

Rethinking Open-Vocabulary Segmentation of Radiance Fields in 3D Space

Hyunjee Lee*, Youngsik Yun*, Jeongmin Bae, Seoha Kim, Youngjung Uh†

Yonsei University
 {hyunji12, bbangsik, jaymin.bae, hailey07, yj.uh}@yonsei.ac.kr

Abstract

Understanding the 3D semantics of a scene is a fundamental problem for various scenarios such as embodied agents. While NeRFs and 3DGS excel at novel-view synthesis, previous methods for understanding their semantics have been limited to incomplete 3D understanding: their segmentation results are 2D masks and their supervision is anchored at 2D pixels. This paper revisits the problem set to pursue a better 3D understanding of a scene modeled by NeRFs and 3DGS as follows. 1) We directly supervise the 3D points to train the language embedding field. It achieves state-of-the-art accuracy without relying on multi-scale language embeddings. 2) We transfer the pre-trained language field to 3DGS, achieving the first real-time rendering speed without sacrificing training time or accuracy. 3) We introduce a 3D querying and evaluation protocol for assessing the reconstructed geometry and semantics together. Code, checkpoints, and annotations will be available online.

Project page: <https://hyunji12.github.io/Open3DRF>

1 Introduction

Semantically understanding 3D space is important for various computer vision tasks. For instance, it is crucial to segment 3D objects accurately for robot manipulation (Rashid et al. 2023; Zheng et al. 2024). Recently, several works have focused on understanding 3D scenes represented by radiance fields such as Neural Radiance Fields (NeRFs) and 3D Gaussian Splatting (3DGS). LERF (Kerr et al. 2023) locates objects with open-vocabulary querying in a given viewpoint by introducing a language embedding field. The language field is supervised by CLIP embeddings from the multi-scale patches to capture various sizes of objects. Subsequent works (Zhang, Li, and Ahuja 2024; Qin et al. 2023) incorporate SAM masks (Kirillov et al. 2023) to supervise the language field for clear segmentation boundaries.

We revisit the 3D understanding of NeRFs and 3DGS in four aspects: problem setting, supervision, embeddings, and evaluation. The problem setting of previous works leads to limited 3D understanding as they merely produce 2D masks for given camera viewpoints rather than 3D volume. On the other hand, we set the problem to segment 3D volume regarding the semantics of 3D points.

*These authors contributed equally.

†Corresponding author.

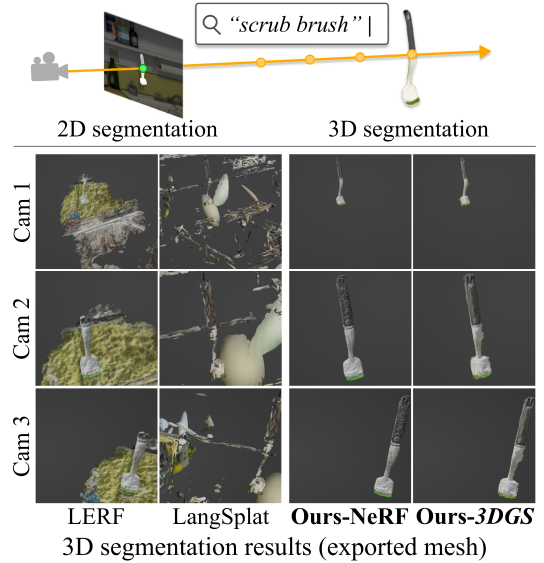


Figure 1: **Teaser.** Previous works segment 2D masks on rendered images to understand radiance fields. Instead, we reformulate the task to segment 3D volumes. Our approach greatly improves 3D and 2D understanding of radiance fields.

The previous methods encourage the rendered embeddings, rather than the embeddings in 3D space, to match the ground truth. On top of that, multi-scale language embeddings require finding the optimal scale for a given viewpoint and a query text, leading to view inconsistency. By contrast, we encourage 3D points to learn language embeddings following our revisited problem setting. Although the ground truth is still language embeddings of 2D images, changing the dimension for computing the loss greatly improves 3D understanding as shown in Figure 1. Furthermore, we remove multi-scale embeddings by encoding masked objects, achieving consistent 2D understanding between viewpoints.

In addition, previous works on understanding 3DGS explicitly add language embeddings for each Gaussian and jointly train them (Zhang, Li, and Ahuja 2024; Qin et al. 2023). However, directly appending a 512-dimensional language embedding to each Gaussian and rendering them ex-

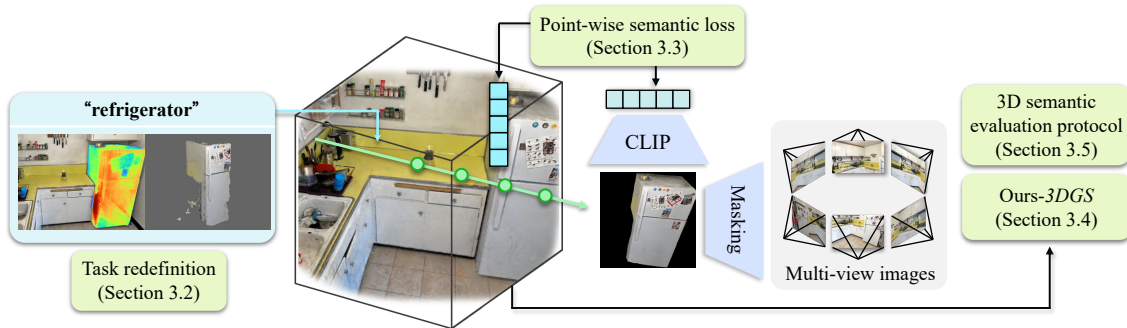


Figure 2: **Overview.** We propose 3D segmentation as a more practical problem setting, segmenting the 3D volume for a given text query (Section 3.2). Then we propose point-wise semantic loss to supervise the sampled point embeddings (Section 3.3). Furthermore, the learned language fields can be transferred into 3DGS for faster rendering speeds (Section 3.4). Lastly, our 3D evaluation protocol measures the 3D segmentation performance both in reconstructed geometry and semantics (Section 3.5).

ceeds the maximum GPU shared memory by $10\times$. Hence, they either 1) compress language embeddings to low-dimensional features, significantly sacrificing the accuracy (Shi et al. 2023; Qin et al. 2023), or 2) utilize GPU global memory which is slow for training (Zhou et al. 2023; Ye and Kanazawa 2023). To overcome these drawbacks, we transfer our pre-trained language field into 3DGS to enable the first real-time rendering of the accurate language field.

Our final aspect is evaluation. Existing works only report their accuracy in mIoU for given viewpoints as in 2D segmentation, or measure mIoU on the *ground truth* point clouds instead of *reconstructed* geometry (Engelmann et al. 2024). However, 3D understanding should know the correct 3D volume of a target semantics. Therefore, the evaluation method should assess semantics and the reconstructed geometry together. Accordingly, we propose to evaluate the accuracy of 3D understanding as the agreement between estimated volume in mesh and ground truth mesh, measured in F1-score. The proposed evaluation is applicable to both NeRF and 3DGS.

In the experiments, we demonstrate our superiority regarding 1) 3D and 2D segmentation accuracy, 2) training and rendering time, and 3) consistency across viewpoints.

In short, our contributions are:

- We propose a practical problem setting for 3D understanding of NeRFs and 3DGS.
- We propose to directly supervise 3D points before volume rendering to learn a language embedding field. It achieves the state-of-the-art accuracy in 3D and 2D segmentation.
- We propose to transfer the language field to 3DGS. It achieves the first real-time rendering of semantics which is $28\times$ faster than the previous fastest method.
- We propose a 3D evaluation protocol between estimated volume and ground truth volume represented as meshes.

2 Related Work

In this section, we briefly review neural 3D scene representations, especially NeRF and 3DGS, and then discuss the

semantic understanding of 3D scenes.

Radiance Fields NeRF (Mildenhall et al. 2020) reconstructs a scene as a continuous function that maps 3D points to radiance and density values. Volume rendering pipeline renders the points to determine the color of each pixel on an image from an arbitrary perspective. Rather than volumetric rendering, 3DGS (Kerbl et al. 2023) achieves fast rendering speeds by projecting 3D Gaussians onto the camera plane followed by depth-sorted alpha-blending. Each 3D Gaussian stores location, rotation, scale, opacity, and color. Recent studies use these representations to semantically understand 3D scenes for various applications such as robotics (Rashid et al. 2023; Wang et al. 2023; Zheng et al. 2024; van Oort et al. 2024; Li and Pathak 2024).

Semantic Understanding of Radiance Field The most straightforward way to understand 3D scenes represented by neural radiance fields is by adding an auxiliary branch for semantic segmentation (Zhi et al. 2021; Bing, Chen, and Yang 2022; Liu et al. 2023; Ye et al. 2023). This approach allows synthesizing semantic masks from novel views but requires a pre-defined list of target classes before training, called *closed-set*. Therefore, retraining or using additional models for untrained queries is necessary, limiting application for open-vocabulary scenarios. As vision-language model (VLM) features (*e.g.*, CLIP (Radford et al. 2021)) expand the semantic understanding to *open-set*, it becomes a widely used approach to distill CLIP features into 3D scenes.

LERF (Kerr et al. 2023) pre-computes multi-scale patches to prepare multi-scale ground-truth CLIP features. Instead of using patches, LEGaussians (Shi et al. 2023) distills pixel-level CLIP features. Similarly, OpenNeRF (Engelmann et al. 2024) computes pixel-level CLIP features using OpenSeg (Ghiasi et al. 2022). As it produces mixed CLIP features to contain multiple objects into a patch, LangSplat (Qin et al. 2023) utilizes SAM (Kirillov et al. 2023) to create patches for a single object. Since studies using multi-scale feature maps select the most relevant scale, different scales can be chosen at different views with the same query, leading to multi-view inconsistent results (Kerr et al. 2023; Qin et al.

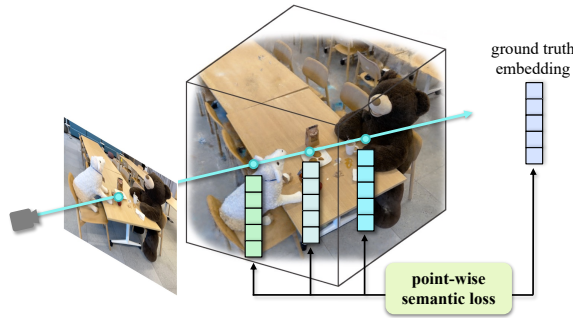


Figure 3: **Point-wise semantic loss** supervises the language embeddings of sampled points directly in 3D space, ensuring precise semantics.

2023; Shi et al. 2023). On the other hand, we use SAM to make the field free from scales.

Language Embedded 3DGS Directly embedding high-dimensional CLIP features into each Gaussian is infeasible due to the limited GPU shared memory. Therefore, recent studies for understanding open-vocabulary 3D segmentation using 3D as a backbone address this problem in two ways. LEGaussians and LangSplat compress the CLIP features using an autoencoder (Shi et al. 2023; Qin et al. 2023). However, they need to optimize the autoencoder for each scene and endure performance degradation. Others using 3DGS without feature compression (Zhou et al. 2023; Bhalgat et al. 2024; Labe et al. 2024) utilize global memory, suffering for longer training time. Meanwhile, we effectively optimize 3DGS transferred from our learned language field, without the above trade-off.

3 Methods

Firstly, we redefine *3D segmentation* and explain how to query 3D and 2D segmentation results. Subsequently, we construct an additional language field to NeRF as in LERF. Then, we introduce point-wise semantic loss to supervise embeddings of sampled ray points in the language field and describe our scale-free embeddings. Additionally, we propose to transfer our pre-trained language field into 3DGS for real-time rendering. Lastly, we introduce the evaluation protocol for 3D segmentation. Figure 2 illustrates an overview.

3.1 Preliminary: LERF

For open-vocabulary segmentation, LERF (Kerr et al. 2023) builds an additional language field on top of iNGP (Müller et al. 2022). The language field F_{lang} is jointly trained with the radiance field by querying point embeddings $F_{\text{lang}}(\mathbf{x}, s)$ at N sample points' position \mathbf{x} and scale s along each ray. The 3D point language embeddings are then accumulated via the volume rendering technique (Max 1995) to obtain rendered language embedding $\hat{\phi}_{\text{lang}}^s$ in 2D pixel space: $\hat{\phi}_{\text{lang}}^s = \sum_{i=0}^N w_i F_{\text{lang}}(\mathbf{x}_i, s)$. The weight of each sampled point is calculated as $w_i = T_i(1 - \exp(-\sigma_i \delta_i))$, where $T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j)$ is transmittance, δ is the dis-

tance between adjacent samples, and σ is the volume density. Then, the rendered embedding is normalized to the unit sphere as in CLIP: $\phi_{\text{lang}}^s = \hat{\phi}_{\text{lang}}^s / \|\hat{\phi}_{\text{lang}}^s\|$. To train the language field, LERF crops the training dataset into multi-scale patches, creating ground truth language embedding $\phi_{\text{lang}}^{gt_s}$ and maximizing the cosine similarity between the rendered language embedding ϕ_{lang}^s :

$$L_{\text{lang}} = - \sum_s \lambda_{\text{lang}} \phi_{\text{lang}}^s \cdot \phi_{\text{lang}}^{gt_s}. \quad (1)$$

Furthermore, LERF builds an additional branch for the DINO (Caron et al. 2021) feature field as an extra regularizer for achieving clearer object boundaries. Similar to the above, the DINO branch is trained to maximize the similarity between rendered DINO embedding $\hat{\phi}_{\text{dino}}$ and the DINO ground truth ϕ_{dino}^{gt} . The DINO branch is not used during inference.

3.2 Task Redefinition

Existing methods have limited 3D semantic understanding as they produce 2D masks for given texts and viewpoints. Instead, we propose *3D segmentation* as a more practical problem setting: segmenting the 3D volume for a given text. While following the basic elements of language embeddings and relevancy scores, we compute relevancy scores of the language embeddings queried on 3D points \mathbf{x} instead of the ones rendered on 2D images.

$$\min_i \frac{\exp(F_{\text{lang}}(\mathbf{x}) \cdot \phi_{\text{text}})}{\exp(F_{\text{lang}}(\mathbf{x}) \cdot \phi_{\text{text}}) + \exp(F_{\text{lang}}(\mathbf{x}) \cdot \phi_{\text{canon}}^i)}, \quad (2)$$

where ϕ_{canon}^i represents predefined canonical texts such as *photo* and *image*. In NeRF representation, we compute the relevancy scores on the points along the rays through pixels. For 3DGS representation, we compute the relevancy scores on the center positions of the Gaussians.

For 2D segmentation, we follow the existing works to compute the relevancy score using the rendered embeddings without scale $\hat{\phi}_{\text{lang}}$.

$$\min_i \frac{\exp(\hat{\phi}_{\text{lang}} \cdot \phi_{\text{text}})}{\exp(\hat{\phi}_{\text{lang}} \cdot \phi_{\text{text}}) + \exp(\hat{\phi}_{\text{lang}} \cdot \phi_{\text{canon}}^i)}. \quad (3)$$

The regions where the computed relevancy scores surpass the selected threshold are considered object regions.

3.3 Supervising Semantics in 3D Space

It is a reasonable choice to minimize the error between the rendered colors and the ground truth colors on 2D images for novel view synthesis. On the other hand, the similar objective for language embeddings (Eq. (1)) harms correctly understanding 3D semantics as shown in Figure 1.

To address this issue, we propose a point-wise semantic loss which uses CLIP embeddings ϕ_{lang}^{gt} as direct ground truth supervision for the embeddings of 3D points on the ray $F_{\text{lang}}(\mathbf{x})$. Specifically, we maximize the similarity between

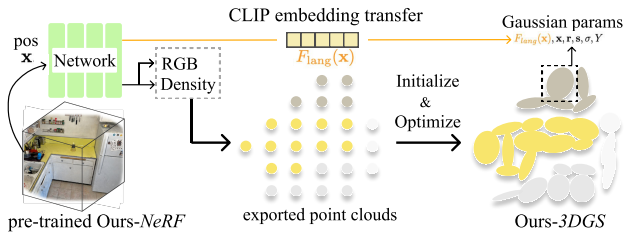


Figure 4: **Transferring Ours-NeRF into 3DGS.** We initialize 3DGS using the point cloud exported from our pre-trained NeRF, then optimize the attributes of 3DGS except for position. The language features obtained by querying the language field at the Gaussian center positions are then transferred to 3DGS.

point embeddings $F_{\text{lang}}(\mathbf{x})$ and the ground truth CLIP embedding $\phi_{\text{lang}}^{\text{gt}}$:

$$L_{\text{PS}} = - \sum_{i=0}^N (w_i F_{\text{lang}}(\mathbf{x}_i) \cdot \phi_{\text{lang}}^{\text{gt}}). \quad (4)$$

We also apply the same approach described above for DINO regularization. We show that this improves not only 3D segmentation performance but also 2D segmentation performance thanks to a better understanding of 3D scenes.

We observe that existing works (Kerr et al. 2023; Qin et al. 2023) on multi-scale CLIP embedding fields require different optimal scales for a text query based on the viewpoint, leading to multi-view inconsistency. To address this, we use SAM (Kirillov et al. 2023) to obtain object masks and create scale-free ground truth CLIP embeddings. This ensures view-consistent language fields without the need to determine the optimal scale. We note that previous methods with SAM still model multi-scale CLIP embedding fields (Zhang, Li, and Ahuja 2024; Qin et al. 2023).

3.4 Transferring Language Field into 3DGS

A straightforward approach for the same task with 3DGS is to jointly optimize the Gaussians and their additional language embeddings, which requires vast memory consumption larger than GPU shared memory. It suffers from slow training ($47\times$) by allocating GPU global memory¹ (Zhou et al. 2023; Bhalgat et al. 2024), or requires an additional model for feature compression, which degrades accuracy due to compression loss (Shi et al. 2023; Qin et al. 2023). Appendix B.1 provides more details.

To tackle this problem, we propose to transfer our pre-trained language field into 3DGS, removing the need for independent optimization. We simply append the language embeddings queried at the center of Gaussians. It runs instantly (22ms) and still allows fast rendering¹.

In addition, we propose to initialize the coordinates and language embeddings of 3DGS from the pre-trained radiance-language field in order to ensure the agreement between the geometry and semantics. We collect ray points

¹Training on GPU global memory slows due to the gradient computation in the backward pass. The forward pass is still fast.

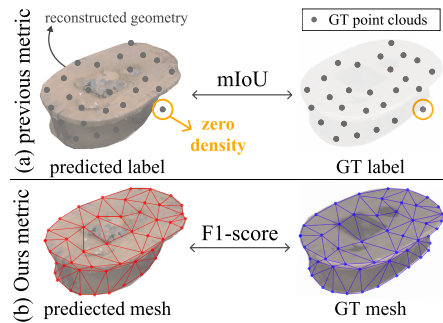


Figure 5: **Comparison of 3D Evaluation.** (a) Existing methods predict the labels at the *ground truth* point cloud. It is misleading when the language embeddings capture the object area while the *reconstructed* geometry does not cover that object area. (b) To address this problem, we extract 3D meshes from the segmented points of the reconstructed scene to measure the F1-score between the exported mesh and the ground truth mesh.

and extract the top 1M points regarding density following NeRFstudio (Tancik et al. 2023). During training, we freeze the center position of the Gaussians \mathbf{x} , without densification or pruning. We then concatenate language embeddings $F_{\text{lang}}(\mathbf{x})$ at the center of Gaussians.

Similar to previous methods, we modify the rasterizer to use global memory instead of limited shared memory to render high-dimensional language features. The rendered language embedding $\hat{\phi}_{\text{lang}}$ is obtained by α -blending:

$$\hat{\phi}_{\text{lang}} = \sum_{i \in N} F_{\text{lang}}(\mathbf{x}_i) \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (5)$$

where N denotes the number of Gaussians overlapping on the pixel. The alpha value is calculated as $\alpha_i = \sigma_i G(\mathbf{x}_i)$, where G denotes the Gaussian kernel and σ denotes the opacity of Gaussians. We use Eq. (2) for 3D segmentation.

3.5 3D Semantic Evaluation Protocol

A previous work (Engelmann et al. 2024) follows the 2D segmentation protocol for 3D segmentation: mIoU on the ground truth point clouds instead of 2D pixels. It might incorrectly classify the points with zero density and correct embedding, as shown in Figure 5(a).

To address this problem, we propose to measure the F1-score between the exported mesh from the segmented results to ground truth mesh, inspired by surface reconstruction literature (Knapitsch et al. 2017; Li et al. 2023). Precision computes the ratio of correct points among the estimated mesh vertices. Recall computes the ratio of covered points among the ground truth mesh vertices. A point is considered *correct* or *covered* if there exists a point within a radius in the counterpart. We note that this protocol can be generally applied to various neural representations, such as NeRF or 3DGS.

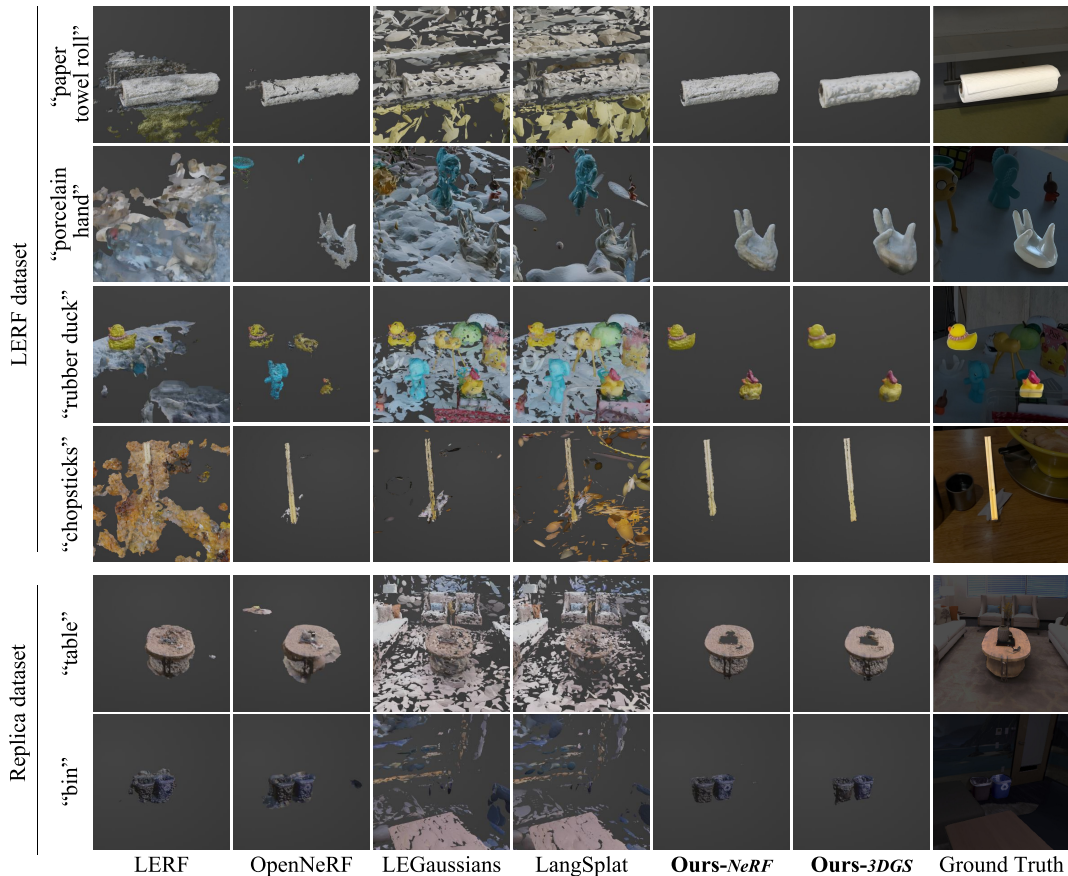


Figure 6: **Qualitative comparisons of 3D segmentation on LERF and Replica datasets.** We show an exported mesh of 3D querying results for the given text query. Unlike competitors, our method produces more clear boundaries in 3D segmentation results.

Method	Replica	LERF	3D-OVS
LERF	0.50	0.50	0.50
OpenNeRF	0.50	0.50	0.50
Ours- <i>NeRF</i>	0.55	0.60	0.55
LEGaussians	0.50	0.50	0.50
LangSplat	0.40	0.40	0.40
Ours- <i>3DGS</i>	0.60	0.60	0.55

Table 1: **Selected Thresholds in Quantitative Results**

4 Experiments

In this section, we evaluate our methods on various datasets and compare them to competitors in terms of 3D and 2D segmentation with given text queries. We choose LERF, OpenNeRF, LEGaussians, and LangSplat as competitors, which are open-sourced. We used the official code without modification. Appendix A.2 provides details of the datasets.

Evaluation For quantitative comparison in 3D segmentation, we export meshes using the Poisson surface reconstruction. In LERF, OpenNeRF, and Ours-*NeRF*, we use Nerfstu-

Backbone	Method	Replica F1-score \uparrow
NeRF	LERF	0.0845
	OpenNeRF	0.0361
	Ours-<i>NeRF</i>	0.1520
3DGS	LEGaussians	0.0067
	LangSplat	0.0087
	Ours-<i>3DGS</i>	0.1353

Table 2: **3D Segmentation Accuracy Comparison on the Replica dataset.** Red and orange highlights indicate the 1st, and 2nd-best model.

dio (Tancik et al. 2023) to export meshes with 30K sampled points. In LEGaussians, LangSplat and Ours-*3DGS*, we use SuGaR (Guédon and Lepetit 2023) to export meshes. For 2D segmentation, we evaluate the mIoU and mAP between the ground truth mask and the predicted mask. For qualitative comparison in 3D segmentation, we use 50K points for Nerfstudio in LERF, OpenNeRF, and Ours-*NeRF*. We use the thresholds in Table 1 for quantitative results.

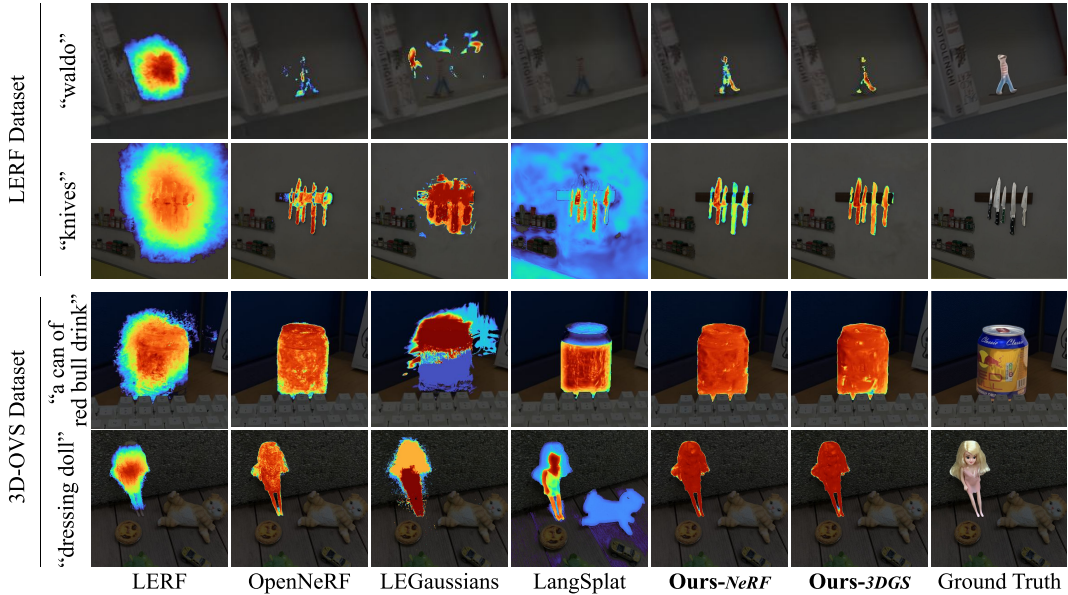


Figure 7: **Comparisons of 2D Segmentation on LERF and 3D-OVS Dataset.** We show a heatmap of the similarity for the given text query. We dim the background except for the target object, for better visualization. Our method achieves accurate segmentation results in 2D compared to competitors.

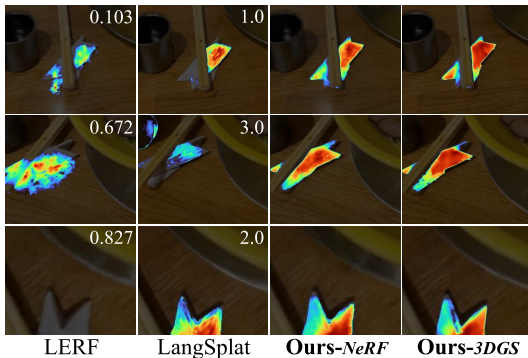


Figure 8: **Comparison of view consistency.** The figure shows the optimal scales of LERF and LangSplat for each viewpoint on `Napkin` in LERF dataset, highlighting significant variation and view inconsistency in competitors. Our method avoids this by utilizing scale-free embeddings.

4.1 Segmentation Accuracy

3D Segmentation We compare qualitative results of 3D segmentation by exporting meshes for the region obtained through 3D querying (Section 3.2). In Table 2, both of our models surpass the competitors.

In Figure 6, unlike competitors, *Ours-NeRF* and *Ours-3DGS* produce clear segmentation boundaries of the target object from the 3D scene. Notably, both of our models accurately segment complex shape objects and multiple objects like `porcelain hand` and `rubber ducks`. In the `rubber duck` query, LERF, LEGaussians, and LangSplat

Backbone	Method	LERF		3D-OVS	
		mIoU \uparrow	mAP \uparrow	mIoU \uparrow	mAP \uparrow
NeRF	LERF	31.88	30.44	52.60	57.03
	OpenNeRF	26.52	25.84	75.12	75.35
	Ours-NeRF	46.37	45.86	77.46	84.65
3DGS	LEGaussians	21.43	20.81	47.98	50.86
	LangSplat	37.53	36.39	74.54	79.49
	Ours-3DGS	44.37	44.57	77.51	84.86

Table 3: **Quantitative Results of 2D Segmentation on LERF and 3D-OVS datasets.**

completely fail to localize the target object, while OpenNeRF incorrectly segments unrelated objects (e.g., a toy elephant and Jake) along with the target object.

2D Segmentation In Table 3, both of our models show the highest 2D segmentation performance on LERF and 3D-OVS datasets. Also, we present qualitative results of 2D segmentation by comparing a heatmap for given text queries through 2D querying (Section 3.2). As shown in Figure 7, *Ours-NeRF* and *Ours-3DGS* show clear segmentation boundaries. LangSplat and LEGaussians occasionally fail to find the target object, as seen with `waldo`.

Figure 8 shows relevancy maps and optimal scales along different viewpoints. The float value on the top right of the image represents the optimal scale for LERF and LangSplat from each viewpoint. In Figure 8, our method renders view consistent relevancy maps with `napkin` query. However, LERF and LangSplat show view inconsistent segmentation results due to the changing optimal scale.

Backbone	Method	Training ↓	Rendering				
			Render ↓	Decode ↓	Post-Process ↓	FPS ↑	Real-Time
NeRF	LERF	40 mins	23688. ms	-	-	0.04	✗
	OpenNeRF	40 mins	5944.7 ms	-	-	0.17	✗
	Ours-NeRF	40 mins	2337.7 ms	-	-	0.43	✗
3DGS	LEGaussians	90 mins	15.665 ms	384.6 ms	-	2.50	✗
	LangSplat	100 mins	17.249 ms	0.935 ms	10473 ms	0.10	✗
	Ours-3DGS	40+5 mins	14.257 ms	-	-	70.1	✓

Table 4: **Computational Time.** We measure the computational cost on *waldo kitchen* scene on an RTX A5000.

Method	3D Segmentation F1-score ↑	2D Segmentation		Computational Cost		
		mIoU ↑	mAP ↑	Training Time ↓	FPS ↑	# of Gaussians ↓
Ours-3DGS	0.1354	44.37	44.57	5 mins	70.1	1M
w/o NeRF init	0.1108	36.14	36.74	13 mins	41.2	2M

Table 5: **Ablation Study of Initializing NeRF into 3DGS.** The ablation is conducted on Replica dataset for 3D segmentation and on LERF dataset for 2D segmentation. We measure the computational cost on *waldo_kitchen* scene with RTX A5000.

4.2 Computational Time

Table 4 shows the computational times of the LERF *waldo_kitchen* scene with a single RTX A5000. We report the training time and rendering time for both ours and the competitors. LangSplat and LEGaussians train an autoencoder for each scene and require decoding during rendering. Note that LangSplat trains 3DGS from scratch, optimizes the language-embedded 3DGS *per scale*, and requires additional post-processing. Ours-3DGS takes 40 minutes to train Ours-NeRF, 5 minutes to optimize 3DGS, and 22 milliseconds to query and transfer language field into 3DGS.

Ours-3DGS achieves real-time rendering of language embeddings for the first time, $28\times$ faster than LEGaussians. Note that LangSplat heavily depends on post-processing as shown in Appendix B.2.

4.3 Ablation Study

Point-wise semantic loss We demonstrate the necessity of our point-wise semantic loss. Figure 9 qualitatively shows improvements due to point-wise semantic loss compared to 2D supervision with Eq. (1). point-wise semantic loss removes meaningless floaters related to depth ambiguity. Table 6 shows that leveraging point-wise semantic loss instead of Eq. (1) makes better understanding in 3D space, improving the F1-score ($+0.0995$). Furthermore, point-wise semantic loss also indirectly improves 2D segmentation.

Scale-Free Embedding Our scale-free language embedding field from SAM-segmented objects produce greatly improved accuracy with clear object boundaries compared to multi-scale embedding field. Figure 9 shows that scale-free embedding helps covering the head of the toy elephant which was lost with multi-scale embedding. In Table 6, the quantitative performance of the 3D and 2D segmentation greatly improves using scale-free embedding.

NeRF Initialization We demonstrate the necessity of initializing 3DGS with the extracted point clouds from our pre-trained NeRF. As shown in Table 5, training 3DGS with SfM (Schönberger and Frahm 2016) initialization with the

Method	3D Segmentation	2D Segmentation	
	F1-score ↑	mIoU ↑	mAP ↑
Ours full	0.1532	46.37	45.86
w/o point-wise semantic loss	0.0537	44.50	44.05
w/o scale-free embedding	0.1114	30.03	28.85

Table 6: **Ablation Study of Supervising Semantics in 3D Space.** The datasets are Replica and LERF for 3D and 2D segmentation, respectively.

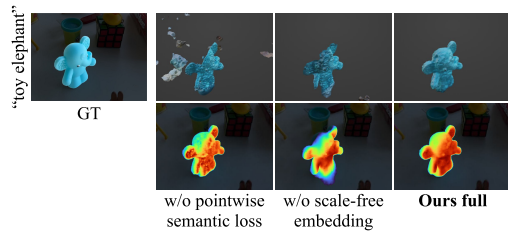


Figure 9: **Qualitative Results of Ablation Study on figurines in LERF dataset.**

densification and pruning leads to performance degradation in both 3D and 2D segmentation. This indicates that aligning geometry and semantics is essential for transferring the language field to 3DGS. Moreover, NeRF initialization produces fewer Gaussians ($0.50\times$) which lead to less training time ($0.38\times$) and faster rendering speed ($1.70\times$).

5 Conclusion

We revisit the current literature on 3D understanding of NeRF and 3DGS and revise the problem setting. We reformulate the task to produce 3D segmented volumes instead of 2D rendered masks and propose a 3D evaluation protocol. We achieve state-of-the-art segmentation accuracy in both 3D and 2D. Moreover, by transferring the language field to 3DGS, we enable the first real-time rendering speed. We hope this paper drives forward a better 3D understanding of radiance fields by reconsidering the problem set.

References

- Bhalgat, Y.; Laina, I.; Henriques, J. F.; Zisserman, A.; and Vedaldi, A. 2024. N2F2: Hierarchical Scene Understanding with Nested Neural Feature Fields. *arXiv:2403.10997*.
- Bing, W.; Chen, L.; and Yang, B. 2022. DM-NeRF: 3D Scene Geometry Decomposition and Manipulation from 2D Images. *arXiv preprint arXiv:2208.07227*.
- Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; and Joulin, A. 2021. Emerging Properties in Self-Supervised Vision Transformers. *arXiv:2104.14294*.
- Engelmann, F.; Manhardt, F.; Niemeyer, M.; Tateno, K.; and Tombari, F. 2024. OpenNerf: Open Set 3D Neural Scene Segmentation with Pixel-Wise Features and Rendered Novel Views. In *The Twelfth International Conference on Learning Representations*.
- Ghiasi, G.; Gu, X.; Cui, Y.; and Lin, T.-Y. 2022. Scaling Open-Vocabulary Image Segmentation with Image-Level Labels. *arXiv:2112.12143*.
- Guédon, A.; and Lepetit, V. 2023. SuGaR: Surface-Aligned Gaussian Splatting for Efficient 3D Mesh Reconstruction and High-Quality Mesh Rendering. *arXiv preprint arXiv:2311.12775*.
- Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics*, 42(4).
- Kerr, J.; Kim, C. M.; Goldberg, K.; Kanazawa, A.; and Tancik, M. 2023. LERF: Language Embedded Radiance Fields. In *International Conference on Computer Vision (ICCV)*.
- 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. 2023. Segment Anything. *arXiv:2304.02643*.
- Knapitsch, A.; Park, J.; Zhou, Q.-Y.; and Koltun, V. 2017. Tanks and Temples: Benchmarking Large-Scale Scene Reconstruction. *ACM Transactions on Graphics*, 36(4).
- Labe, I.; Issachar, N.; Lang, I.; and Benaim, S. 2024. DGD: Dynamic 3D Gaussians Distillation. *arXiv:2405.19321*.
- Li, Y.; and Pathak, D. 2024. Object-Aware Gaussian Splatting for Robotic Manipulation. In *ICRA 2024 Workshop on 3D Visual Representations for Robot Manipulation*.
- Li, Z.; Müller, T.; Evans, A.; Taylor, R. H.; Unberath, M.; Liu, M.-Y.; and Lin, C.-H. 2023. Neuralangelo: High-Fidelity Neural Surface Reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, K.; Zhan, F.; Zhang, J.; Xu, M.; Yu, Y.; Saddik, A. E.; Theobalt, C.; Xing, E.; and Lu, S. 2023. Weakly Supervised 3D Open-vocabulary Segmentation. *arXiv preprint arXiv:2305.14093*.
- Max, N. 1995. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2): 99–108.
- Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.
- Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.*, 41(4): 102:1–102:15.
- Qin, M.; Li, W.; Zhou, J.; Wang, H.; and Pfister, H. 2023. LangSplat: 3D Language Gaussian Splatting. *arXiv preprint arXiv:2312.16084*.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; Krueger, G.; and Sutskever, I. 2021. Learning Transferable Visual Models From Natural Language Supervision. *arXiv:2103.00020*.
- Rashid, A.; Sharma, S.; Kim, C. M.; Kerr, J.; Chen, L. Y.; Kanazawa, A.; and Goldberg, K. 2023. Language Embedded Radiance Fields for Zero-Shot Task-Oriented Grasping. In *7th Annual Conference on Robot Learning*.
- Schönberger, J. L.; and Frahm, J.-M. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Shi, J.-C.; Wang, M.; Duan, H.-B.; and Guan, S.-H. 2023. Language Embedded 3D Gaussians for Open-Vocabulary Scene Understanding. *arXiv:2311.18482*.
- Tancik, M.; Weber, E.; Ng, E.; Li, R.; Yi, B.; Kerr, J.; Wang, T.; Kristoffersen, A.; Austin, J.; Salahi, K.; Ahuja, A.; McAllister, D.; and Kanazawa, A. 2023. Nerfstudio: A Modular Framework for Neural Radiance Field Development. In *ACM SIGGRAPH 2023 Conference Proceedings, SIGGRAPH '23*.
- van Oort, T.; Miller, D.; Browne, W. N.; Marticorena, N.; Haviland, J.; and Suenderhauf, N. 2024. Open-Vocabulary Part-Based Grasping. *arXiv:2406.05951*.
- Wang, Y.; Li, Z.; Zhang, M.; Driggs-Campbell, K.; Wu, J.; Fei-Fei, L.; and Li, Y. 2023. D³Fields: Dynamic 3D Descriptor Fields for Zero-Shot Generalizable Robotic Manipulation. *arXiv preprint arXiv:2309.16118*.
- Ye, M.; Danelljan, M.; Yu, F.; and Ke, L. 2023. Gaussian Grouping: Segment and Edit Anything in 3D Scenes. *arXiv preprint arXiv:2312.00732*.
- Ye, V.; and Kanazawa, A. 2023. Mathematical Supplement for the `gsplat` Library. *arXiv:2312.02121*.
- Zhang, H.; Li, F.; and Ahuja, N. 2024. Open-NeRF: Towards Open Vocabulary NeRF Decomposition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 3456–3465.
- Zheng, Y.; Chen, X.; Zheng, Y.; Gu, S.; Yang, R.; Jin, B.; Li, P.; Zhong, C.; Wang, Z.; Liu, L.; Yang, C.; Wang, D.; Chen, Z.; Long, X.; and Wang, M. 2024. GaussianGrasper: 3D Language Gaussian Splatting for Open-vocabulary Robotic Grasping. *arXiv:2403.09637*.
- Zhi, S.; Laidlow, T.; Leutenegger, S.; and Davison, A. J. 2021. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 15838–15847.
- Zhou, S.; Chang, H.; Jiang, S.; Fan, Z.; Zhu, Z.; Xu, D.; Chari, P.; You, S.; Wang, Z.; and Kadambi, A. 2023. Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields. *arXiv preprint arXiv:2312.03203*.

<i>figurines</i>	green apple, ice cream cone, jake, miffy, old camera, pikachu, pink ice cream, porcelain hand, quilted pumpkin, rabbit, red apple, rubber duck, rubics cube, spatula, tesla door handle, toy cat statue, toy chair, toy elephant, twizzlers, waldo
<i>ramen</i>	bowl, broth, chopsticks, egg, glass of water, green onion, napkin, nori, pork belly, ramen, sake cup, wavy noodles
<i>teatime</i>	bag of cookies, bear nose, coffee, coffee mug, cookies on a plate, dall-e, hooves, paper napkin, plate, sheep, spill, spoon handle, stuffed bear, tea in a glass, yellow pouf
<i>waldo_kitchen</i>	blue hydroflask, coffee grinder, cookbooks, cooking tongs, copper-bottom pot, dish soap, faucet, knives, olive oil, paper towel roll, pepper mill, pour-over vessel, power outlet, red mug, scrub brush, sink, spice rack, utensils, vegetable oil, waldo

Table S1: Text Query of LERF Dataset

A Experiments Details

A.1 Implementation Details

We train the NeRF and the CLIP branch jointly from scratch with the same configurations of LERF except the learning rate on Replica datasets which is $1e-4$. The CLIP features are extracted using OpenCLIP ViT-B/16 trained on LAION-2B. We use an automatic mask generator from ViT-H SAM with the default configuration.

We optimize 3DGS in Section 3.4 with the same settings as the official 3DGS, except for the following modifications: 1) freezing the position of the Gaussian center \mathbf{x} , 2) setting the rotation learning rate to $1e-4$, and 3) running for 10K iterations without a densification process.

We train our models and competitors with a single RTX A5000.

A.2 Details of Datasets

We use three datasets for evaluation: Replica dataset for 3D segmentation, and LERF and 3D-OVS datasets for 2D segmentation. All of these datasets are widely used to evaluate the segmentation performance of radiance fields.

Replica dataset. Replica dataset is a high-quality indoor 3D mesh dataset with semantic labels. We measure 3D segmentation performance with labeled 3D mesh, using class names as the text queries.

LERF dataset. LERF dataset is a collection of real-world scenes captured by iPhone. This includes text queries and their ground truth 2D bounding boxes. In LERF dataset, most competitors (Qin et al. 2023; Shi et al. 2023) report their performance by building their own ground truth mask datasets corresponding to their partly chosen text queries. However, it is a challenge to compare because the performance of open-vocabulary segmentation can significantly vary depending on the text query.

Hence, starting from the given text queries and 2D bounding boxes, we use SAM to create auxiliary masks and manually correct errors. SAM(Kirillov et al. 2023) receives an image and a user prompt (bounding box or point) as input and creates three masks. Among these candidate masks, we select the ground truth mask that corresponds most closely

to the queried object. Then, we manually label a few inaccurate masks. The masks will be publicly available. We use four scenes *figurines*, *ramen*, *teatime*, *waldo_kitchen* from LERF dataset. We choose three evaluation views that are not blurry. These masks will be publicly available online.

3D-OVS dataset. 3D-OVS dataset is a collection of real-world scenes, text queries, and their ground truth pixel-level labels. The dataset consists of 10 scenes, each with 20-30 viewpoints, including five evaluation viewpoints. Several scenes in 3D-OVS dataset show poor reconstruction quality from the evaluation view shown in Figure S1. Therefore, we report the performance only for the scenes with successful reconstruction: (*bed*, *bench*, *lawn*, and *office_desk*).

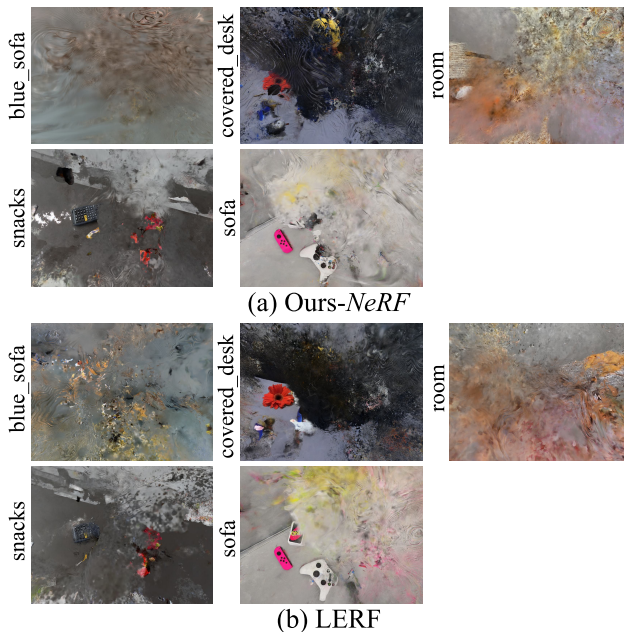


Figure S1: Poor reconstruction quality in evaluation view of 3D-OVS dataset.

B Details of Previous Language Embed 3DGS

B.1 Dimensions of Embeddings on 3DGS

Rendering 512-dimensional CLIP embeddings requires 522 KB of shared memory per streaming multiprocessor (SM). Specifically, 3DGS (Kerbl et al. 2023) requires 10 KB of shared memory and needs an additional 1 KB/SM for each increase in the dimension of the language feature. 3DGS allocates static shared memory, and CUDA limits this allocation to 48 KB for architectural compatibility. Even with a dynamic allocation that enables allocation above the 48 KB limit, 522 KB of shared memory exceeds the maximum shared memory of current top GPUs (*e.g.*, the H100 has a maximum shared memory of 228 KB).

Consequently, existing methods either 1) compress high-dimensional features or 2) allocate variables for the backward pass to global memory instead of shared memory. Each approach has trade-offs: 1) compressing features reduces segmentation performance but accelerates training, while 2) allocating to global memory avoids performance degradation but slows down training. Compressing the features produces poor mIoU (-7.30) but speeds up the training time 47 \times , as shown in Table S2. However, Ours-3DGS achieves fast training without feature compression.

Method	LERF Dataset		
	mIoU \uparrow	mAP \uparrow	Training Time \downarrow
3dim-feature	31.54	30.09	30 mins
512dim-feature	38.84	38.20	1400 mins

Table S2: **Trade-off of Feature Compression:** We compress the feature following LangSplat. For a fair comparison, we trained each model with the same number of Gaussians.

B.2 Post-process of LangSplat

We show the performance difference with and without post-processing in LangSplat², as shown in Figure S2. Without post-processing, it achieves a fast rendering speed of 55 fps (550 \times) but produces a poor mIoU (-5.5) as shown in Table S3. Note that Ours-3DGS achieve faster rendering speed even compared to LangSplat without post-processing, as shown in Table 4.

LangSplat	mIoU \uparrow	mAP \uparrow
w/ post-process	37.53	36.39
w/o post-process	32.01	30.64

Table S3: **Post-process of LangSplat**

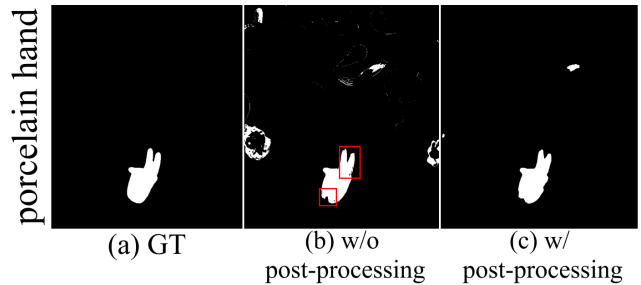


Figure S2: LangSplat includes post-processing, which significantly improves performance by reducing floaters.

C More Results

We report the quantitative results of each scene on Replica dataset, LERF dataset, and 3D-OVS dataset, in Table S4-Table S6.

²<https://github.com/minghanqin/LangSplat>

Method	<i>office_0</i> F1-score ↑	<i>office_1</i> F1-score ↑	<i>office_2</i> F1-score ↑	<i>office_3</i> F1-score ↑
LERF	0.0774	0.0191	0.0644	0.0621
OpenNeRF	0.0486	0.0169	0.0318	0.0356
Ours-NeRF	0.1434	0.0434	0.1493	0.1255
LEGaussians	0.0013	0.0043	0.0109	0.0031
LangSplat	0.0004	0.0123	0.0044	0.0038
Ours-3DGS	0.1097	0.0406	0.1395	0.1181

Method	<i>office_4</i> F1-score ↑	<i>room_0</i> F1-score ↑	<i>room_1</i> F1-score ↑	<i>room_2</i> F1-score ↑
LERF	0.0645	0.1124	0.1607	0.1158
OpenNeRF	0.0521	0.0180	0.0432	0.0423
Ours-NeRF	0.1722	0.1654	0.1934	0.2232
LEGaussians	0.0111	0.0080	0.0066	0.0085
LangSplat	0.0107	0.0139	0.0122	0.0118
Ours-3DGS	0.1629	0.1636	0.1529	0.1959

Table S4: Per-scene quantitative results on Replica dataset.

Method	<i>figurines</i>		<i>ramen</i>		<i>teatime</i>		<i>waldo.kitchen</i>	
	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑
LERF	40.67	38.75	24.56	24.35	36.85	35.06	25.44	23.62
OpenNeRF	16.24	16.02	23.05	23.57	39.08	38.26	27.69	25.50
Ours-NeRF	59.79	59.00	37.44	38.38	47.59	46.92	40.67	39.16
LEGaussians	21.02	20.57	24.90	24.67	24.50	24.09	15.28	13.89
LangSplat	42.23	40.55	48.17	48.71	37.74	36.26	21.98	20.05
Ours-3DGS	56.23	55.79	38.28	40.81	41.59	41.16	42.60	41.79

Table S5: Per-scene quantitative results on LERF dataset.

Method	<i>bed</i>		<i>bench</i>		<i>lawn</i>		<i>office_desk</i>	
	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑	mIoU ↑	mAP ↑
LERF	45.29	57.74	50.88	51.59	54.19	57.40	60.04	61.41
OpenNeRF	82.25	82.40	70.24	71.14	78.70	79.27	69.29	69.79
Ours-NeRF	66.71	78.65	66.69	78.81	96.74	96.98	79.72	84.18
LEGaussians	33.23	44.94	63.86	64.25	43.48	44.20	51.36	50.05
LangSplat	65.74	73.01	82.75	88.06	82.41	87.95	67.26	68.93
Ours-3DGS	67.40	79.33	67.08	79.20	96.25	96.92	78.89	84.36

Table S6: Per-scene quantitative results on 3D-OVS dataset.