# GLS: Geometry-aware 3D Language Gaussian Splatting

Jiaxiong Qiu[1] , Liu Liu[1] , Xinjie Wang[1] , Tianwei Lin[1] , Wei Sui[2] , Zhizhong Su[1]
[1]Horizon Robotics, Beijing, China    [2]D-Robotics, Beijing, China

## Abstract

*Recently, 3D Gaussian Splatting (3DGS) has achieved impressive performance on indoor surface reconstruction and 3D open-vocabulary segmentation. This paper presents GLS, a unified framework of 3D surface reconstruction and open-vocabulary segmentation based on 3DGS. GLS extends two fields by improving their sharpness and smoothness. For indoor surface reconstruction, we introduce surface normal prior as a geometric cue to guide the rendered normal, and use the normal error to optimize the rendered depth. For 3D open-vocabulary segmentation, we employ 2D CLIP features to guide instance features and enhance the surface smoothness, then utilize DEVA masks to maintain their view consistency. Extensive experiments demonstrate the effectiveness of jointly optimizing surface reconstruction and 3D open-vocabulary segmentation, where GLS surpasses state-of-the-art approaches of each task on MuSHRoom, ScanNet++ and LERF-OVS datasets. Project webpage: https://jiaxiongq.github.io/GLS_ProjectPage.*

## 1. Introduction

Surface reconstruction[10, 22, 24, 57–59] and 3D open-vocabulary segmentation[23, 39, 44, 49, 52, 60] based on 3DGS[29], being widely applied in AR/VR[27] and embodied intelligence[8, 38] due to its capabilities of efficient training and real-time rendering. Recently, notable works have achieved significant progress in both areas. For surface reconstruction, SuGaR[22] proposes regularization terms to align Gaussians and scene surface, then uses Possion reconstruction[28] to extract mesh from Gaussians. For 3D open-vocabulary segmentation, LangSplat[39] and OpenGaussian[49] successfully introduces SAM[31] and CLIP[26] to 3DGS. Gaussian grouping[52] utilizes a universal temporal propagation model DEVA[14] to obtain consistent object masks across views and propose a 3D regularization term to grouping Gaussians. However, these methods only focus on one task and suffer from unstable performance in complex indoor scenes as Fig. 1 shows.

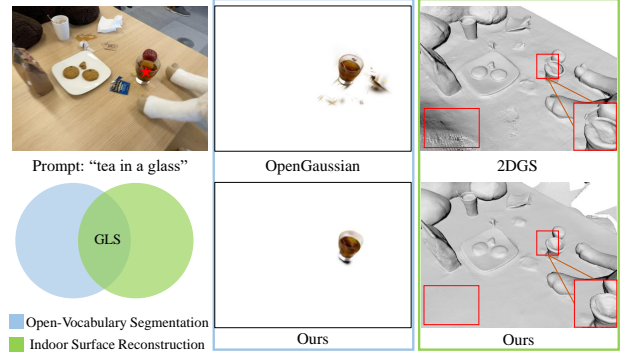In this work, we aim to achieve efficient and robust in-



Figure 1. Indoor surface reconstruction and 3D open-vocabulary segmentation. Shadow and high-light regions make state-of-the-art methods struggle in indoor scenes. Our proposed GLS jointly optimizes two tasks based on 3DGS and achieves much better results than OpenGaussian[49] and 2DGS[24].

door surface reconstruction and 3D open-vocabulary segmentation based on 3DGS. Our goal has two main motivations. On the one hand, the 2D open-vocabulary supervision[26, 31] is naturally view-inconsistent, which easily results in noise body and blur boundary of the segmented object from Gaussians. The accurate scene surface consists of sharp and smooth object surfaces, illustrating that Gaussians are mainly distributed on objects and preserve the object boundaries. This property can make the object segmentation results from Gaussians cleaner and sharper. On the other hand, due to the complex materials and lighting conditions of the indoor scene, the shadow and high-light regions on texture-less and reflective objects always result in noisy surfaces. Fortunately, accurate object masks erase the interference details on the object's surface. Hence, the object segmentation results can supply the smoothness prior to reducing the reconstruction noises of these objects. In general, the optimization goals of the two tasks can be considered to be the same.

Based on above motivations, we introduce GLS, which leverages complementary between surface reconstruction and 3D open-vocabulary segmentation to boost the performance of 3DGS in two tasks. The framework is presented in Fig. 3. Specifically, we first introduce the normal prior[3] to regularize the surface normal estimated from rendered depth. Then we analyzed different situations of rendered
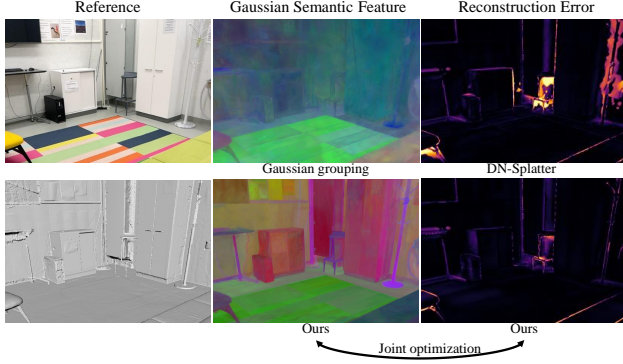
Figure 2. Effect of joint optimization between indoor surface reconstruction and 3D open-vocabulary segmentation. Our method significantly improves the quality of Gaussian semantic features and surface reconstruction.

depth under different normal errors and propose a regularization term to enhance the sharpness of the scene surface. To integrate the open-vocabulary information, we add Gaussian semantic features into vanilla 3DGS and then utilize the consistent object masks from DEVA and image features from CLIP to supervise them. In addition, we consider the clip features as the smoothness prior to strengthen the accuracy of texture-less and reflective surfaces. Finally, we adopt TSDF fusion[37] from rendered depth to extract scene mesh, then compute the similarity between input text embeddings and learned semantic embeddings to acquire object masks. We conduct extensive experiments on MuSHRoom[41], ScanNet++[54] and LERF-OVS [30, 39] datasets. As Fig. 2 shows, the superior performance of our model on both tasks demonstrates the effectiveness of connecting indoor surface reconstruction and 3D open-vocabulary segmentation in 3DGS. GLS makes the reconstructed scene surface interactive. The application demos can be visualized in the supplementary videos. In summary, our technical contributions can be listed as follows:

1. We design a novel framework based on 3DGS, by jointly optimizing surface reconstruction and 3D open-vocabulary segmentation in complex indoor scenes.
2. We propose two novel regularization terms with the help of geometric and semantic cues, to facilitate the sharpness and smoothness of reconstructed scene surfaces and segmented objects.
3. Our method inherits the training and rendering efficiency of 3DGS and achieves state-of-the-art accuracy on surface reconstruction and 3D open-vocabulary segmentation tasks.

## 2. Related Work

### 2.1. Neural Rendering

Neural Radiance Fields (NeRF)[35] based on implicit representations have achieved remarkable advancements in novel view synthesis. Some NeRF-based approaches[4–6] concentrate on extending NeRF to more challenging scenes, while they are still limited by low training and rendering efficiency. Alternatively, other NeRF-based methods[13, 17, 36, 40, 45, 55] utilize explicit representations such as voxels, grids, point clouds and meshes. They faithfully reduce the large computational cost of neural networks. The recent technologies[29, 33, 56] of neural rendering based on alpha-blend rendering have further advanced the rendering speed and quality, by optimizing the attributes of the explicit 3D Gaussians. Our framework is built upon 3DGS and focuses on optimizing it to tackle both surface reconstruction and 3D open-vocabulary segmentation tasks.

### 2.2. 3DGS-Based Surface Reconstruction

Surface reconstruction from multi-view images[19–21, 42, 43, 57] is a challenging task in computer vision and graphics. Recently, extracting scene surfaces from Gaussian primitives has become a popular research topic. SuGaR[22] utilizes Poisson reconstruction[28] to extract mesh from sampled Gaussians and encourages Gaussians attach to the scene surface by several regularization terms, for enhancing reconstruction accuracy. However, due to the disorder of the Gaussians and the complexity of the scene, fitting Gaussians on geometric surfaces is challenging. 2DGS[24] replaces 3D Gaussian primitives with 2D Gaussian disks for efficient training. They also introduce two regularization terms of depth distortion and normal consistency to reconstruct a smooth scene surface. PGSR[10] estimates the unbiased depth to disable the degeneration of the blending coefficient, and introduces the multi-view regularization term to achieve global-consistent reconstruction. DN-Splatter[46] and VCR-GauS[11] leverage normal or depth priors from generalizeable models[3, 7, 18], to learn effective surfel representations. FDS[12] introduce the optical flow into 3DGS for more accurate geometry.

Due to the lack of object attributes, these methods struggle to reconstruct sharp objects against complex lighting conditions of an indoor scene. To address these issues, we propose combining 3D open-vocabulary segmentation and surface reconstruction. In contrast with previous related works[47, 48, 51], our method releases the need for ground-truth segmentation supervision and achieves highly efficient training speed.

### 2.3. 3DGS-Based open-vocabulary segmentation

Recent works[25, 50, 61] about open-vocabulary scene understanding have significant progress by integrating 2D foundational models with 3D point cloud representations. These methods project a 3D point cloud into 2D space for zero-shot learning of aligned features. LERF[30] designs a pipeline of 3D open-vocabulary segmentation by distilling features from CLIP into NeRF. 3GS-based methods add the
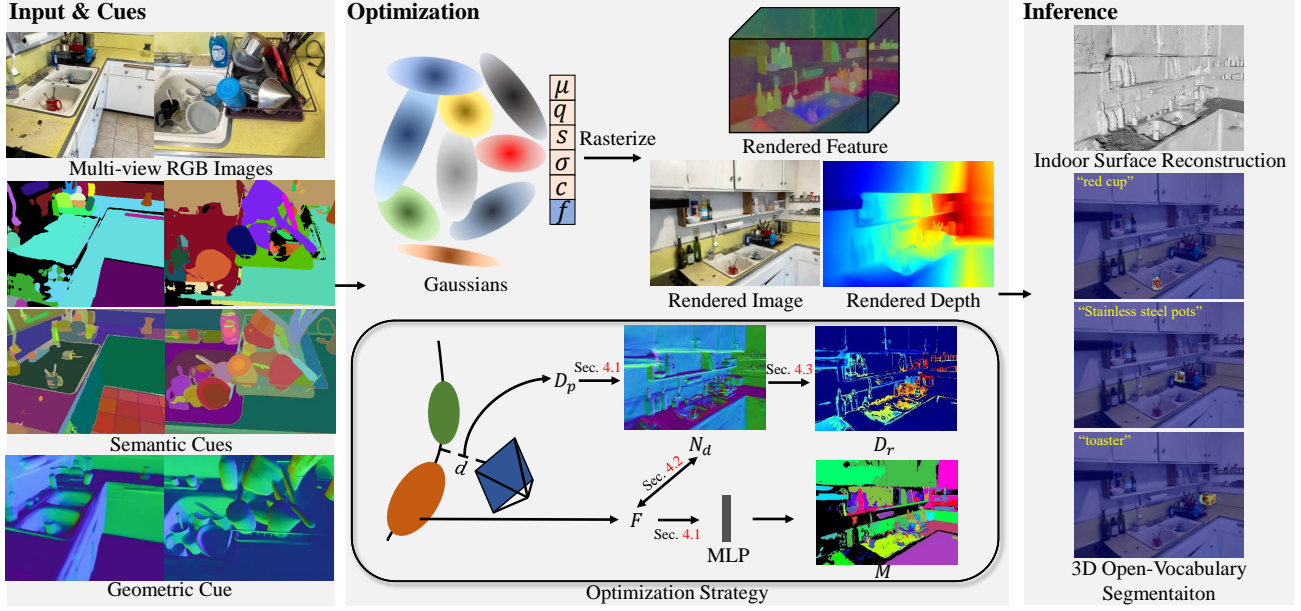
Figure 3. Overview of our framework. We adopt geometric and semantic cues produced by generalizable models to jointly strengthen the reconstruction and segmentation quality of 3DGS. Semantic cues consist of view-consistent segmentation results of DEVA[14] and object-level CLIP features generated from LangSplat[39]. Geometric cue comes from DSINE[3]. Our framework achieves accurate indoor surface reconstruction and 3D open-vocabulary segmentation via effective regularization terms, which faithfully enable Gaussian primitives to distribute along the object surface.

semantic attribute into 3D Gaussians and are supervised by 2D scene priors[26, 31] to overcome large computational cost of NeRF-based methods. LangSplat[39] designs a per-scene autoencoder to reduce the dimension of CLIP features on multi-level latent spaces, this scheme generates clear boundaries of rendered features. LEGaussians[44] integrate uncertainty with CLIP and DINO[9] image features to Gaussians. To learn high-dimensional semantic features, Feature3DGS[60] develops a parallel Gaussian rasterizer. Gaussian Grouping[52] introduces view-consistent masks from DEVA and proposes a 3D local consistency regularization term to strengthen the segmentation accuracy. OpenGaussian[49] focuses on 3D point-level open-vocabulary understanding via CLIP features and proposes a two-level codebook scheme with SAM masks to refine the rendered mask. Similar to these methods, we introduce CLIP features and SAM masks to supervise the semantic features of 3DGS.

## 3. Preliminary

Initialized with point clouds and colors produced by SfM[43], 3DGS[29] adopt differentiable 3D Gaussian primitives $\{G\}$ to explicitly represent a scene. Each primitive is parameterized by the Gaussian function:

$$G(x|\mu, \Sigma) = e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}, \quad (1)$$

where $\mu \in \mathbb{R}^3$ and $\Sigma \in \mathbb{R}^{3\times3}$ are the center and the covariance matrix of spatial points $x$ respectively. $\Sigma$ consists of a

scaling vector $s \in \mathbb{R}^3$ and a quaternion vector $q \in \mathbb{R}^4$.

3DGS enables an efficient alpha-blending procedure for real-time rendering. Given a camera view, 3D Gaussian primitives are projected into viewing space to be 2D Gaussians, which are sorted by the z-buffer strategy and alpha-composited by the volume rendering equation[35] to generate pixel colors $C$:

$$C = \sum_{i\in N} c_i \alpha_i \prod_{j=1}^{i-1}(1-\alpha_j), \quad (2)$$

where $N$ is the number of 3D Gaussian primitives, $c \in \mathbb{R}^3$ is the color of a Gaussian primitive estimated from spherical harmonics and the viewing direction. $\alpha$ is the blending co-efficient determined by the opacity $\sigma$. Similar to rendered color $C$, the rendered alpha $A$, depth $D$ and rendered semantic features $F$ can be denoted by:

$$A = \sum_{i\in N} \alpha_i \prod_{j=1}^{i-1}(1-\alpha_j),$$

$$D = \sum_{i\in N} d_i \alpha_i \prod_{j=1}^{i-1}(1-\alpha_j), \quad (3)$$

$$F = \sum_{i\in N} f_i \alpha_i \prod_{j=1}^{i-1}(1-\alpha_j).$$

where $d$ is the distance between the 2D Gaussian point and the camera center, $f$ is the semantic feature of each Gaussian.

3

## 4. Method

Given multi-view RGB images captured by a camera in an indoor scene, our goal is to jointly reconstruct the scene and open-vocabulary objects. To achieve this goal, we introduce GLS, a novel framework based on 3DGS. As shown in Fig. 3, our framework consists of three procedures. In the input procedure, we use the generalizable model SAM[31], DEVA[14] and CLIP[26] to produce 2D consistent semantic masks $\hat{M}$ and object-level features $\hat{F}$. Then we adopt the generalizable model[3] of surface normal estimation to acquire the geometric cue $\hat{N}$. In the optimization procedure, we utilize the semantic and normal priors for regularization. We first follow previous approaches to regularize the rendered color, depth and semantic feature (Sec. 4.1). Then we propose a novel smoothness term (Sec. 4.2) to tackle texture-less regions and a novel constraint by analyzing the normal error of Gaussians (Sec. 4.3) to refine object structures. In the inference procedure, our model reconstructs the indoor surface and selects the target object by the open-vocabulary text simultaneously.

### 4.1. Leveraging 2D Semantic and Geometric Cues

Previous works[46, 49] either only consider surface reconstruction, or only perform the 2D open-vocabulary segmentation. We propose combining two tasks to jointly reconstruct the scene surface and segment objects of the scene.

For indoor surface reconstruction, in contrast to DN-splatter[46], we utilize the TSDF Fusion[34, 37] to extract the mesh. Hence, we focus on optimizing the rendered depth $D$. As demonstrated in 2DGS[24], the local smoothness of rendered depth is challenging and important for the surface reconstruction of 3DGS. To smooth rendered depth, we introduce the normal prior $\hat{N}$. We follow the manner of PGSR[10] to estimate the gradients of 3D points projected from rendered depth $D$, and take them as the local surface normal $N_d$. Then we leverage $\hat{N}$ to regularize it by:

$$\mathcal{L}_n = \sum_{i,j} A(1 - N_d^T \hat{N}), \qquad (4)$$

For 3D open-vocabulary segmentation, we directly use the view-consistent segmentation results of DEVA[14] as the supervision of the rendered mask $M$. Specifically, given the rendered features $F$, we use an MLP layer to increase its feature dimension to the number of total categories first. The softmax function and standard cross-entropy loss $\mathcal{L}_m$ are adopted for final classification, which is defined as: $\mathcal{L}_m = -\sum_{i=1}^{S} y_i \log(\hat{y}_i)$, where $S$ is the number of classes, $y_i$ is the true label (0 or 1) for the $i$-th class, and $\hat{y}^i$ is the predicted probability for the $i$-th class. To enhance the open-vocabulary capabilities of our model, we employ the CLIP features $\hat{F}$ to supervise $F$ by:

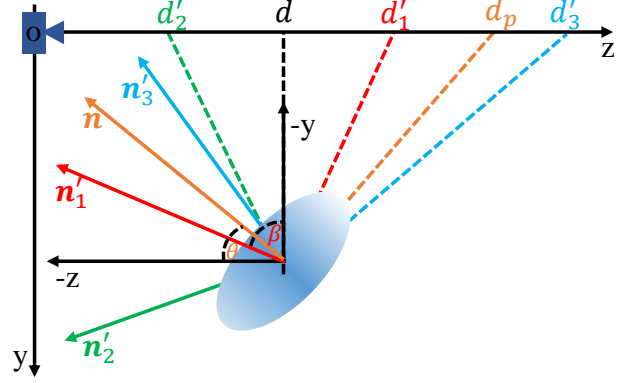$$\mathcal{L}_{clip} = ||F - \hat{F}||^1, \qquad (5)$$



Figure 4. Depth refinement guided by the normal error between the rendered normal $n$ and the ideal normal $n'$. We analyze three conditions of $n'$. The ideal depth $d'$ of each condition is the unbiased depth proposed by PGSR[10].

### 4.2. Semantic-feature Guided Normal Smoothing

Big surfaces like floors and desktops are generally over-smoothing in an indoor scene. Although $\mathcal{L}_n$ supply local smoothness for reconstruction, it still struggles in shadow and high-light regions of these surfaces. High-weight $\mathcal{L}_n$ causes an over-smoothing effect and then makes some small objects disappear, while low-weight $\mathcal{L}_n$ disables the local smoothness of rendered depth.

To resolve this issue, we propose introducing clip features to help smooth big surfaces, then we can employ $\mathcal{L}_n$ with low weight for protecting small object boundaries. Essentially, clip features supply high-dimension object-level smoothness which can implicitly reduce the inference of big surfaces. Concretely, we first introduce SAM to select objects that occupy the top-$k$ area by the mask $M_o$. and design a novel regularization term $\mathcal{L}_s$ for smoothing them as:

$$\mathcal{L}_s = M_o(|\partial_x N_d|e^{-||\partial_x \hat{F}||^1} + |\partial_y N_d|e^{-||\partial_y \hat{F}||^1}), \qquad (6)$$

### 4.3. Normal-error Guided Depth Refinement

We take the unbiased depth $D_p$ from PGSR[10] to replace the rendered depth and as the input of TSDF fusion procedure, which denoted by: $D_p = \frac{D}{cos(\theta)}$, where $\theta$ is the angle between the direction of the intersecting ray and the rendered normal $n$[15]. However, due to the ambiguity of the shortest Gaussian axis, $n$ always occurs large error when compared to the ideal normal $n'$. As Fig. 4 shows, given $n$ and the corresponding unbiased depth $d_p$, there are three main conditions of the ideal normal $n'$ as follows:

1. $n'_1$ exists between $n$ and $-z$. The corresponding ideal depth $d'_1 \in [d, d_p]$.
2. $n'_2$ exists between $-z$ and $y$. The corresponding ideal depth $d'_2 \in [0, d]$.
3. $n'_3$ exists between $n$ and $-y$. The corresponding ideal depth $d'_3 \in [d_p, d]$.

To determine the $i$-th $n'$, we propose using the angle $\alpha$ between $n'_i$ and $-y$ along with $\theta' = 90° - \theta$. Then we
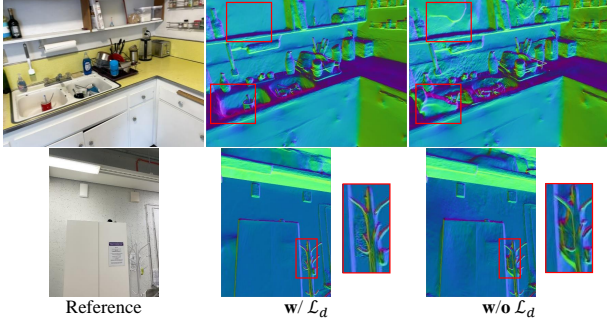
Figure 5. Effect of $\mathcal{L}_d$. The surface normal is estimated from the unbiased depth. $\mathcal{L}_d$ faithfully enhances the sharpness and smoothness of the reconstructed surface.

have $cos(\alpha) = \boldsymbol{n}' \cdot -\boldsymbol{y}$ and $cos(\theta') = \boldsymbol{n} \cdot -\boldsymbol{y}$. Consequently, we can get $\boldsymbol{n}_1' \in \{\boldsymbol{n}'|cos(\alpha)>cos(\theta')>0\}$, $\boldsymbol{n}_2' \in \{\boldsymbol{n}'|cos(\alpha)<0\}$ and $\boldsymbol{n}_3' \in \{\boldsymbol{n}'|0<cos(\alpha) \le cos(\theta')\}$.

In practice, we take $\hat{N}$ to approximate $\boldsymbol{n}'$ in 2D space. Then we can acquire three masks ($M_1$, $M_2$ and $M_3$) through the above three conditions. Ulteriorly, to make $d_p$ close to $d'$, we refine $D_p$ by reconstituting depth in each mask. Specifically, for the first condition, we introduce the rendered alpha $A$ to integrate $D$ and $D_p$ in $M_1$ by: $M_1(AD_p + D - AD)$. For the second condition, we choose $M_2D$ as the target depth. For the last condition, we choose $M_3D_p$ as the target depth. Finally, the whole target depth $D_r$ can be denoted by:

$$D_r = M_1(AD_p + D - AD) + M_2D + M_3D_p, \quad (7)$$

Then we conduct a novel regularization term $\mathcal{L}_d$ between $D_r$ and $D_p$ by:

$$\mathcal{L}_d = 1 - e^{-||D_p - D_r||^1}, \quad (8)$$

Furthermore, we choose the pixels under $\{N_d^T \hat{N}<0.9\}$ to compute the value of $\mathcal{L}_d$. As Fig. 5 demonstrates, $\mathcal{L}_d$ encourages the Gaussians closely attach to the object surface and then improves the sharpness and smoothness of the reconstructed surface.

### 4.4. Optimization

We adopt the photometric loss $L_c$ from vanilla 3DGS[29]. All loss functions are simultaneously optimized by training from scratch. The total loss function $\mathcal{L}$ can be defined as:

$$\mathcal{L} = \mathcal{L}_c + \alpha_n\mathcal{L}_n + \alpha_m\mathcal{L}_m + \alpha_{clip}\mathcal{L}_{clip} + \alpha_d\mathcal{L}_d + \alpha_s\mathcal{L}_s. \quad (9)$$

We set $\alpha_n = 0.07$, $\alpha_m = 0.3$, $\alpha_{clip} = 1.0$, $\alpha_d = 0.01$ and $\alpha_s = 0.5$ in our experiments.

## 5. Experiments

### 5.1. Settings

**Datasets & Metrics.** For 3D open-vocabulary segmentation, we follow LangSplat[39] to conduct experiments on the LERF-OVS dataset[30]. We adopt the evaluation metrics from Gaussian Grouping. Specifically, we use the text query to select 3D Gaussians, then calculate the average IoU (mIoU) and boundary IoU (mBIoU) accuracy between rendered masks and annotated object masks. For indoor surface reconstruction, we conduct experiments on two RGBD datasets: MuSHRoom[41] (only contains: 'coffee_room', 'computer', 'honka', 'kokko' and 'vr_room') and ScanNet++[54] (only contains: '8b5caf3398' and 'b20a261fdf'). We and the same tool of DN-Splatter[46] to evaluate mesh quality through five metrics: Accuracy, Completion, Chamfer-L1 distance, Normal Consistency and F-scores. For novel view synthesis, we follow standard PSNR, SSIM and LPIPS metrics for rendered images.

**Baselines.** We compare our model against a series of baseline approaches on two tasks. a) 3D open-vocabulary segmentation: Langsplat[39], Gaussian Grouping[52] and OpenGaussian[49]. We deploy the scheme of OpenGaussian to render the selected objects. Specifically, we compute the similarity between 3D Gaussian semantic features and text features, then select Gaussians with high similarity and obtain the rendered object by the rasterizer. b) 3DGS-based indoor surface reconstruction: 2DGS[24], PGSR[10], DN-Splatter[46] and FDS[12]. For fair comparison among all methods, we take the LiDAR points of iPhone as the default input point cloud on both datasets, and re-ran their source codes released in GitHub in 5 scenes. To align with the setting of DN-Splatter, we introduce the sensor depth to supervise $D_p$ by Mean Absolute Error loss.

**Implementation details.** Our code is built based on PGSR[10]. The densification strategy is adopted from AbsGS[53]. We train GLS for 30k iterations, consuming about 40 minutes on a single NVIDIA RTX 4090 GPU. For indoor surface reconstruction, we first generate rendered depth in each training view, followed by performing the TSDF Fusion[34, 37] to extract the mesh in the TSDF field. Subsequently, for 3D open-vocabulary segmentation, we first compute the similarity between 3D Gaussian semantic features and text features, then filter Gaussians with high similarity to render and reconstruct the selected object. More details can be seen in the supplementary materials.

### 5.2. Indoor Surface Reconstruction

**Quantitative Comparisons.** For indoor surface reconstruction, we conduct comparisons on two real-world datasets, including MuSHRoom[41] and ScanNet++[54]. We report the metric values in Tab. 1 and Tab. 7. It can be seen that our method outperforms other 3DGS-based approaches[10, 24] without sensor depth among all metrics. When adopting sensor depth as the prior information of scene scale, our model also achieves better performance of sharpness and smoothness than DN-Splatter[46] according to the Accuracy and Normal Consistency metric. Com-

Table 1. Quantitative results of indoor surface reconstruction on MuSHRoom dataset (only contains: 'coffee_room', 'computer', 'honka', 'kokko' and 'vr_room'). Our method outperforms other methods on most metrics. The best metrics are **highlighted**.

| Methods | Sensor Depth | Accuracy↓ | Completion↓ | Chamfer$-L_1$↓ | Normal Consistency↑ | F-score↑ | Time |
|---|---|---|---|---|---|---|---|
| DN-Splatter[46] | ✓ | 0.0402 | **0.0211** | 0.0306 | 0.8449 | 0.8590 | 1.0h |
| Ours | ✓ | **0.0288** | 0.0254 | **0.0271** | **0.8640** | **0.8924** | 0.7h |
| 2DGS[24] | ✗ | 0.1728 | 0.1619 | 0.1674 | 0.7113 | 0.3433 | 0.3h |
| PGSR[10] | ✗ | 0.4957 | 0.7672 | 0.6315 | 0.5566 | 0.1178 | 0.7h |
| FDS[12] | ✗ | 0.0639 | **0.0528** | 0.0584 | 0.8169 | **0.6998** | 1.3h |
| Ours | ✗ | **0.0538** | 0.0582 | **0.0560** | **0.8357** | 0.6922 | 0.7h |

Table 2. Quantitative results of indoor surface reconstruction on ScanNet++ dataset (only contains: '8b5caf3398' and 'b20a261fdf'). The best metrics are **highlighted**.

| Methods | Sensor Depth | Accuracy↓ | Completion↓ | Chamfer$-L_1$↓ | Normal Consistency↑ | F-score↑ | Time |
|---|---|---|---|---|---|---|---|
| DN-Splatter[46] | ✓ | 0.0977 | 0.0431 | 0.0704 | 0.8272 | 0.7094 | 1.0h |
| Ours | ✓ | 0.0640 | **0.0272** | **0.0444** | 0.9064 | **0.8623** | 0.6h |
| 2DGS[24] | ✗ | 0.2440 | 0.4362 | 0.3401 | 0.6343 | 0.1838 | 0.2h |
| PGSR[10] | ✗ | 0.1670 | 0.2188 | 0.1929 | 0.7622 | 0.2227 | 0.6h |
| Ours | ✗ | **0.0861** | **0.1009** | **0.0935** | **0.8578** | **0.4799** | 0.6h |

pared with the state-of-the-art method FDS[12], our results have better smoothness, accuracy and training efficiency.

**Qualitative Comparisons.** As shown in Fig. 6, there are surface reconstruction results produced by different methods on the MuSHRoom dataset. It can be observed that DN-Splatter[46] generates severe noise on scene surfaces and even destroys object structures because of texture-less regions. On the contrary, our model reconstructs smooth and sharp scene surfaces and recovers more thin structures than ground-truth scene surfaces. This observation demonstrates the joint optimization of surface reconstruction and 3D open-vocabulary segmentation can significantly promote the reconstruction quality. For 3DGS-based approaches without sensor depth as supervision, due to the complex camera motion and light conditions, PGSR presents unstable performance and fails in most scenes. When compared to 2DGS, our model handles shadow and high-light regions better and generates cleaner scene surfaces.

Moreover, we further evaluate the generalization ability of all methods on ScanNet++ dataset as Fig. 7 shows. These scenes encode various lighting conditions, which makes reconstructing scene surfaces more challenging. Hence, more noises occur in DN-splatter results, while our results are still sharp and smooth. For 3DGS-based approaches without sensor depth as supervision, PGSR and 2DGS produce meaningless results while our results successfully tackle these scenes well.

### 5.3. 3D open-vocabulary segmentation

**Quantitative Comparisons.** 3D open-vocabulary segmentation drops predefined labels of objects and uses texts as prompts to segment target objects. We conduct comparisons of 3D open-vocabulary segmentation on the widely-used LERF-OVS dataset. Following Gaussian grouping[52], we employ mIoU and mBIoU as the evaluation metrics, which represent the segmentation quality of the selected object. The estimated metric scores are reported in Tab. 3. Our model outperforms state-of-the-art method OpenGaussian[49] over 4% on mIoU and achieves over 20% gain when compared with Gaussian grouping[52] on mBIoU. This fact illustrates that the joint optimization of two tasks improves the smoothness and sharpness of scene surfaces then makes segmentation results smoother and sharper.

**Qualitative Comparisons.** Fig. 8 shows some 3D open-vocabulary segmentation results of different methods. Without the optimization of surface reconstruction, OpenGaussian[49] easily treats other objects as targets and Gaussian grouping[52] exhibits noisy results. Our model successfully identifies the 3D Gaussians relevant to the query text and generates object selection with a clearer boundary.

### 5.4. Ablation study

We conduct sufficient ablation experiments to study the effect of different regularization terms, by disabling each one and enabling others. The results of indoor surface reconstruction and 3D open-vocabulary segmentation are reported in Tab. 4 and Tab. 5 respectively.

**Effect of $\mathcal{L}_n$.** The Normal prior supplies the regions with smooth surfaces, where shadow and highlight exist in an indoor scene. For indoor surface reconstruction, $\mathcal{L}_n$ remarkably improves surface smoothness in terms of the Normal Consistency metric. For 3D open-vocabulary segmentation, $\mathcal{L}_n$ helps the full model to produce accurate boundary of segmentation results by the joint optimization according to the metric mBIoU.
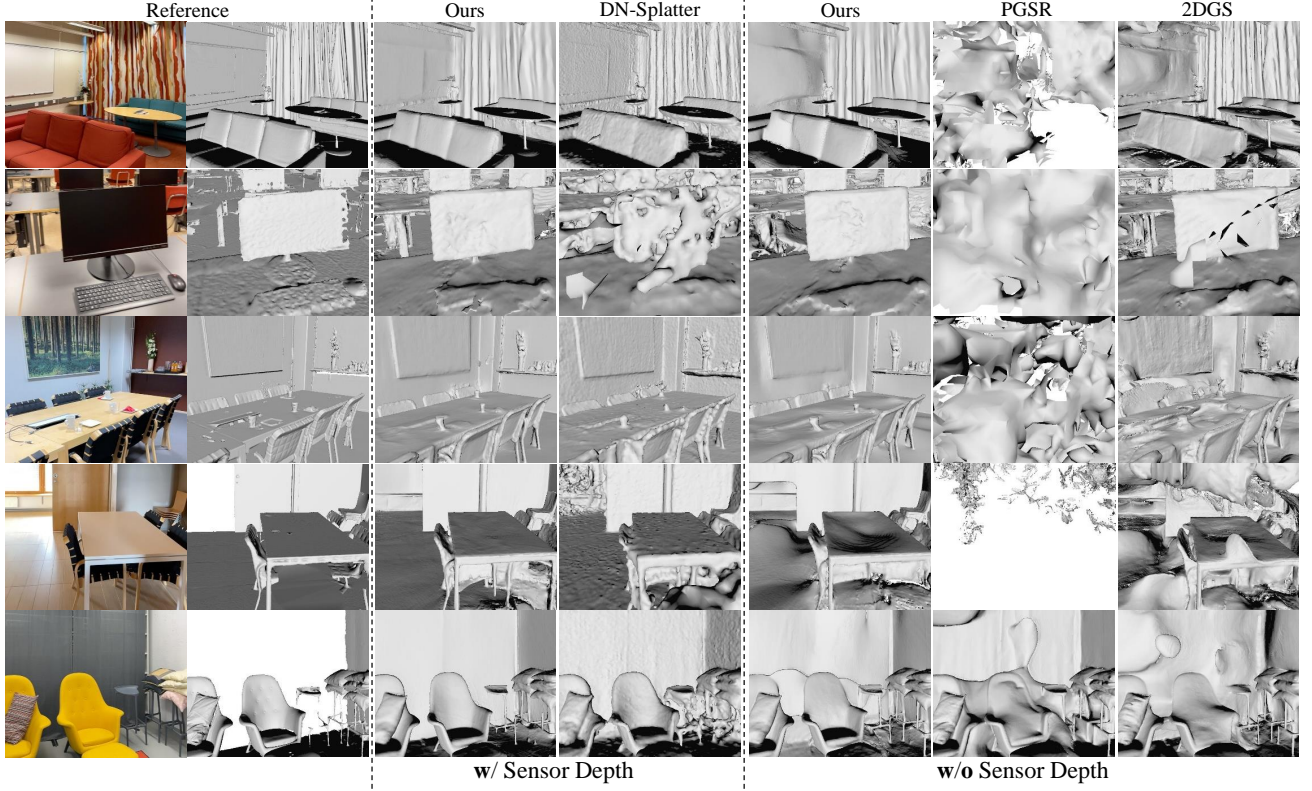
Figure 6. Qualitative comparisons of indoor surface reconstruction on MuSHRoom dataset. PGSR generates unstable results by the default hyperparameters.
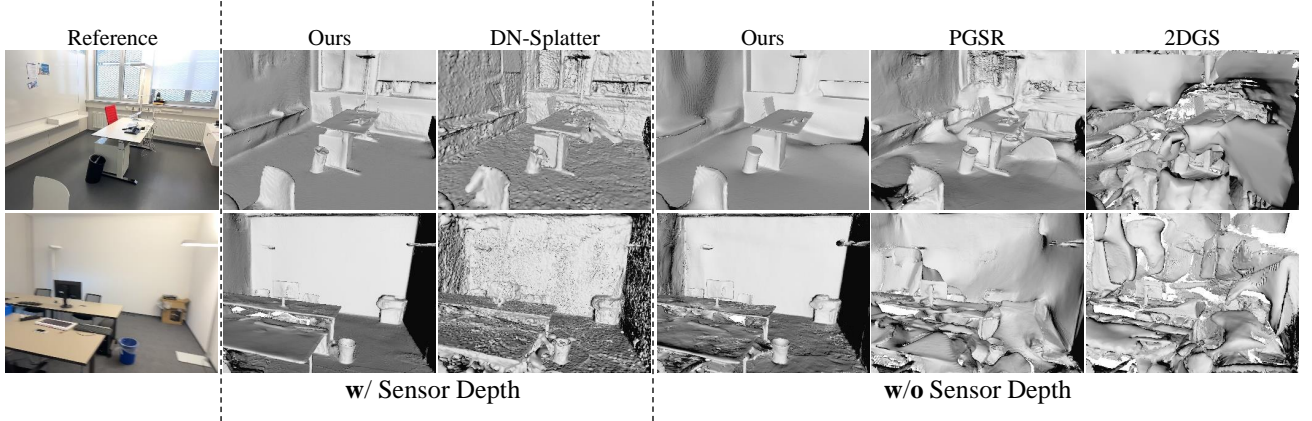


Figure 7. Qualitative comparisons of indoor surface reconstruction on ScanNet++ dataset. The lighting conditions change violently in these scenes.

**Effect of $\mathcal{L}_d$.** Fig. 5 illustrates the visual effects of $\mathcal{L}_d$ for indoor surface reconstruction. For 3D open-vocabulary segmentation, the performance of our model degrades when disabling $\mathcal{L}_d$. Although $\mathcal{L}_d$ is designed for refining the unbiased depth, it also strengthens the quality of segmentation results through joint optimization.

**Effect of $\mathcal{L}_s$.** To determine whether $\mathcal{L}_s$ is necessary for two tasks, we disable it from our full model. $\mathcal{L}_s$ enhances the performance of our model among all metrics, by connecting the surface normal estimated from the unbiased depth and CLIP features.

necting the surface normal estimated from the unbiased depth and CLIP features.

**Effect of $\mathcal{L}_{clip}$.** Then we disable $\mathcal{L}_{clip}$, to verify it is a significant regularization term for two tasks. We introduce CLIP features to strengthen the representation ability of Gaussian semantic features, and which plays the most important role in 3D open-vocabulary segmentation. It also significantly maintains the smoothness and sharpness of reconstructed indoor surfaces via joint optimization.

Table 3. Quantitative results of 3D open-vocabulary segmentation on LERF-OVS dataset. We follow the metrics of LangSplat[39] and OpenGaussian[49] from the paper of OpenGaussian. The best metrics are **highlighted**.

| Methods | mIoU↑ | | | | | mBIoU↑ | | | | | Time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | figurines | ramen | teatime | waldo_kitchen | **Mean** | figurines | ramen | teatime | waldo_kitchen | **Mean** | |
| LangSplat[39] | 10.16 | 7.92 | 11.38 | 9.18 | 9.66 | - | - | - | - | - | 2.2h |
| Gaussian grouping[52] | 15.53 | 17.49 | 22.27 | 26.51 | 20.45 | 13.71 | 13.86 | 19.10 | 18.85 | 16.38 | 0.7h |
| OpenGaussian[49] | 39.29 | 31.01 | **60.44** | 22.70 | 38.36 | - | - | - | - | - | 1.0h |
| Ours | **49.73** | **31.21** | 58.01 | **32.15** | **42.78** | 47.97 | 27.90 | 52.96 | 23.93 | 38.19 | 0.7h |

Table 4. Ablation study of indoor surface reconstruction without the sensor depth. The best metrics are **highlighted**. The red box means that the loss function is designed for 3D open-vocabulary segmentation.

| Settings | Accuracy↓ | Normal Consistency↑ | F-score↑ |
|---|---|---|---|
| No $\mathcal{L}_n$ | 0.1850 | 0.6893 | 0.2322 |
| No $\mathcal{L}_d$ | 0.1077 | 0.8023 | 0.4440 |
| No $\mathcal{L}_s$ | 0.1303 | 0.7910 | 0.4356 |
| No $\mathcal{L}_{clip}$ | 0.0966 | 0.8393 | 0.4614 |
| No $\mathcal{L}_m$ | 0.0951 | 0.8425 | 0.4671 |
| All | **0.0814** | **0.8474** | **0.5127** |

Table 5. Ablation study of 3D open-vocabulary segmentation. The best metrics are **highlighted**. The red box means that the loss function is designed for indoor surface reconstruction.

| Settings | mIoU↑ | mBIoU↑ |
|---|---|---|
| No $\mathcal{L}_n$ | 32.90 | 27.61 |
| No $\mathcal{L}_d$ | 31.81 | 27.37 |
| No $\mathcal{L}_s$ | 33.98 | 29.08 |
| No $\mathcal{L}_{clip}$ | 29.14 | 25.01 |
| No $\mathcal{L}_m$ | 28.71 | 22.62 |
| All | **42.78** | **38.19** |

Table 6. Quantitative results of novel view synthesis on Scan-Net++ and LERF-OVS dataset. The best metrics are **highlighted**.

| Dataset | Methods | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|---|
| ScanNet++ | DN-Splatter[46] | **20.73** | 0.8476 | 0.1996 |
| | 2DGS[24] | 15.60 | 0.7791 | 0.2581 |
| | PGSR[10] | 19.47 | 0.8275 | 0.2060 |
| | Ours | 20.69 | **0.8496** | **0.1709** |
| LERF-OVS | OpenGaussian[49] | 23.78 | 0.8508 | **0.2295** |
| | Ours | **24.13** | **0.8537** | 0.2427 |

rexonstruction and segmentation results.

## 5.5. Novel view synthesis

Finally, we evaluate the performance of our model in the task of novel view synthesis. The results are presented in Tab. 6. Our model achieves comparable performance with DN-Splatter, and better performance against other geometry-aware approaches in the field of indoor surface reconstruction. We also compare our method with Open-Gaussian in the field of 3D open-vocabulary segmentation, our method renders more accurate novel views than it.

## 6. Conclusion

In this work, we present GLS, a novel 3DGS-based framework that effectively combines indoor surface reconstruction and 3D open-vocabulary segmentation. We propose leveraging 2D geometric and semantic cues to optimize the performance of 3DGS on two tasks jointly. We design two novel regularization terms to enhance the sharpness and smoothness of the scene surface, and then improve the segmentation quality. Comprehensive experiments on both 3D open-vocabulary segmentation and indoor surface reconstruction tasks illustrate that GLS outperforms state-of-the-art methods quantitatively and qualitatively. Besides, the ablation study explores the effectiveness of each regularization term on two tasks.

**Limitation.** The proposed method follows the natural limitation of TSDF fusion, whose completeness relies on the number of captured views. Our model generates empty geometry of selected objects in unseen views of the scene. Introducing image-to-3D models to supply unseen information can help solve this issue.



Figure 8. Qualitative comparisons of 3D open-vocabulary segmentation on LERF-OVS dataset, between state-of-the-art methods (OpenGaussian[49] and Gaussian grouping[52]) and our model.

**Effect of** $\mathcal{L}_m$. We further explore the effect of $\mathcal{L}_m$ by dropping the segmentation results of DEVA as supervision. As Tab. 5 reports, $\mathcal{L}_m$ faithfully increases the sharpness of

8

# References

[1] *nerfview: a minimal* web viewer for interactive NeRF rendering*, 2024. 12

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 12

[3] Gwangbin Bae and Andrew J. Davison. Rethinking inductive biases for surface normal estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1, 2, 3, 4, 12

[4] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 2

[5] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022.

[6] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023. 2

[7] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 2

[8] Kazii Botashev, Vladislav Pyatov, Gonzalo Ferrer, and Stamatios Lefkimmiatis. Gsloc: Visual localization with 3d gaussian splatting. *arXiv preprint arXiv:2410.06165*, 2024. 1

[9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3

[10] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. Pgsr: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *arXiv preprint arXiv:2406.06521*, 2024. 1, 2, 4, 5, 6, 8, 11, 12

[11] Hanlin Chen, Fangyin Wei, Chen Li, Tianxin Huang, Yunsong Wang, and Gim Hee Lee. Vcr-gaus: View consistent depth-normal regularizer for gaussian surface reconstruction. *arXiv preprint arXiv:2406.05774*, 2024. 2, 12, 14

[12] Lin-Zhuo Chen, Kangjie Liu, Youtian Lin, Zhihao Li, Siyu Zhu, Xun Cao, and Yao Yao. Flow distillation sampling: Regularizing 3d gaussians with pre-trained matching priors. In *ICLR*, 2025. 2, 5, 6, 11

[13] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16569–16578, 2023. 2

[14] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *ICCV*, 2023. 1, 3, 4

[15] Kai Cheng, Xiaoxiao Long, Kaizhi Yang, Yao Yao, Wei Yin, Yuexin Ma, Wenping Wang, and Xuejin Chen. Gaussianpro: 3d gaussian splatting with progressive propagation. In *Forty-first International Conference on Machine Learning*, 2024. 4

[16] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 12

[17] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5501–5510, 2022. 2

[18] Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuexin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and Xiaoxiao Long. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image. In *European Conference on Computer Vision*, pages 241–258. Springer, 2025. 2

[19] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009. 2

[20] Clement Godard, Peter Hedman, Wenbin Li, and Gabriel J Brostow. Multi-view reconstruction of highly specular surfaces in uncontrolled environments. In *2015 International Conference on 3D Vision*, pages 19–27. IEEE, 2015.

[21] Shubham Goel, Georgia Gkioxari, and Jitendra Malik. Differentiable stereopsis: Meshes from multiple views using differentiable rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8635–8644, 2022. 2

[22] Antoine Guédon and Vincent Lepetit. 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*, pages 5354–5363, 2024. 1, 2

[23] Jun Guo, Xiaojian Ma, Yue Fan, Huaping Liu, and Qing Li. Semantic gaussians: Open-vocabulary scene understanding with 3d gaussian splatting. *arXiv preprint arXiv:2403.15624*, 2024. 1

[24] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *SIGGRAPH 2024 Conference Papers*. Association for Computing Machinery, 2024. 1, 2, 4, 5, 6, 8, 12, 13

[25] Tianyu Huang, Bowen Dong, Yunhan Yang, Xiaoshui Huang, Rynson WH Lau, Wanli Ouyang, and Wangmeng Zuo. Clip2point: Transfer clip to point cloud classification with image-depth pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22157–22167, 2023. 2

[26] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Open-

clip, 2021. If you use this software, please cite it as below. 1, 3, 4

[27] Ying Jiang, Chang Yu, Tianyi Xie, Xuan Li, Yutao Feng, Huamin Wang, Minchen Li, Henry Lau, Feng Gao, Yin Yang, et al. Vr-gs: a physical dynamics-aware interactive gaussian splatting system in virtual reality. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–1, 2024. 1

[28] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 1, 2

[29] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42 (4), 2023. 1, 2, 3, 5

[30] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 2, 5, 15, 16

[31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 1, 3, 4

[32] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 11

[33] Zhihao Liang, Qi Zhang, Wenbo Hu, Ying Feng, Lei Zhu, and Kui Jia. Analytic-splatting: Anti-aliased 3d gaussian splatting via analytic integration. *arXiv preprint arXiv:2403.11056*, 2024. 2

[34] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Seminal graphics: pioneering efforts that shaped the field*, pages 347–353. 1998. 4, 5

[35] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 3

[36] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. 2

[37] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE international symposium on mixed and augmented reality*, pages 127–136. Ieee, 2011. 2, 4, 5, 14

[38] Vaishakh Patil and Marco Hutter. Radiance fields for robotic teleoperation. *arXiv preprint arXiv:2407.20194*, 2024. 1

[39] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20051–20060, 2024. 1, 2, 3, 5, 8, 14

[40] Christian Reiser, Rick Szeliski, Dor Verbin, Pratul Srinivasan, Ben Mildenhall, Andreas Geiger, Jon Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *ACM Transactions on Graphics (TOG)*, 42(4):1–12, 2023. 2

[41] Xuqian Ren, Wenjia Wang, Dingding Cai, Tuuli Tuominen, Juho Kannala, and Esa Rahtu. Mushroom: Multi-sensor hybrid room dataset for joint 3d reconstruction and novel view synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4508–4517, 2024. 2, 5, 14, 15

[42] Xuqian Ren, Matias Turkulainen, Jiepeng Wang, Otto Seiskari, Iaroslav Melekhov, Juho Kannala, and Esa Rahtu. Ags-mesh: Adaptive gaussian splatting and meshing with geometric priors for indoor room reconstruction using smartphones. In *International Conference on 3D Vision (3DV)*, 2025. 2, 11

[43] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pages 501–518. Springer, 2016. 2, 3

[44] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5333–5343, 2024. 1, 3

[45] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5459–5469, 2022. 2

[46] Matias Turkulainen, Xuqian Ren, Iaroslav Melekhov, Otto Seiskari, Esa Rahtu, and Juho Kannala. Dn-splatter: Depth and normal priors for gaussian splatting and meshing, 2024. 2, 4, 5, 6, 8, 11

[47] Qianyi Wu, Xian Liu, Yuedong Chen, Kejie Li, Chuanxia Zheng, Jianfei Cai, and Jianmin Zheng. Object-compositional neural implicit surfaces. In *European Conference on Computer Vision*, 2022. 2

[48] Qianyi Wu, Kaisiyuan Wang, Kejie Li, Jianmin Zheng, and Jianfei Cai. Objectsdf++: Improved object-compositional neural implicit surfaces. In *ICCV*, 2023. 2

[49] Yanmin Wu, Jiarui Meng, Haijie Li, Chenming Wu, Yahao Shi, Xinhua Cheng, Chen Zhao, Haocheng Feng, Errui Ding, Jingdong Wang, et al. Opengaussian: Towards point-level 3d gaussian-based open vocabulary understanding. *arXiv preprint arXiv:2406.02058*, 2024. 1, 3, 4, 5, 6, 8, 11, 14

[50] Le Xue, Ning Yu, Shu Zhang, Artemis Panagopoulou, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, et al. Ulip-2: Towards scalable multimodal pre-training for 3d understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27091–27101, 2024. 2

[51] Hongji Yang, Jiao Liu, Shaoping Lu, and Bo Ren. Self-supervised implicit 3d reconstruction via rgb-d scans. In

*2023 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1115–1120. IEEE, 2023. 2

[52] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *ECCV*, 2024. 1, 3, 5, 6, 8

[53] Zongxin Ye, Wenyu Li, Sidun Liu, Peng Qiao, and Yong Dou. Absgs: Recovering fine details in 3d gaussian splatting. In *ACM Multimedia 2024*, 2024. 5, 14

[54] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–22, 2023. 2, 5, 14, 15

[55] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. 2

[56] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2

[57] Zehao Yu, Torsten Sattler, and Andreas Geiger. Gaussian opacity fields: Efficient adaptive surface reconstruction in unbounded scenes. *ACM Transactions on Graphics*, 2024. 1, 2

[58] Baowen Zhang, Chuan Fang, Rakesh Shrestha, Yixun Liang, Xiaoxiao Long, and Ping Tan. Rade-gs: Rasterizing depth in gaussian splatting. *arXiv preprint arXiv:2406.01467*, 2024.

[59] Wenyuan Zhang, Yu-Shen Liu, and Zhizhong Han. Neural signed distance function inference through splatting 3d gaussians pulled on zero-level set. In *Advances in Neural Information Processing Systems*, 2024. 1

[60] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21676–21685, 2024. 1, 3

[61] Xiangyang Zhu, Renrui Zhang, Bowei He, Ziyu Guo, Ziyao Zeng, Zipeng Qin, Shanghang Zhang, and Peng Gao. Pointclip v2: Prompting clip and gpt for powerful 3d open-world learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2639–2650, 2023. 2

Unlike recent methods of indoor surface reconstruction [12, 42, 46] and 3D open-vocabulary segmentation [49] only concentrate on optimizing one task, our goal is to strengthen two tasks via our joint optimization scheme. Here, we provide more details of our method and experiments. Specifically, we provide analysis through a Q&A section to answer four important questions (Sec. A), additional algorithm details (Sec. B), additional details of experiments (Sec. C), additional results (Sec. D) and the future work (Sec. E).

## A. Q&A

Q1. *Please explain the benefits of the proposed joint optimization.*

A1. We have analyzed the benefits of the proposed joint optimization of indoor surface reconstruction and 3D open-vocabulary segmentation from Line_037 to Line_056 of the manuscript. Theoretically, the optimization goals of the two tasks are the **smoothness** and **sharpness** of generated results (including the surface mesh and the segmentation mask). Then we can optimize two tasks simultaneously through the gradient descent strategy. The marked metrics in Table 4 and Table 5 of the manuscript demonstrate the effectiveness of the proposed joint optimization in two fields. The qualitative results of the ablation study are shown in Fig. 11. The performance of GLS degenerates reasonably when we disable a regularization term and enable others.

Q2. *Why does GLS only focus on indoor scenes?*

A2. Considering the unique challenges of indoor scenes and the limitation of geometric and semantic cues in outdoor scenes, there are two main reasons as follows:

   i. For surface reconstruction, we follow DN-Splatter [46] and FDS [12], then focus on achieving high-fidelity reconstruction in indoor scenes. Unlike outdoor scenes (e.g., the Tanks and Temples dataset [32]), indoor scenes suffer from more challenging conditions. On the one hand, indoor scenes are always captured by free camera trajectories, which causes the imbalanced allocation of spatial representation capacity. On the other hand, there are a lot of texture-less regions (such as white walls and floors) in indoor scenes. These issues cause extreme ambiguity of multi-view reconstruction and large performance degeneration of the state-of-the-art method PGSR [10]. However, GLS tackles these issues well because of the robust joint optimization.

   ii. As Fig. 14 shows, geometric cues predicted from the pre-trained monocular normal model, suffer from inconsistent noises across different views in outdoor scenes. The semantic cues only consist of a few instances and encode large ambiguity of object structures. These issues are harmful to the joint optimization of GLS.

   Hence, we focus on the surface indoor scenes where the quality of geometric cues is more stable and semantic distribution is highly predictable. To illustrate the generalization ability of GLS, We also evaluate the performance of our model in a large-scale indoor scene ('Meetingroom') of the Tanks and Temples dataset, the result is presented in Fig. 10. Our method successfully handles the high-light and dark regions and reconstructs a smoother scene surface than PGSR [10].

Q3. *How about the robustness of GLS under weak semantic priors?*

"I want to make a piece of toast"    "I want to cut a tomato"    "I want to wash my hands"    "I want to store my ice cream"

Figure 9. A demo of **GLS + nerfview + GPT-4V**. Our tool can find the object that solves the user's request (bottom) and extract the scene geometry.
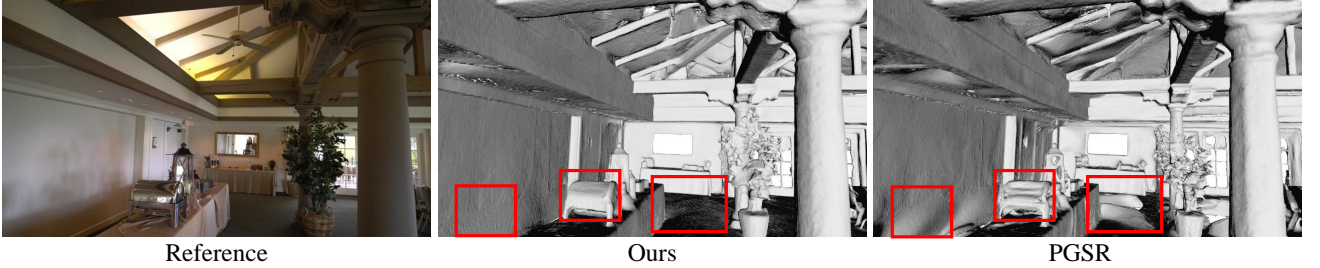


Reference    Ours    PGSR

Figure 10. Qualitative results in a large-scale indoor scene. Our method successfully handles the high-light and dark regions.

A3. Most view-consistent segmentation priors we used on the MuSHRoom dataset are noisy or incomplete, such as Fig. 13 shows. It demonstrates that our method is certainly robust. We also evaluate the reconstruction result under different quality of CLIP features, by disabling the denoising procedure of pre-processing to obtain the low-quality CLIP features. The result is shown in Fig. 12. High-quality CLIP features reduce shadow interference. This result also illustrates the effectiveness of the proposed joint optimization.

Q4. *How about applications of GLS?*

A4. Thanks to the high-quality surface reconstruction and segmentation results of GLS, we develop several tools to explore applications of GLS:

  i. **GLS + Blender**: With the help of Blender [16], we can achieve the scene-editing effects, including moving the target object and adding objects.

 ii. **GLS + nerfview**: We develop an interactive tool based on nerfview [1], to show the scene geometry (depth and normal) and open-vocabulary attention in free views.

iii. **GLS + nerfview + GPT-4V**: As Fig. 9 shows, we also add the GPT-4V [2] to 'ii' to achieve intelligent interaction. Given a user's request, our tool can find the object that solves the request and extract the scene geometry used to help the user plan the path to the target object. We believe this tool has potential value in the field of embodied intelligence.

Our supplementary videos consist of two application demos ('application_i.mp4' and 'application_ii.mp4'). We

will release all tools when this paper is accepted.

## B. Additional Algorithm Details

**About $N_d$.** We follow the manner of PGSR [10] to estimate $N_d$. Given a pixel point p and its four neighboring pixels, we unproject these 2D points into 3D points $\{P_i | i = 0, ..., 4\}$ by $D_p$, then calculate the local normal $N_d$ of p via:

$$N_d(p) = \frac{(P_1 - P_0) \times (P_3 - P_2)}{|(P_1 - P_0) \times (P_3 - P_2)|} \quad (10)$$

**About $\mathcal{L}_n$.** Inspired by the 2DGS [24] and considering semi-transparent surfels, we adaptively give high weights to actual surfaces in $\mathcal{L}_n$ by $A$.

**About $\mathcal{L}_s$.** $\mathcal{L}_s$ is a powerful regularization term for $N_d$ and easily causes over-smoothing effect for small objects. In practice, we adaptively select big objects that occupy top-3 areas in a view by $M_o$. Then we adopt $\mathcal{L}_s$ to constrain them.

**About $D_r$ and $\mathcal{L}_d$.** Essentially, $D_r$ is a refined depth based on the probability rather than an actual depth. It is built on the error analysis between the rendered normal and the ideal normal. We only consider the normal prior [3] as the reference normal, because its quality is unstable as mentioned by VCR-GauS [11]. Hence, the ideal normal cannot be acquired and the error analysis is necessary. We adopt a total of six loss functions during training, to balance the value of each loss function and avoid the abrupt value of $D_r$, we propose the function $y = 1 - e^{(-|x|)}$ to generate $\mathcal{L}_d$.
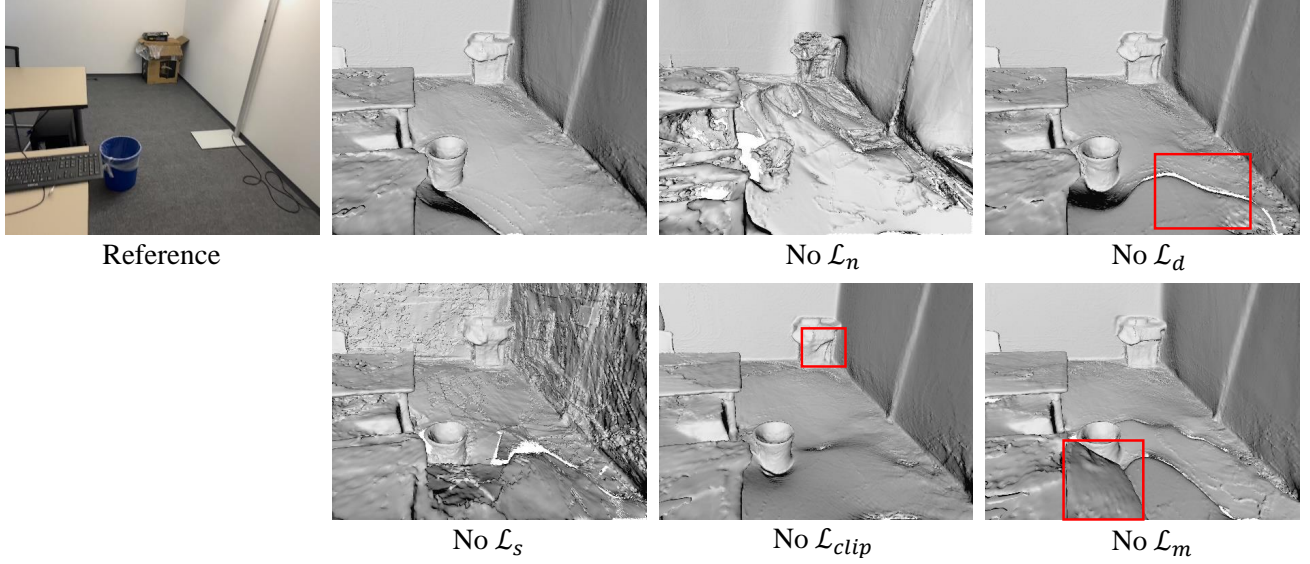
12

Figure 11. Qualitative results of the ablation study on 'b20a261fdf' of ScanNet++ dataset. The performance of GLS degenerates reasonably when we disable a regularization term and enable others.
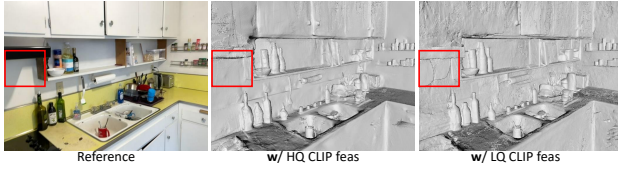


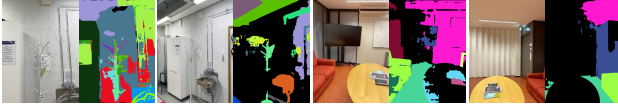Figure 12. Effect of different quality of CLIP features.



Figure 13. Example of the view-consistent segmentation results of DEVA [11]. We directly use these results as the supervision of the rendered mask.

## C. Additional Experimental Details

Tab. 8 presents metrics of datasets used in our experiments. We list additional optimization hyperparameters below:

```
def __init__(self, parser):
    self.iterations = 30_000
    self.position_lr_init = 0.00016
    self.position_lr_final = 0.0000016
    self.position_lr_delay_mult = 0.01
    self.position_lr_max_steps = 30_000
    self.feature_lr = 0.0025
    self.opacity_lr = 0.05
    self.scaling_lr = 0.005
```
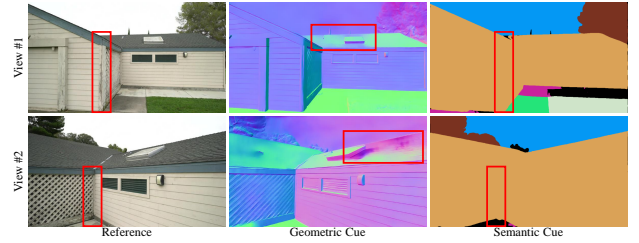


Figure 14. Limitation of geometric and semantic cues in outdoor scenes. The geometric cues predicted from the pretrained monocular normal model, suffer from inconsistent noises across different views. The semantic cues encode large ambiguity of object structures.
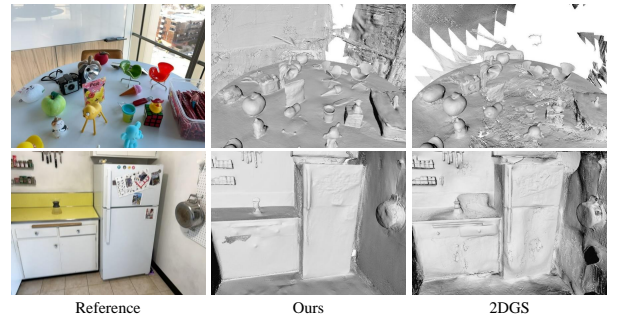


Figure 15. Results of indoor surface reconstruction on LERF-OVS dataset. Our model reconstructs sharper and smoother indoor surfaces than 2DGS [24].

```
self.rotation_lr = 0.001
self.percent_dense = 0.01
self.lambda_dssim = 0.2
self.densification_interval = 100
self.opacity_reset_interval = 3000
self.densify_from_iter = 500
self.densify_until_iter = 15_000
self.densify_grad_threshold = 0.0006
self.opacity_cull_threshold = 0.05
self.densify_abs_grad_threshold = 0.0008
self.abs_split_radii2D_threshold = 20
self.max_abs_split_points = 50_000
self.max_all_points = 6000_000
```

For Open-vocabulary Segmentation, we follow LangSplat [39] to decrease the last dimension of original CLIP features from 512 to 16 by the encoder part of an encoder-decoder network. Then we use the decoder part of the same network to increase the the last dimension of Gaussian semantic features from 16 to 512, to compute the relevancy between them and the original CLIP features. The extraction scheme for SAM masks and CLIP features is also aligned with LangSplat [39], while we follow OpenGaussian [49] only to extract the large layer of SAM masks. The learning rate of the MLP layer is 0.00005.

For indoor surface reconstruction, we use the iPhone sequences with COLMAP registered poses on both MuSH-Room [41] and ScanNet++ [54] datasets. Besides, we adopt the densification strategy of AbsGS [53]. We disable the exposure compensation strategy of PGSR to maintain fair comparison.

## D. Additional results

Per-scene quantitative results of GLS on the MuSHRoom and ScanNet++ datasets are reported in Tab. 7. As Fig. 16 and Fig. 17 show, GLS also can attach the semantic information to the resconstructed mesh, by replacing the color image with the encoded semantic mask during TSDF fusion [37]. Please note that our model is trained without any manual semantic annotations. Fig. 18 shows more open-vocabulary segmentation results of GLS. Our model can accurately segment target objects selected by the corresponding text.

Fig. 15 shows additional reconstruction results of our model and 2DGS on LERF-OVS dataset. Our model not only produces smooth and sharp object surfaces, but also successfully reconstructs the surfaces of highly- reflective whiteboard and transparent glass.

## E. Future Work

In the future, we have two plans:
i. To address the limitation of TSDF Fusion, we can first segment object mesh based on GLS to supply the ob-

ject scale and orientation, then leverage the image-to-3D models to obtain the whole object body from the segmented object appearance.

ii. To extend GLS to outdoor scenes, we can follow VCR-GauS [11] to estimate the confidence of geometric and semantic cues to solve their issues.

Table 7. Per-scene quantitative results of GLS on the MuSHRoom and ScanNet++ datasets.

| | Scenes | Accuracy ↓ | Completion ↓ | Chamfer−$L_1$ ↓ | Normal Consistency ↑ | F-score ↑ |
|---|---|---|---|---|---|---|
| | coffee_room | 0.0231 | 0.0275 | 0.0253 | 0.8702 | 0.9020 |
| | computer | 0.0394 | 0.0277 | 0.0336 | 0.8756 | 0.8445 |
| | honka | 0.0264 | 0.0284 | 0.0274 | 0.8742 | 0.9090 |
| w/ Sensor Depth | kokko | 0.0305 | 0.0272 | 0.0444 | 0.9064 | 0.8623 |
| | vr_room | 0.0244 | 0.0237 | 0.0241 | 0.8885 | 0.8802 |
| | 8b5caf3398 | 0.0936 | 0.0241 | 0.0588 | 0.8657 | 0.8741 |
| | b20a261fdf | 0.0335 | 0.0270 | 0.0303 | 0.9219 | 0.8841 |
| | coffee_room | 0.0684 | 0.0669 | 0.0676 | 0.7992 | 0.6103 |
| | computer | 0.0963 | 0.0874 | 0.0918 | 0.7708 | 0.4739 |
| | honka | 0.0810 | 0.0750 | 0.0780 | 0.7899 | 0.5503 |
| w/o Sensor Depth | kokko | 0.1102 | 0.1041 | 0.1072 | 0.7515 | 0.3969 |
| | vr_room | 0.0956 | 0.0852 | 0.0904 | 0.8035 | 0.5453 |
| | 8b5caf3398 | 0.0618 | 0.0580 | 0.0599 | 0.8732 | 0.6045 |
| | b20a261fdf | 0.1103 | 0.1438 | 0.1270 | 0.8423 | 0.3553 |

Table 8. Metrics of datasets used in our experiments.

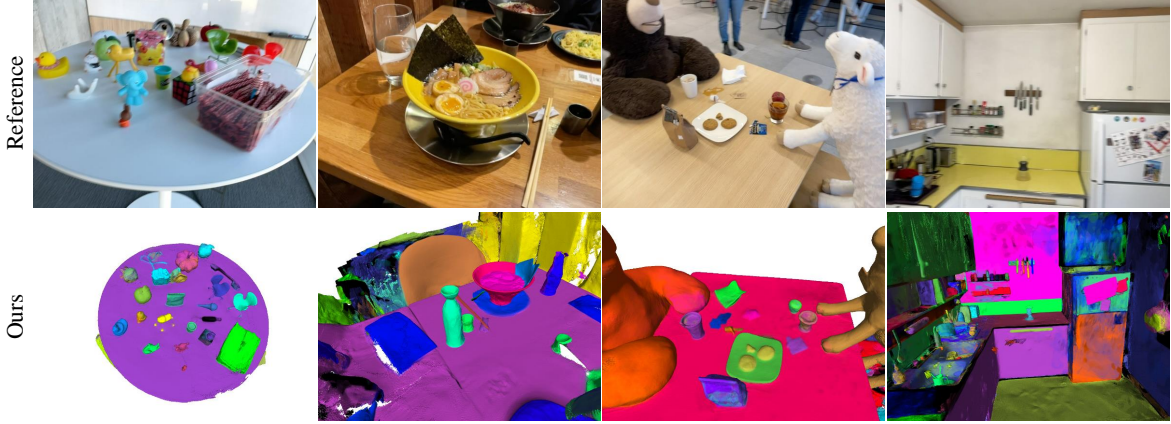| Dataset | LERF-OVS [30] | | | | MuSHRoom [41] | | | | | ScanNet++ [54] | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Scene | figurines | ramen | teatime | waldo_kitchen | coffee_room | computer | honka | kokko | vr_room | 8b5caf3398 | b20a261fdf |
| Resolution | 986 × 728 | 988 × 731 | 988 × 730 | 985 × 725 | 738 × 994 | | | | | 1920 × 1440 | |
| Training Views | 299 | 131 | 177 | 187 | 353 | 455 | 320 | 348 | 418 | 126 | 59 |
| Initial points | 65k | 27k | 23k | 15k | 1000k | | | | | 111k | 111k |



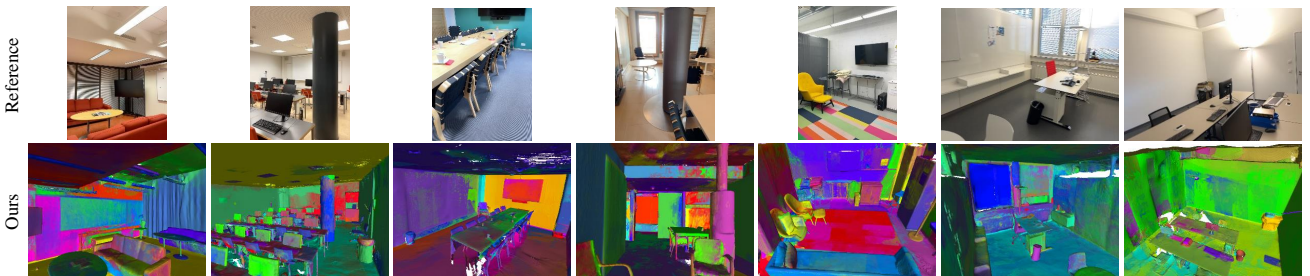Figure 16. Semantic mesh results of GLS on LERF-OVS [30].



Figure 17. Semantic mesh results of GLS on MuSHRoom [41] and ScanNet++ [54].
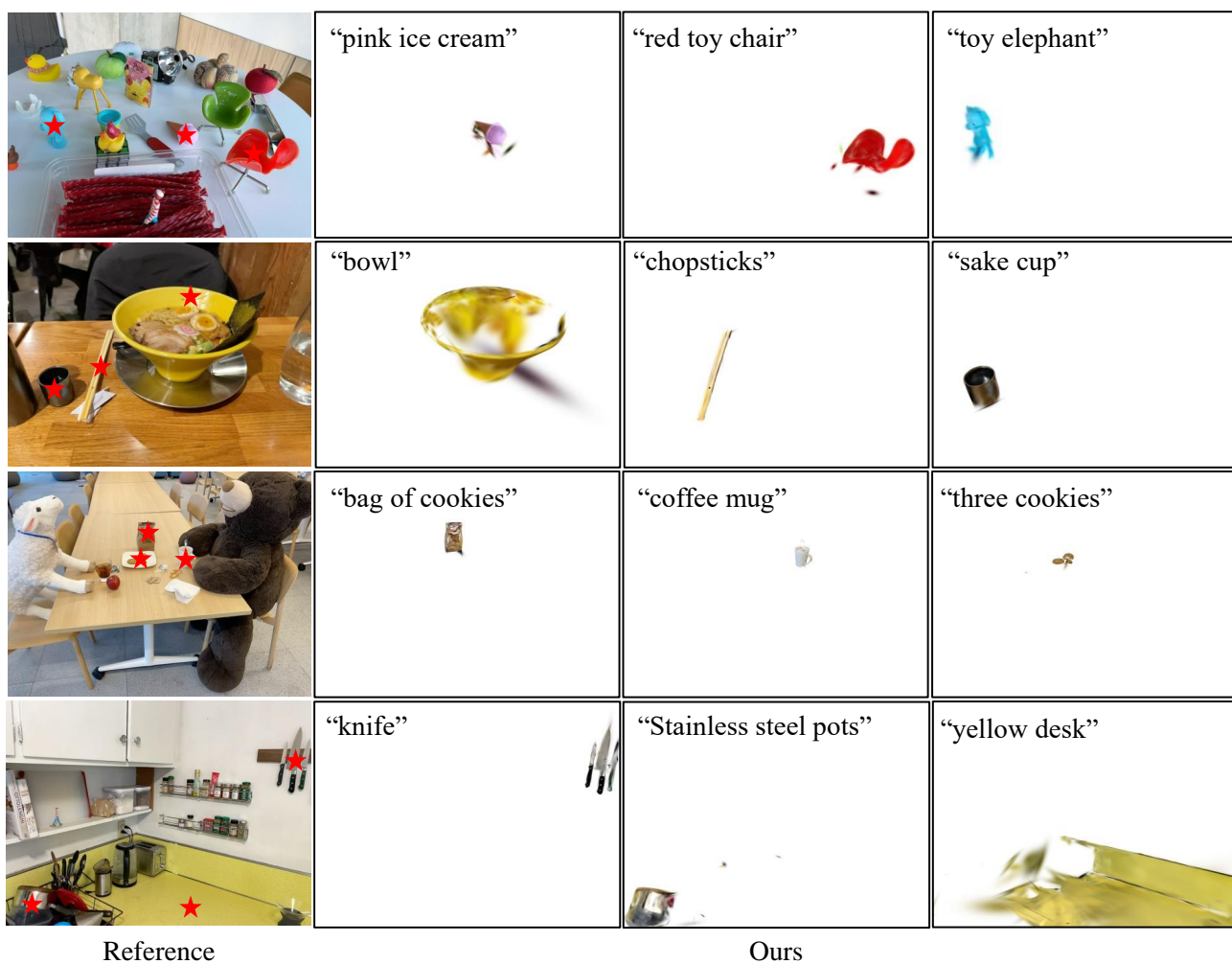
15

Figure 18. More open-vocabulary segmentation results of GLS on LERF-OVS [30].