floors, or complex object structures that are divided into multiple parts.

Additionally, LERF and LangSplat primarily focus on rendering 2D feature maps from novel viewpoints, without imposing constraints to ensure that semantic features are accurately aligned with the true surfaces of objects. Consequently, the extracted semantic features are not spatially consistent with object surfaces in 3D space, which dramatically limits the application performance in downstream tasks like 3D querying, segmentation, and editing.

To address the aforementioned limitations, we propose LangSurf (Language-Embedded Surface Field). Unlike previous methods, LangSurf prioritizes the alignment of semantic features with the actual surfaces of objects in the 3D scene, ensuring a more spatially coherent semantic field. Specifically, to overcome the representation limitation of the local semantic features, our approach introduces a Hierarchical-Context Awareness Module, which first extracts pixel-level semantic features for the entire image and then applies SAM’s mask to perform mask pooling within the corresponding regions, yielding context-aware semantic feature for each

mask. Such a module enriches each mask’s semantic feature by supplementing it with global context information, especially beneficial for low-texture areas or objects with intricate structures. Additionally, by retaining LangSplat’s hierarchical structure, LangSurf enables the model to perceive objects at varying levels of granularity. Through this approach, LangSurf creates a more accurate and contextually aligned 3D semantic field, enabling more effective downstream applications.

LangSurf’s model architecture employs a joint learning strategy that synchronizes geometry and semantic information. We enhance the geometric quality of the semantic field through multi-view normal vector constraints, ensuring precise alignment with object surfaces. Additionally, we employ a self-supervised semantic grouping strategy that assigns language features to Gaussian points in both the 2D feature maps and 3D representations, rather than relying solely on 2D feature map supervision as in previous methods. To differentiate between objects while preserving their unique language characteristics, we implement an instance-aware training scheme that maximizes semantic distances between objects. Together, these strategies ensure precise semantic field distribution within 3D space, improving downstream task performance. Our main contributions are as follows:

1. We introduce LangSurf, a model that emphasizes aligning semantic features with the actual surfaces of objects within 3D scenes. This alignment ensures a more spatially coherent semantic field, improving the accuracy of downstream tasks such as 3D querying, segmentation, and editing.
2. We develop a Hierarchical-Context Awareness Module that extracts pixel-level semantic features from entire images and applies mask pooling using SAM’s masks. This process enriches each mask’s semantic feature with global context information, particularly benefiting low-texture areas and objects with intricate structures.
3. The proposed method is evaluated on the LERF and ScanNet datasets, demonstrating superior performance in open-vocabulary 2D/3D semantic segmentation tasks compared to current SOTA methods based on 2D images and 3D language fields. Additionally, the model showcases potential in 3D editing / removal applications, highlighting its versatility and effectiveness.

2. Related Works

2.1. 3D Gaussian Models

3D Gaussian Splatting (3DGS) (Kerbl et al., 2023) has emerged as a significant advancement in 3D reconstruction, offering high-resolution real-time rendering capabilities that surpass traditional Neural Radiance Field (NeRF) methods

(Mildenhall et al., 2021; Pumarola et al., 2021; Qin et al., 2024b; Park et al., 2021; Barron et al., 2021; 2023; Xu et al., 2022). This efficiency facilitates various downstream applications (Zhang et al., 2024b; Cai et al., 2024b; Liu et al., 2024b; Qin et al., 2024a; Wu et al., 2024a; Cai et al., 2024a; Zhang et al., 2024a; Xu et al., 2024; Ren et al., 2024; Yu et al., 2024; Keetha et al., 2024; Liu et al., 2024a). For dynamic scenes, 4D-GS (Wu et al., 2024a) has been adapted to handle temporal changes efficiently, enabling real-time rendering of deformable scenes. In the field of 3D generation (Liang et al., 2024; Cai et al., 2024c; He et al., 2024), Lucidreamer (Liang et al., 2024) incorporates 3DGS into text-to-3D generation pipeline to achieve high-quality 3D generated results. Some methods (Qin et al., 2024a; Qiu et al., 2024; Kim et al., 2024; Shi et al., 2024; Ji et al., 2024; Zuo et al., 2024; Yue et al., 2025; Wu et al., 2024c) focus on reconstructing 3D feature field, LangSplat (Qin et al., 2024a) utilizes language-embedded Gaussians to enable precise and efficient open-vocabulary querying within 3D spaces. Despite its capabilities, 3DGS struggles with accurately reconstructing 3D object surfaces. Several subsequent methods (Guédon & Lepetit, 2024; Huang et al., 2024; Chen et al., 2024) enhance this process by incorporating geometric regularization. Sugar (Guédon & Lepetit, 2024) introduces a regularization term that encourages Gaussians to align closely with the surface of the scene. This alignment is then utilized to extract a mesh from the Gaussians through Poisson reconstruction. PGSR (Chen et al., 2024) employs an unbiased depth rendering method and supervises Gaussian primitives using both single-view and multi-view loss, leading to detailed reconstruction of 3D surfaces and meshes. Unlike these methods, our paper focuses on the surface of 3D language field, which is critical for open-vocabulary 2D / 3D segmentation, 3D removal, and editing.

2.2. 3D Scene Understanding

The integration of natural language processing with 3D scene understanding has garnered significant attention in recent years. Recent advancements in open-world 3D scene understanding have sought to combine Neural-based 3D representations (Kerr et al., 2023; Kim et al., 2024; Qin et al., 2024a; Ye et al., 2025; Shi et al., 2024) with SAM or CLIP. LangSplat (Qin et al., 2024a) advances the field by utilizing a collection of 3D Gaussians and training a scene-wise language autoencoder, each encoding language feature distilled from SAM and CLIP, to represent the scene-specific latent language field. Gaussian Grouping (Ye et al., 2025) enhances each Gaussian with a compact identity encoding by leveraging the mask predictions by SAM and spatial consistency regularization, enabling grouping based on object instance or category within the 3D scene. However, current methods provide limited surface reconstruction, leading

to inaccurate meshes. In contrast, we propose a language-embedded surface field that embeds 3D language fields with the surface to enhance surface reconstruction quality, boosting the accuracy of mesh segmentation and other geometric tasks.

3. Preliminaries

LangSplat (Qin et al., 2024a) integrate semantic into 3DGS for 3D scene understanding. It additionally adds latent space $\mathbf{f}^{lang} \in \mathbb{R}^3$ into Gaussian attributes to represent the language field in 3D. Such a method facilitates real-time alpha blending of numerous Gaussians to render novel-view RGB images and language maps:

$$\begin{cases} \mathbf{C}(v) = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \\ \mathbf{F}(v) = \sum_{i \in \mathcal{N}} \mathbf{f}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \end{cases} \quad (1)$$

where $\mathbf{C}(v)$ and $\mathbf{F}(v)$ represent the rendered color and feature at pixel v , \mathcal{N} is the number of Gaussians that the ray passes through, α is the opacity of the Gaussian, and $\mathbf{c}_i \in \mathbb{R}^3$ is the view-dependent colors represented as a series of sphere harmonics coefficients in the practice of 3DGS. Although it can synthesize 2D feature maps from novel viewpoints, it lacks constraints to ensure that semantic features are accurately aligned with the true surfaces of objects, causing inaccurate language representation at the 3D level.

4. Methodology

Given input views $\{\mathbf{I}_i \in \mathbb{R}^{3 \times H \times W}\}$, our main objective is to reconstruct the language-embedded surface field, which is denoted as a set of Gaussians $\{(\mathbf{x}_i, \boldsymbol{\Sigma}_i, \mathbf{S}_i, \alpha_i, \mathbf{f}_i^{lang}, \mathbf{f}_i^{ins})\}$, where $\mathbf{f}_i^{ins} \in \mathbb{R}^3$ denotes instance features. Such a field allows performing text queries and manipulations (*i.e.* removal, editing) at the instance level. Our framework can be divided into two stages. **Hierarchical-Context Awareness Module** extracts language-pixel aligned features $\{\mathbf{L}_i^{lang, h}, |h = (s, m, l)\}$ on different hierarchies from image sets to facilitate the subsequent training (see Sec. 4.1). **Language-Embedded Surface Field Training** constructs the language-embedded surface field by utilizing a multi-step training strategy with 2D and 3D semantic supervision (see Sec. 4.2).

4.1. Hierarchical-Context Awareness Module

In order to correctly capture visual-language aligned features, we propose a simple but efficient method named Hierarchical-Context Awareness Module. Unlike previous work (Qin et al., 2024a) extracted language features for the masked objects segmented by SAM, for each input image \mathbf{I}_i ,

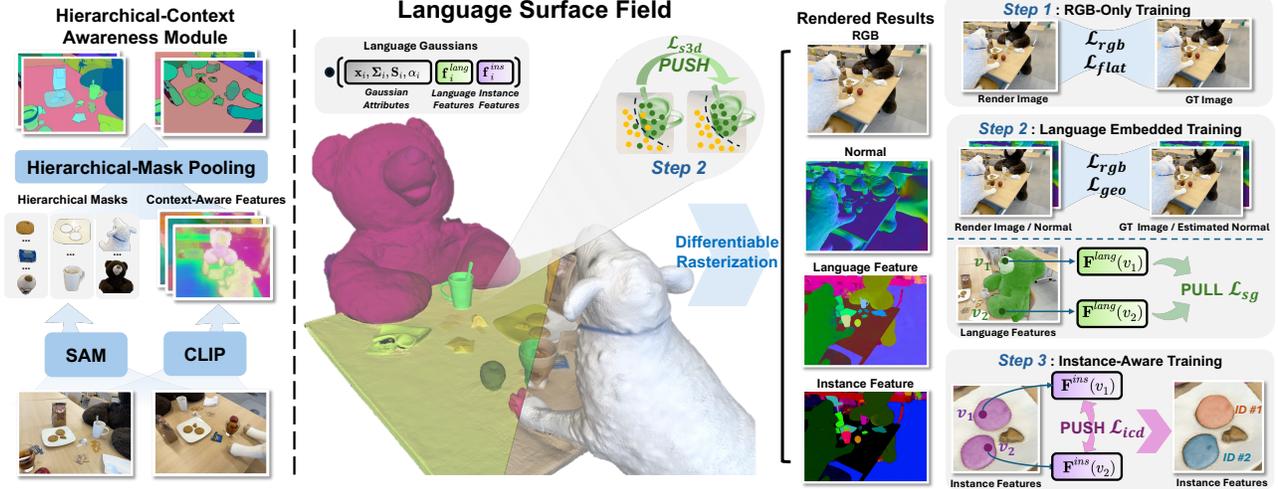


Figure 2: **Overview of proposed LangSurf.** Given input views, we reconstruct a language-embedded surface field to enable 2D / 3D open-vocabulary segmentation as well as downstream tasks. Our pipeline contains two main steps: 1) Hierarchical-Context Awareness Module extracts context-aware features with multiple hierarchies (Sec. 4.1); 2) Language-Embedded Training utilizes a joint training strategy to construct language-embedded surface field (Sec. 4.2).

we apply a pre-trained image encoder (Ghiasi et al., 2022) to produce image-wise features $\mathbf{L}_i^{lang} \in \mathbb{R}^{D \times H \times W}$. It considers contextual information to enable accurate visual-language alignment of the features, especially for obscured objects and objects with low textures and shapes.

However, simply adopting \mathbf{L}_i^{lang} to train the Gaussian model will limit the ability of open-vocabulary query with different scales since context-aware features suppress the diversity of the features, making it difficult to distinguish small objects from the large one, such as “bear nose” and “bear”. To this end, we perform Hierarchical-Mask Pooling for the image features, which uses the multi-hierarchy masks $\{\mathbf{M}_i^h, |h = s, m, l\}$ segmented by SAM to decompose the multi-scale language features from the original features \mathbf{L}_i , where s, m, l represents the small, medium, and large hierarchy and each hierarchy obtains multi masks $\mathbf{M}_i^h = \{\mathbf{M}_{i,j}^h, (j = 1, \dots, M)\}$. For each hierarchy, image features perform masked average pooling to enhance semantic consistency within the masks:

$$\mathbf{L}^{lang,h}(v) = \frac{\sum \mathbf{L}^{lang}(v) \cdot \mathbf{M}^h(v)}{\sum \mathbf{M}^h(v)}, h = \{s, w, l\}, \quad (2)$$

where v represents the pixel within the mask region at the hierarchy h .

In the end, we adopt an end-to-end autoencoder that compresses the language features $\{\mathbf{H}^{lang,h}\}$ into low-dimensional latent space $\{\mathbf{L}^{lang,h}\}$ during training and decodes the latent space into original features during inference, where $\{\mathbf{H}_i^{lang,h}\}$ is a 3-channels feature map. It reduces memory consumption and improves efficiency.

4.2. Language-Embedded Surface Field Training

We decompose the training procedure into three steps: 1) basic RGB supervision is deployed to obtain basic 3D representation; 2) both geometry and semantic supervision are deployed to optimize the language Gaussians with not only accurate spatial semantic distributions; 3) well-trained language features initialize the instance features of Gaussians, then instance-level training are deployed to distinguish objects from the language space. These stages are as follows:

Step 1: RGB-Only Training. In this step, we aim to obtain the basic 3D field by deploying basic RGB supervision \mathcal{L}_{rgb} followed by a Gaussian flatten supervision \mathcal{L}_{flat} :

$$\begin{cases} \mathcal{L}_{rgb} = \|\mathbf{C}_i - \mathbf{I}_i\|_1, \\ \mathcal{L}_{flat} = \|\min(s_1, s_2, s_3)\|_1, \end{cases} \quad (3)$$

where $\{s_1, s_2, s_3\} = \text{diag}(\mathbf{S}_i)$ are the scale factors of Gaussian. It compresses the Gaussians and makes them flatten into the object planes.

Step 2: Language-Embedded Training. Here, we perform joint training using geometry and semantic supervision to construct a language-embedded surface field. Firstly, following the strategy of PGSR (Chen et al., 2024), we adopt geometry regularization constraints \mathcal{L}_{geo} to optimize the geometry representation of our model. Apart from regular L2 semantic loss between \mathbf{F}^{lang} and \mathbf{H}^{lang} , we further propose Semantic Grouping terms \mathcal{L}_{sg} , which groups the rendered features $\mathbf{F}^{lang}(\cdot)$ within the same mask \mathbf{M}_j of the image by minimizing their semantic distance. Such a strategy maintains semantic consistency within the object and creates clearer boundaries between different objects, the

Table 1: **2D Quantitative Results on LERF Dataset.** We report the open-vocabulary localization accuracy (%) and 2D semantic segmentation (IoU scores). LSeg (Li et al., 2022a) and CAT-Seg (Cho et al., 2024) are 2D open-vocabulary segmentation networks, while other methods (Qin et al., 2024a; Kerr et al., 2023; Ye et al., 2025) are language field models. We denote ‘‘GS-Group’’ as Gaussian-Grouping. The **bold** denotes the best results.

Scene	LSeg		CAT-Seg		LERF		LangSplat		GS-Group		Ours	
	mAcc \uparrow	mIou \uparrow										
Teatime	28.07	15.37	69.49	34.39	71.69	38.76	88.10	65.10	79.60	58.20	84.75	73.57
Ramen	8.45	2.16	53.52	16.77	54.71	21.54	56.34	46.52	30.90	24.30	63.40	47.03
Kitchen	44.00	20.87	68.00	28.24	64.15	29.19	72.73	50.83	54.50	39.40	81.82	54.99
Overall	26.84	12.80	63.67	26.46	63.51	29.83	72.39	54.15	55.00	40.63	76.65	58.53

formulation is shown below:

$$\mathcal{L}_{sg} = \frac{1}{M} \sum_{j=1}^M \sum_{v_1, v_2 \in \mathbf{M}_j} \|\mathbf{F}^{lang}(v_1) - \mathbf{F}^{lang}(v_2)\|_2, \quad (4)$$

where M is the total numbers of segmentation masks $\mathbb{M} = \{\mathbf{M}_j, (j = 1, \dots, M)\}$ of the image. Meanwhile, given the fact that the outlier language Gaussians seriously disturb the 3D localization of the objects as well as downstream tasks, we additionally propose a Spatial-Aware Semantic Supervision \mathcal{L}_{s3d} to increase the spatial constraint of language Gaussians to suppress the outlier of language Gaussians. It utilizes KL-divergence supervision to align the semantic features with the top-k nearest Gaussians:

$$\mathcal{L}_{s3d} = \sum_{j=1}^N \sum_{k=1}^{N_k} \mathbf{f}_j^{lang} \cdot \left(\log \left(\mathbf{f}_j^{lang} / \mathbf{f}_k^{lang} \right) \right), \quad (5)$$

where N denotes the total number of Gaussian points and N_k is the Top-K nearest number of the i -th Gaussian point. Such strategies enable our language-embedded surface field to fit into the surface of the scene accurately, achieving comprehensive scene understanding.

Step 3: Instance-Aware Training. We expand our language-embedded surface field for perceiving instance level since there are multiple objects with the same categories in the scenes. To this end, we introduce new instance-aware parameters to assign each Gaussian to its corresponding instance or stuff in the 3D scene while preserving all attributes of the language Gaussians for text-aligned querying. In detail, we use well-trained language features \mathbf{f}_i^{lang} to initialize the instance features \mathbf{f}_i^{ins} for each Gaussian. After that, we compute the mean feature \mathbf{z}_i^{ins} for each masked region \mathbf{M}_i on the render instance map \mathbf{F}_i^{ins} :

$$\mathbf{z}_i^{ins} = \frac{1}{|\mathbf{M}_i|} \sum_{v \in \mathbf{M}_i} \mathbf{F}_i^{ins}(v), \quad (6)$$

where $|\cdot|$ represents the area of the mask. In the end, an Instance Contrastive Decomposition supervision \mathcal{L}_{icd} is proposed to decompose the objects by maximizing the distance

Table 2: **2D Quantitative Results on ScanNet Dataset.** We report the open-vocabulary 2D semantic segmentation (IoU scores). The **bold** denotes the best results.

Scene	LSeg	CATSeg	LangSplat	GS-Group	Ours
085	32.75	42.36	36.69	43.41	66.89
114	19.11	33.66	16.81	20.78	33.25
616	33.83	45.57	24.29	23.84	51.40
617	15.53	24.51	11.49	19.42	40.09
Mean	25.30	36.52	22.32	26.86	47.91

Table 3: **3D Quantitative Results on ScanNet Dataset.** We report the average open-vocabulary Semantic F-Score \uparrow . The **bold** denotes the best results.

Scene	LangSplat	GS-Group	Ours
085	6.19	18.92	52.70
114	13.19	13.35	30.62
616	08.12	09.47	36.44
617	11.39	10.62	33.05
Mean	9.72	13.09	38.20

between instance features of different masks:

$$\mathcal{L}_{icd} = \sum_{j=1}^M \sum_{k \neq j}^M \text{ReLU}(D_{min} - \|\mathbf{z}_j^{ins} - \mathbf{z}_k^{ins}\|_2), \quad (7)$$

where D_{min} is the minimum distance between instances. Notably, we only train the instance features of the Gaussians instead of all parameters.

5. Experiments

5.1. Implementation Details

To extract image features of each image, we utilize the backbone of OpenSeg (Ghiasi et al., 2022) as the extractor with a vallina CLIP for the text encoder. For SAM (Kirillov et al., 2023), we use the ViT-H model to segment 2D masks. During Language-Embedded Surface Field training, we train 7,000 iterations for Step 1, 23,000 iterations for Step 2, and 10,000 for Step 3. To evaluate the comprehensive scene

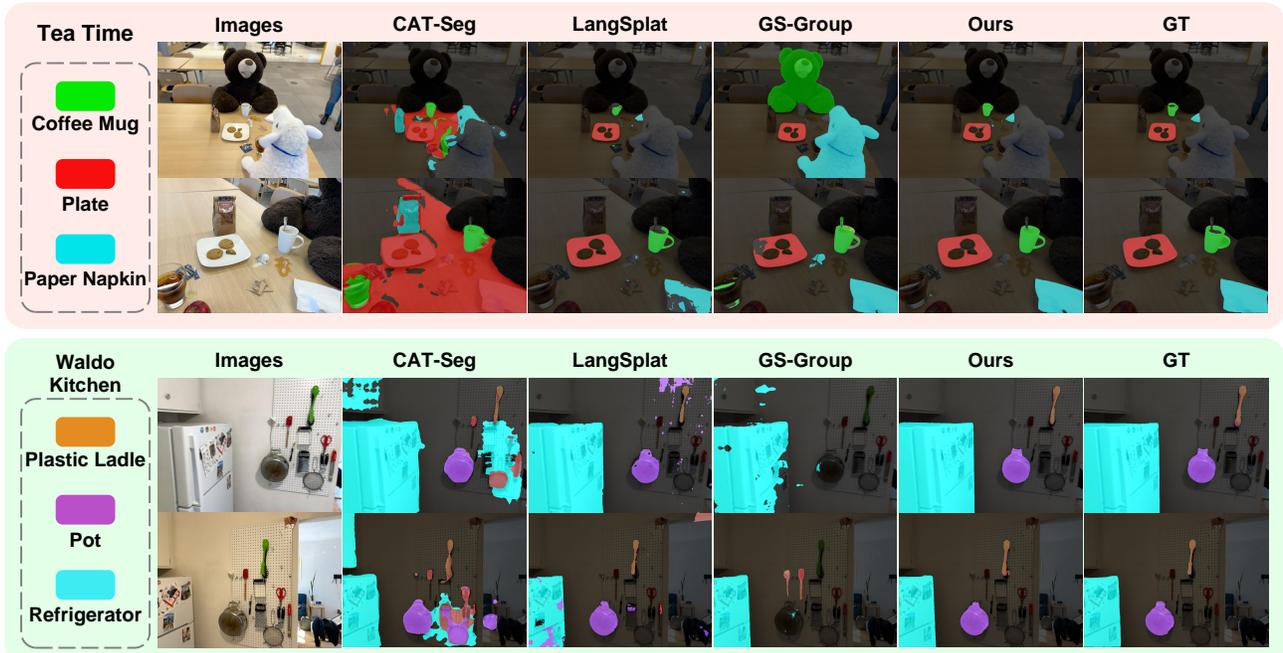


Figure 3: **2D Qualitative Results on LERF Datasets.** Here we showcase two scenes (*i.e.* Teatime, Waldo Kitchen) with multiple text-query segmentation masks. On the left, we present the images alongside the queried texts. On the right, we display the rendered results of our method and other methods, along with the corresponding ground truth annotations.

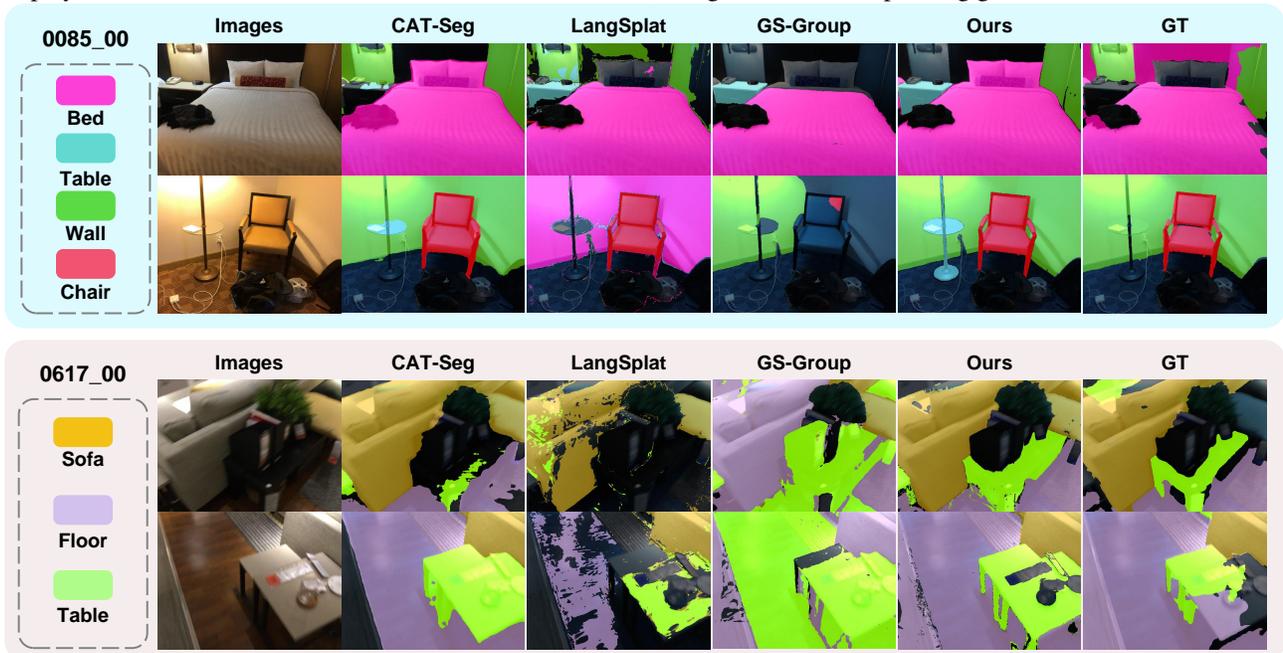


Figure 4: **2D Qualitative Results on ScanNet (Dai et al., 2017) Dataset.** Here we showcase two scenes (*i.e.* scene0085_00, scene0617_00) with multiple text-query segmentation masks. The masks predicted by ours contain more comprehensive regions and sharper boundaries than other methods, such as the “Table” prompt, which also surpasses the GT masks.

understanding of our proposed language-embedded surface field, we employ two datasets, LERF (Kerr et al., 2023) and ScanNet (Dai et al., 2017). The LERF dataset is an in-the-wild dataset captured by a handheld device, which

consists of a greater number of object elements. It also consists of precise annotations on both segmentation masks and bounding boxes at 2D level. The ScanNet dataset is a large dataset captured by RGB-D devices in complex indoor

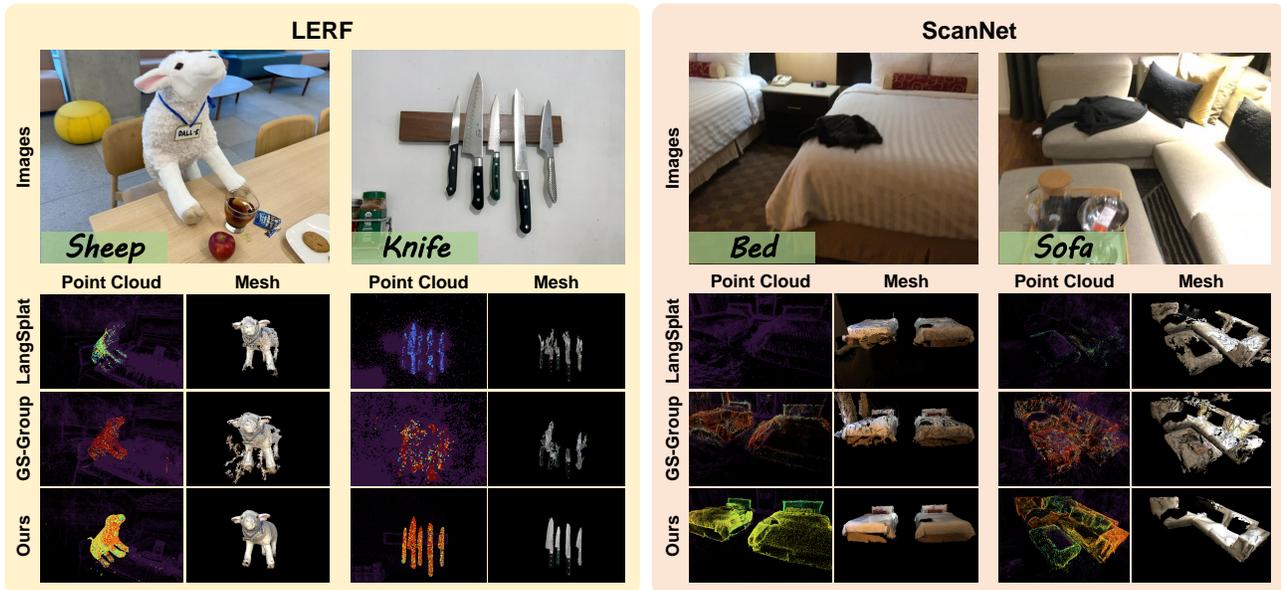


Figure 5: **3D Qualitative Results on both LERF (Kerr et al., 2023) and ScanNet (Dai et al., 2017) Datasets.** We compare our method with other GS-based methods (Qin et al., 2024a; Ye et al., 2025). We show the queried point clouds with activated score and mesh corresponding to the given text.

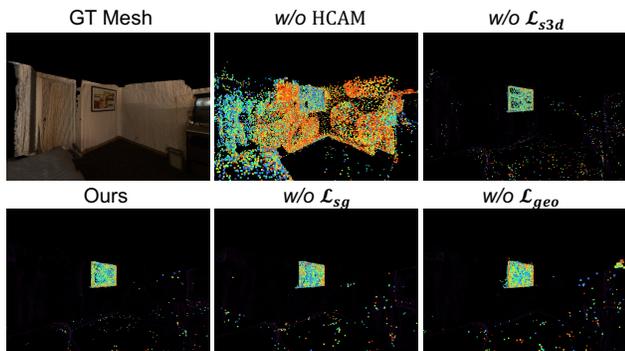


Figure 6: **Qualitative Results of Ablations on ScanNet (Dai et al., 2017) Dataset.** We visualize the 3D open-vocabulary segmentation with “pictures” prompt.

scenes. Each scene includes precise poses and semantic labels at 3D level. These characteristics make it suitable for 3D scene understanding tasks.

5.2. Results of Open-Vocabulary Query Tasks

We conduct experiments on 2D and 3D open-vocabulary segmentation by querying the language feature of each Gaussian primitive with text. Specifically, we take a set of text queries and calculate the cosine similarity between the text embeddings and the Gaussian language features, which yields a confidence score. A threshold is applied to filter out Gaussians with low confidence. We adopt mIoU and mAcc for 2D evaluation, Semantic F-Score for 3D evaluation.

5.2.1. 2D OPEN-VOCABULARY SEMANTIC SEGMENTATION

We evaluate our method on the LERF and ScanNet datasets. Results are reported in Tab. 1 and Tab. 1. By comparing with existing state-of-the-art open-vocabulary 2D segmentation methods, including LSeg and CAT-Seg, as well as 3D language field techniques like LERF, LangSplat, and Gaussian Grouping (GS-Group), our method achieves superior performance in both segmentation and localization tasks, i.e., a 4.38% improvement in terms of mIoU and a 4.26% increase in terms of mAcc on the LERF dataset. On the ScanNet dataset, the improvement upon the best existing method comes to 11.39% in terms of mIoU. The visualization of the 2D segmentation masks is shown in Fig. 3 and Fig. 4. As can be seen, our approach excels at segmenting objects with fine-grained boundaries. This demonstrates the capability of our language-embedded surface field to capture contextual and spatial information to interpret complex scenes.

5.2.2. 3D OPEN-VOCABULARY SEMANTIC SEGMENTATION

The comparison of the 3D open-vocabulary segmentation experiment is conducted on the ScanNet dataset. The results are reported in Tab. 3. As can be seen, our method outperforms the existing methods by a large margin, i.e., 25.11% in the Semantic F-Score metric. The visualization results are shown in Fig. 5, which reflect that our method achieves more accurate text-query-based segmentation in 3D space with fewer outlier activations and is able to gen-

Table 4: **Ablations of proposed module and losses.** We perform ablations on both ScanNet (scene0085) (Dai et al., 2017) and LERF (Teaime) (Kerr et al., 2023). Metric ‘‘Semantic F1-Score’’ shorts for ‘‘S.F.’’.

HCAM	\mathcal{L}_{geo}	\mathcal{L}_{sg}	\mathcal{L}_{s3d}	ScanNet		LERF
				mIoU	S. F.	mIoU
✗	✓	✓	✓	38.67	30.55	64.97
✓	✗	✓	✓	63.88	49.59	70.37
✓	✓	✗	✓	64.57	50.89	72.21
✓	✓	✓	✗	65.47	50.13	73.27
✓	✓	✓	✓	66.89	52.70	73.57

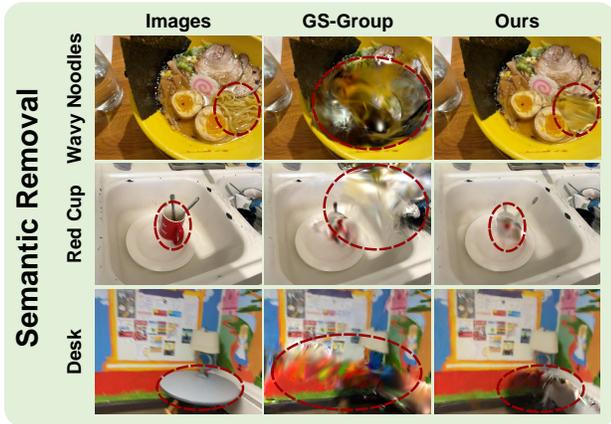


Figure 7: **Qualitative Results of Semantic-Level 3D Objects Removal.** We present the semantic-level removal results on the LERF and ScanNet datasets, comparing them with Gaussian Grouping. The removed area is highlighted with ‘‘○’’.

erate more detailed meshes for the queried objects. The above experiment reflects that our method can effectively align language-embedded surface fields with the surface of objects, obtaining accurate 3D perception capacity.

5.3. Ablation Study

We conduct ablation experiments on both ScanNet and LERF datasets. The results are reported in Tab. 4. As can be seen, our Hierarchical-Context Awareness Module (HCAM) makes significant improvements on both datasets by 28.22% and 8.6% in mIoU. Its effectiveness comes from capturing contextual information for language-aligned feature extraction, which promotes the accurate segmentation of low-texture regions such as walls and pictures on ScanNet, as shown in Fig. 6. Regarding \mathcal{L}_{geo} , it helps reduce Gaussian artifacts and align the language Gaussians closer to the surface of objects, achieving accurate 3D language representation and improving the performance of open-vocabulary query (3.01% and 3.2% in mIoU and 3.11% in S.F.). Meanwhile, \mathcal{L}_{sg} promotes semantic consistency within the same objects and creates sharper boundaries, which brings more

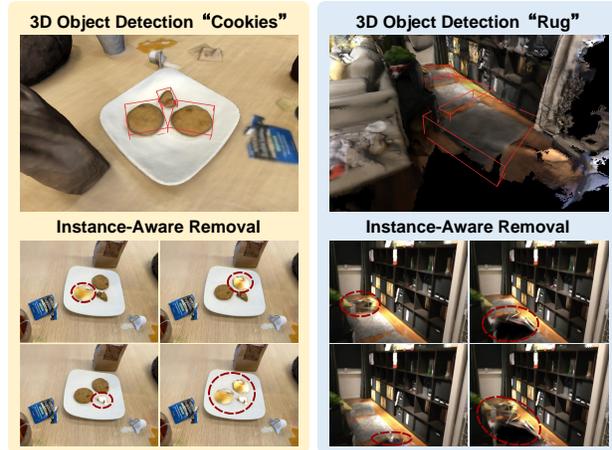


Figure 8: **Qualitative Results of Instance-Level 3D Object Detection and Removal.** Objects are highlighted with ‘‘○’’.

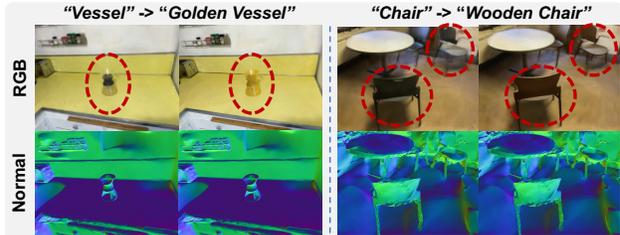


Figure 9: **Qualitative Results of 3D Object Editing** on ScanNet and LERF Datasets. Objects are highlighted with ‘‘○’’.

accurate semantic segmentation, allowing mIoU to improve 2.32% and 1.36% on both datasets. In the end, by introducing spatial supervision to refine the Gaussians that are inconsistent with the semantic distribution of nearby Gaussians, \mathcal{L}_{s3d} effectively reduces the number of outliers during open-vocabulary querying. This resulted in a 2.57% improvement of the S.F. metric.

5.4. Applications

We apply our method to two representative downstream tasks (*i.e.* Object Removal and Editing) to validate the performance of our language-embedded surface field.

3D Scene Objects Removal & Adding. 3D object removal involves completely deleting an object from a 3D scene, serving as an evaluation for the 3D localization accuracy of the proposed method. We employ text-query-based 3D object removal tasks at both semantic and instance levels. As shown in Fig. 7 and Fig. 8, our method effectively removes the text-queried items without affecting nearby scenes compared with Gaussian Grouping. Meanwhile, Fig. 10 demonstrates adding queried objects (*i.e.* cookie bag) from ‘‘Teatime’’ scene to ‘‘Kitchen’’ scene. These experiments reflect our method achieves superior language representation in 3D.

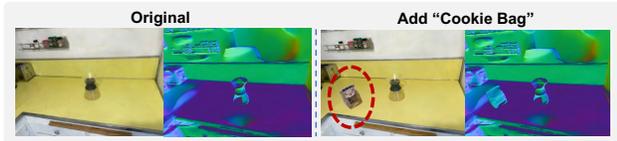


Figure 10: **Qualitative Results of 3D Object Adding on LERF Datasets.** Objects are highlighted with ‘ \circ ’.

3D Scene Objects Editing. Object Editing aims to edit the Gaussians of the queried object in 3D scene. Fig. 9 shows our method enables precisely editing the Gaussians of the object while minimizing any impact on the background, which reflects our language-embedded surface field has accurate language 3D representation.

6. Conclusions

In this paper, we propose Language-Embedded Surface Field that accurately represents the language field in 3D to align the surface of objects. Unlike previous methods that focus on 2D render tasks, we propose a joint training strategy that uses geometry constraint and semantic contrastive losses to align the language field with the surface of objects to enhance the 3D representation of our model. Furthermore, a Hierarchical-Context Awareness Module is proposed to extract context-aware features in different hierarchies for the Gaussian model training. Comprehensive experiments demonstrate that LangSurf achieves significant performance improvements on both 2D and 3D metrics (sometimes $> 10\%$) compared with existing SOTA methods. Based on this, LangSurf can further support multiple downstream tasks, including object removal and editing, which also achieves significant improvements compared with existing methods.

References

- Awais, M., Naseer, M., Khan, S., Anwer, R. M., Cholakkal, H., Shah, M., Yang, M.-H., and Khan, F. S. Foundational models defining a new era in vision: A survey and outlook. *arXiv preprint arXiv:2307.13721*, 2023.
- Barron, J. T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., and Srinivasan, P. P. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5855–5864, 2021.
- Barron, J. T., Mildenhall, B., Verbin, D., Srinivasan, P. P., and Hedman, P. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19697–19705, 2023.
- Caesar, H., Bankiti, V., Lang, A. H., Vora, S., Liong, V. E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., and Beijbom, O. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.
- Cai, Y., Wang, J., Yuille, A., Zhou, Z., and Wang, A. Structure-aware sparse-view x-ray 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11174–11183, 2024a.
- Cai, Y., Xiao, Z., Liang, Y., Qin, M., Zhang, Y., Yang, X., Liu, Y., and Yuille, A. Hdr-gs: Efficient high dynamic range novel view synthesis at 1000x speed via gaussian splatting. *arXiv preprint arXiv:2405.15125*, 2024b.
- Cai, Y., Zhang, H., Zhang, K., Liang, Y., Ren, M., Luan, F., Liu, Q., Kim, S. Y., Zhang, J., Zhang, Z., et al. Baking gaussian splatting into diffusion denoiser for fast and scalable single-stage image-to-3d generation. *arXiv preprint arXiv:2411.14384*, 2024c.
- Chen, D., Li, H., Ye, W., Wang, Y., Xie, W., Zhai, S., Wang, N., Liu, H., Bao, H., and Zhang, G. Pgsr: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *arXiv preprint arXiv:2406.06521*, 2024.
- Cho, S., Shin, H., Hong, S., Arnab, A., Seo, P. H., and Kim, S. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4113–4123, 2024.
- Dai, A., Chang, A. X., Savva, M., Halber, M., Funkhouser, T., and Nießner, M. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5828–5839, 2017.
- Firoozi, R., Tucker, J., Tian, S., Majumdar, A., Sun, J., Liu, W., Zhu, Y., Song, S., Kapoor, A., Hausman, K., et al. Foundation models in robotics: Applications, challenges, and the future. *The International Journal of Robotics Research*, pp. 02783649241281508, 2023.
- Ghiasi, G., Gu, X., Cui, Y., and Lin, T.-Y. Scaling open-vocabulary image segmentation with image-level labels. In Avidan, S., Brostow, G., Cissé, M., Farinella, G. M., and Hassner, T. (eds.), *Computer Vision – ECCV 2022*, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-20059-5.
- Guédon, A. and Lepetit, V. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5354–5363, 2024.
- He, H., Liang, Y., Wang, L., Cai, Y., Xu, X., Guo, H.-X., Wen, X., and Chen, Y. Lucidfusion: Generating 3d gaussians with arbitrary unposed images. *arXiv preprint arXiv:2410.15636*, 2024.

- Huang, B., Yu, Z., Chen, A., Geiger, A., and Gao, S. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–11, 2024.
- Jaritz, M., Gu, J., and Su, H. Multi-view pointnet for 3d scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 0–0, 2019.
- Ji, M., Qiu, R.-Z., Zou, X., and Wang, X. Graspplats: Efficient manipulation with 3d feature splatting. *arXiv preprint arXiv:2409.02084*, 2024.
- Jiang, Y., Yu, C., Xie, T., Li, X., Feng, Y., Wang, H., Li, M., Lau, H., Gao, F., Yang, Y., et al. Vr-gs: a physical dynamics-aware interactive gaussian splatting system in virtual reality. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–1, 2024.
- Keetha, N., Karhade, J., Jatavallabhula, K. M., Yang, G., Scherer, S., Ramanan, D., and Luiten, J. Splatam: Splat track & map 3d gaussians for dense rgb-d slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21357–21366, 2024.
- Kerbl, B., Kopanas, G., Leimkühler, T., and Drettakis, G. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, July 2023.
- Kerr, J., Kim, C. M., Goldberg, K., Kanazawa, A., and Tancik, M. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19729–19739, 2023.
- Kerr, J., Kim, C. M., Wu, M., Yi, B., Wang, Q., Goldberg, K., and Kanazawa, A. Robot see robot do: Imitating articulated object manipulation with monocular 4d reconstruction. *arXiv preprint arXiv:2409.18121*, 2024.
- Kim, C. M., Wu, M., Kerr, J., Goldberg, K., Tancik, M., and Kanazawa, A. Garfield: Group anything with radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21530–21539, 2024.
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C., Lo, W.-Y., Dollár, P., and Girshick, R. Segment anything. *arXiv:2304.02643*, 2023.
- Kwon, O., Park, J., and Oh, S. Renderable neural radiance map for visual navigation. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9099–9108, 2023. doi: 10.1109/CVPR52729.2023.00878.
- Li, B., Weinberger, K. Q., Belongie, S., Koltun, V., and Ranftl, R. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=RriDjddCLN>.
- Li, C., Li, S., Zhao, Y., Zhu, W., and Lin, Y. Rt-nerf: Real-time on-device neural radiance fields towards immersive ar/vr rendering. In *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design*, pp. 1–9, 2022b.
- Li, H., Zhang, D., Dai, Y., Liu, N., Cheng, L., Li, J., Wang, J., and Han, J. Gp-nerf: Generalized perception nerf for context-aware 3d scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21708–21718, 2024.
- Liang, Y., Yang, X., Lin, J., Li, H., Xu, X., and Chen, Y. Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6517–6526, 2024.
- Liu, K., Zhan, F., Xu, M., Theobalt, C., Shao, L., and Lu, S. Stylegaussian: Instant 3d style transfer with gaussian splatting. In *SIGGRAPH Asia 2024 Technical Communications*, pp. 1–4, 2024a.
- Liu, Y., Huang, X., Qin, M., Lin, Q., and Wang, H. Animatable 3d gaussian: Fast and high-quality reconstruction of multiple human avatars. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 1120–1129, 2024b.
- Lu, G., Zhang, S., Wang, Z., Liu, C., Lu, J., and Tang, Y. Manigaussian: Dynamic gaussian splatting for multi-task robotic manipulation. In *European Conference on Computer Vision*, pp. 349–366. Springer, 2025.
- Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., and Ng, R. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- Park, K., Sinha, U., Hedman, P., Barron, J. T., Bouaziz, S., Goldman, D. B., Martin-Brualla, R., and Seitz, S. M. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.
- Pumarola, A., Corona, E., Pons-Moll, G., and Moreno-Noguer, F. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10318–10327, 2021.

- Qin, M., Li, W., Zhou, J., Wang, H., and Pfister, H. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20051–20060, 2024a.
- Qin, M., Liu, Y., Xu, Y., Zhao, X., Liu, Y., and Wang, H. High-fidelity 3d head avatars reconstruction through spatially-varying expression conditioned neural radiance field. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 4569–4577, 2024b.
- Qiu, R.-Z., Yang, G., Zeng, W., and Wang, X. Feature splatting: Language-driven physics-based scene synthesis and editing. *arXiv preprint arXiv:2404.01223*, 2024.
- 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. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Ren, X., Lu, Y., Liang, H., Wu, Z., Ling, H., Chen, M., Fidler, S., Williams, F., and Huang, J. Scube: Instant large-scale scene reconstruction using voxplats. *arXiv preprint arXiv:2410.20030*, 2024.
- Shi, J.-C., Wang, M., Duan, H.-B., and Guan, S.-H. Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5333–5343, 2024.
- Vaswani, A. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Wu, G., Yi, T., Fang, J., Xie, L., Zhang, X., Wei, W., Liu, W., Tian, Q., and Wang, X. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20310–20320, 2024a.
- Wu, J., Bian, J.-W., Li, X., Wang, G., Reid, I., Torr, P., and Prisacariu, V. GaussCtrl: Multi-View Consistent Text-Driven 3D Gaussian Splatting Editing. *ECCV*, 2024b.
- Wu, Y., Meng, J., Li, H., Wu, C., Shi, Y., Cheng, X., Zhao, C., Feng, H., Ding, E., Wang, J., et al. Opengaussian: Towards point-level 3d gaussian-based open vocabulary understanding. *arXiv preprint arXiv:2406.02058*, 2024c.
- Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., and Neumann, U. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5438–5448, 2022.
- Xu, Y., Chen, B., Li, Z., Zhang, H., Wang, L., Zheng, Z., and Liu, Y. Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1931–1941, 2024.
- Ye, M., Danelljan, M., Yu, F., and Ke, L. Gaussian grouping: Segment and edit anything in 3d scenes. In *European Conference on Computer Vision*, pp. 162–179. Springer, 2025.
- Yu, Z., Chen, A., Huang, B., Sattler, T., and Geiger, A. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19447–19456, 2024.
- Yue, Y., Das, A., Engelmann, F., Tang, S., and Lenssen, J. E. Improving 2d feature representations by 3d-aware fine-tuning. In *European Conference on Computer Vision*, pp. 57–74. Springer, 2025.
- Zhang, C., Zou, Y., Li, Z., Yi, M., and Wang, H. Transplat: Generalizable 3d gaussian splatting from sparse multi-view images with transformers. *arXiv preprint arXiv:2408.13770*, 2024a.
- Zhang, C., Ling, Y., Lu, M., Qin, M., and Wang, H. Category-level object detection, pose estimation and reconstruction from stereo images. In *European Conference on Computer Vision*, pp. 332–349. Springer, 2025.
- Zhang, D., Wang, C., Wang, W., Li, P., Qin, M., and Wang, H. Gaussian in the wild: 3d gaussian splatting for unconstrained image collections. *arXiv preprint arXiv:2403.15704*, 2024b.
- Zhao, H., Shi, J., Qi, X., Wang, X., and Jia, J. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2881–2890, 2017.
- Zou, Y., Ding, Y., Qiu, X., Wang, H., and Zhang, H. M2depth: Self-supervised two-frame multi-camera metric depth estimation. In *European Conference on Computer Vision*, pp. 269–285. Springer, 2025.
- Zuo, X., Samangouei, P., Zhou, Y., Di, Y., and Li, M. Fmgs: Foundation model embedded 3d gaussian splatting for holistic 3d scene understanding. *International Journal of Computer Vision*, pp. 1–17, 2024.

Supplementary Material for LANGSURF

A. Failure Cases

We showcase some of our failure cases in Fig. 11, which are primarily summarized by two factors: 1) the words in the text are unable to detect precisely (top of the figure); 2) the datasets exhibit skewed distributions with a long-tail issue, making it difficult to recognize objects of the tail class (bottom of the figure). Our future work will explore training the language surface fields to identify texts in the images for the open-vocabulary task and balancing the data between the head and tail classes.

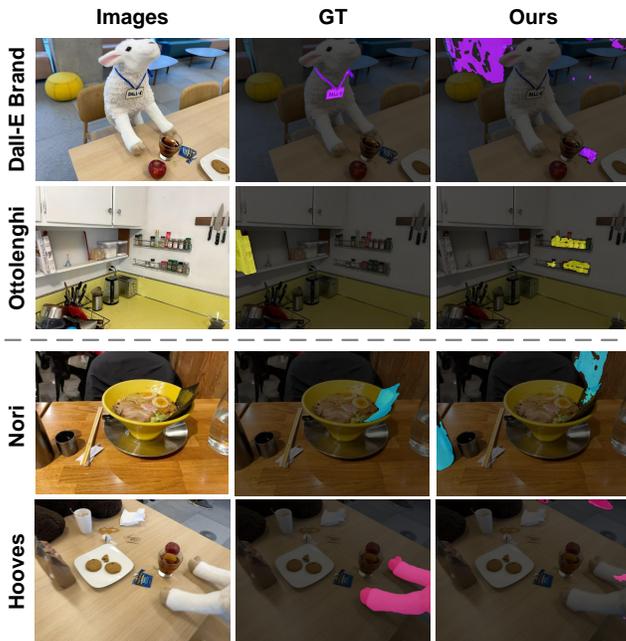


Figure 11: **Visualization of failure cases in LERF dataset (Kerr et al., 2023)**. Top: failure cases where texts in the images cannot be detected. Bottom: failure cases where objects of the tail class cannot be recognized.

B. Implementation Details

B.1. Metrics

2D evaluation. For LERF (Kerr et al., 2023) dataset, we adopt mIoU and localization mAcc for open-vocabulary semantic segmentation. For localization mAcc, we begin by applying a mean convolution filter of size 20 to smooth the values in the relevancy maps, which helps reduce the impact of outliers. We then identify the maximum score within the relevancy maps and consider it accurately localized if its coordinates fall within the ground truth bounding box. In the end, we calculate the average accuracy across all classes.

3D evaluation. For 3D open-vocabulary segmentation on ScanNet, we adopt Semantic F-Score for evaluation. We select each text $t \in T$ to query all points. For the queried points P_t and ground truth points P_t^* , the Semantic F-Score is computed as follows:

$$\begin{aligned}
 Precision_t &= \text{mean}_{p \in P_t} \left(\min_{p^* \in P_t^*} \|p - p^*\| < \tau \right), \\
 Recall_t &= \text{mean}_{p^* \in P_t^*} \left(\min_{p \in P_t} \|p - p^*\| < \tau \right), \\
 S-F-Score_t &= \frac{2 \times Precision_t \times Recall_t}{Precision_t + Recall_t}.
 \end{aligned} \tag{8}$$

where $S-F-Score_t$ is Semantic F-Score for text t and τ is a threshold which we set to 0.05 here.

Table 5: Quantitative Results of ScanNet (Dai et al., 2017). We present Semantic F-Score for each text.

Methods	scene0085_00							scene0616_00						
	bed	chair	curtain	picture	wall	floor	door	wall	floor	chair	table	door	window	picture
LangSplat	08.68	00.00	15.8	02.97	13.44	01.83	00.61	10.34	03.98	28.79	04.55	01.02	02.91	05.26
GS-Group	45.24	06.62	24.32	05.19	15.96	34.99	00.18	04.61	05.56	26.03	14.6	04.16	01.32	10.01
Ours	65.13	87.05	60.39	53.82	29.72	55.17	17.63	12.98	45.09	74.43	66.72	23.03	27.76	06.32

Methods	scene0114_02							scene0617_00						
	wall	floor	cabinet	chair	table	door	window	desk	wall	floor	cabinet	sofa	table	curtain
LangSplat	05.26	01.11	22.15	37.48	00.00	01.33	02.06	36.13	00.49	09.54	04.18	42.01	11.56	00.57
GS-Group	06.01	02.81	00.00	38.21	00.83	01.91	08.97	34.72	00.4	15.38	00.00	30.26	17.07	00.64
Ours	15.19	19.11	48.17	74.93	08.67	01.98	21.55	55.84	08.52	50.81	23.68	72.02	26.97	17.58

B.2. Datasets

In the original LERF dataset, multiple objects with the same categories are ignored in the previous annotations, which leads to inaccurate evaluations. To this end, we refine the annotation errors in LERF datasets for comprehensive evaluation, as shown in Fig. 12. Notably, we don't introduce new categories or new images during the refining procedure.

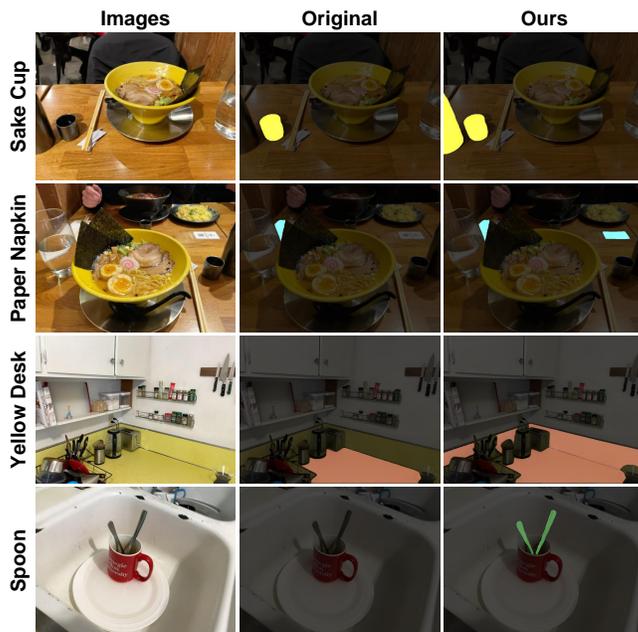


Figure 12: Visualization of our refined annotations in LERF dataset (Kerr et al., 2023). We showcase some differences between the original annotations and our refined annotations.

B.3. Details of Objects Removal

The object removal task can be divided into several steps:

1. We reconstruct the scene with our language-embedded surface field.
2. We perform a text-guided query for language Gaussians and select the Gaussians with high activation scores.
3. With the selected Gaussians, we compute the convex hull of the selected Gaussians and then identify all Gaussians from the original Gaussians that are inside this convex hull. After that, we delete the identified Gaussians.

B.4. Details of Objects Editing

The object editing task can be split into the following steps:

1. We reconstruct the scene with our language-embedded surface field.
2. We perform a text-guided query for language Gaussians and select the Gaussians with high activation scores.
3. We render novel view object masks, RGB images, and depth maps of the edited object, and then feed them into the GaussCtrl (Wu et al., 2024b) to generate edited images.
4. After that, we finetune our Gaussian model for 3,000 steps using the edited images.

C. Additional Experiments

C.1. Qualitative Results.

We visualize the 2D text-queried maps in LERF (Kerr et al., 2023) and ScanNet (Dai et al., 2017) datasets. As shown in Fig. 13, our method performs higher activation scores in the queried objects, which reveals our language-embedded surface field enables better semantic consistency within the same object.

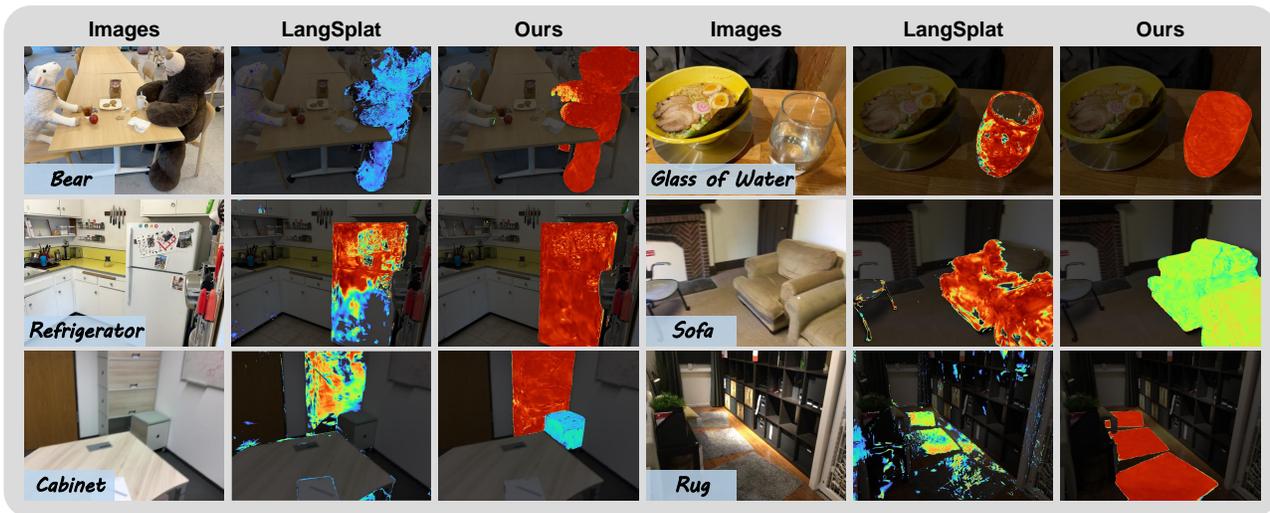


Figure 13: **Visualization of our 2D text-query heatmaps in both datasets.** We showcase our results compared with LangSplat (Qin et al., 2024a).

Additionally, we present more results of 2D open-vocabulary semantic segmentation on both LERF (Kerr et al., 2023) and ScanNet (Dai et al., 2017) datasets. As shown in Fig. 14 and Fig. 15, our method shows a stronger ability to segment 2D objects than other methods, demonstrating the capability of our language-embedded surface field to capture contextual and spatial information.

C.2. Quantitative Results.

In the 3D open-vocabulary segmentation task on ScanNet (Dai et al., 2017), we provide additional details about the Semantic F-Score for each text in Tab. 5. We compare our method with LangSplat and Gaussian-Grouping (GS-Group). The results show that our method outperforms the other techniques, demonstrating its effectiveness in aligning language-embedded surface fields with object surfaces, which leads to improved accuracy in 3D perception.

C.3. Ablations

We also demonstrate the details of our proposed components in ScanNet (Dai et al., 2017). In the 3D open-vocabulary task, we provide the Semantic F-Score for each text in Tab. 6. We observe that the Hierarchical-Context Awareness Module

Table 6: **Ablation of Proposed Components in ScanNet.** We present the Semantic F-Score for each text in the 3D open-vocabulary task.

Components	bed	chair	curtain	picture	wall	floor	door	Overall
w/o HCAM	57.65	73.25	34.16	05.27	24.87	17.04	01.62	30.55
w/o \mathcal{L}_{geo}	59.94	92.39	59.82	50.75	27.91	45.21	11.13	49.59
w/o \mathcal{L}_{sg}	62.18	86.59	62.95	55.58	25.95	52.45	10.54	50.89
w/o \mathcal{L}_{s3d}	64.85	84.28	65.52	50.47	26.00	51.60	08.18	50.13
w/ all	64.80	88.97	59.04	61.91	25.14	51.70	11.51	51.87

(HCAM) makes significant improvements. Additionally, the \mathcal{L}_{geo} , \mathcal{L}_{sg} and \mathcal{L}_{s3d} are complementary, with performance drops when removing either one.

C.4. Video Demo

In Figure 1 of our main paper, we visualize our language surface field and the results of its applications. For visualization of language features, both our method and LangSplat (Qin et al., 2024a) learn a 3-dimensional language embedding for Gaussian primitives, allowing us to maintain color consistency between the two approaches and ensure a fair comparison. We strongly recommend reviewers refer to our video demo in the attached compressed file to observe our language surface fields and applications including 2D/3D open-vocabulary segmentation, object removal and object editing. The video demonstrates that our language surface field accurately represents the language field in 3D to align the surface of objects, and thus it can be transferred to multiple applications.

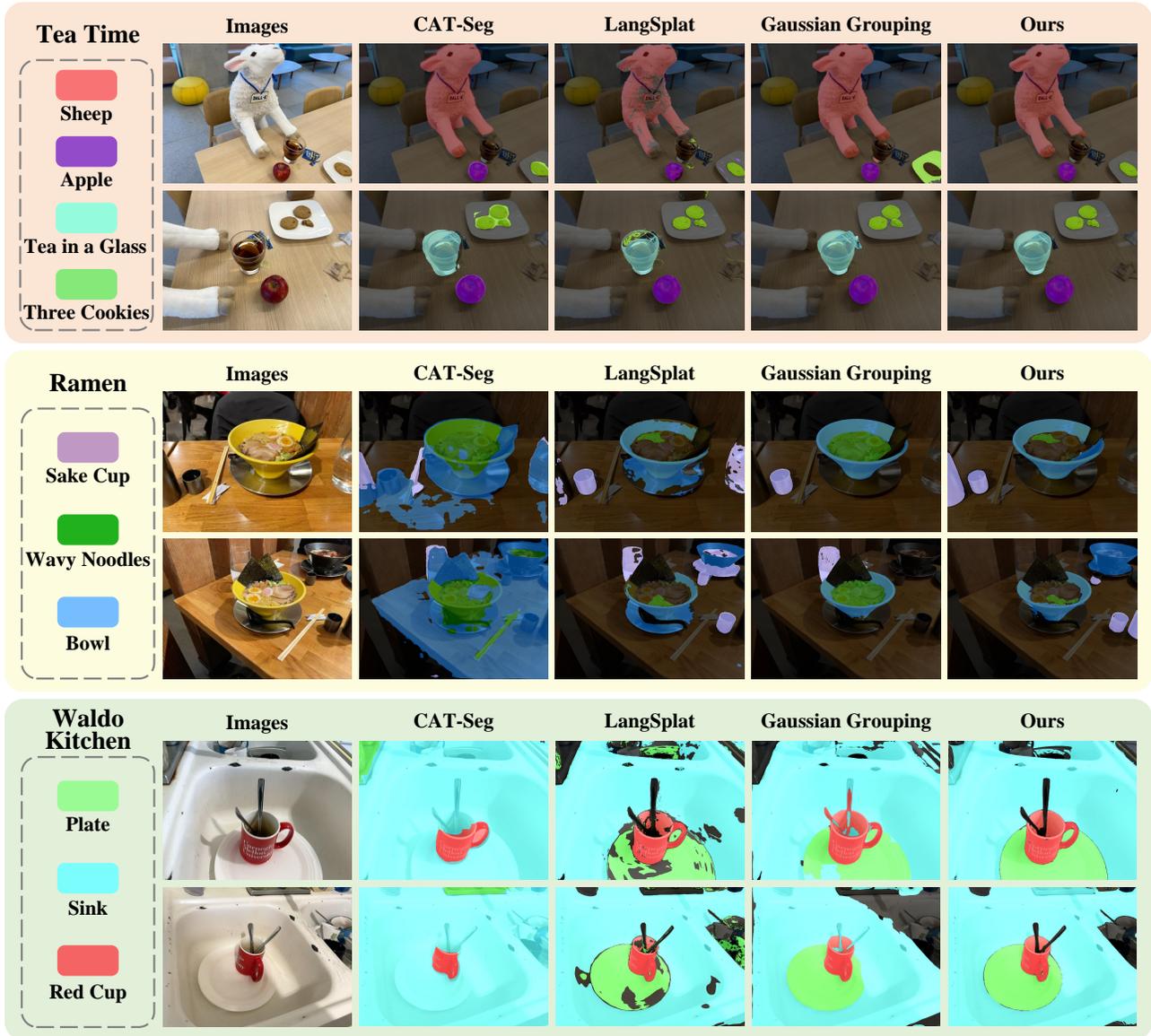


Figure 14: Visualization of our 2D open-vocabulary masks in LERF dataset (Kerr et al., 2023). We showcase 2 images for each scene and present the texts on the left.

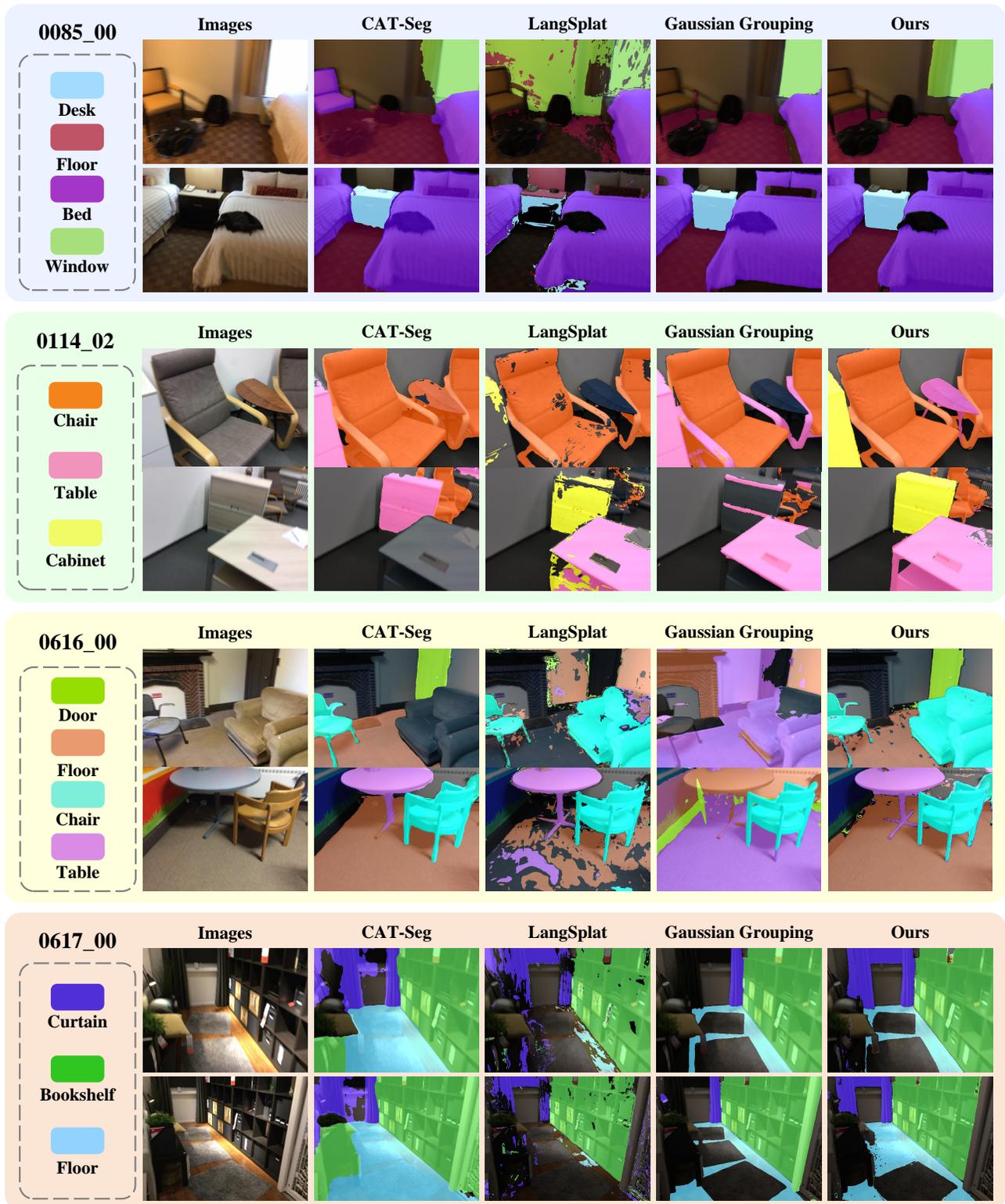


Figure 15: Visualization of our 2D open-vocabulary masks in ScanNet dataset (Dai et al., 2017). We showcase 2 images for each scene and present the texts on the left.