# CoSSegGaussians: Compact and Swift Scene Segmenting 3D Gaussians

**Bin Dou** [1] , **Tianyu Zhang** [1] , **Yongjia Ma** [1] , **Zhaohui Wang** [1] , **Zejian Yuan** [1] *

[1]Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University

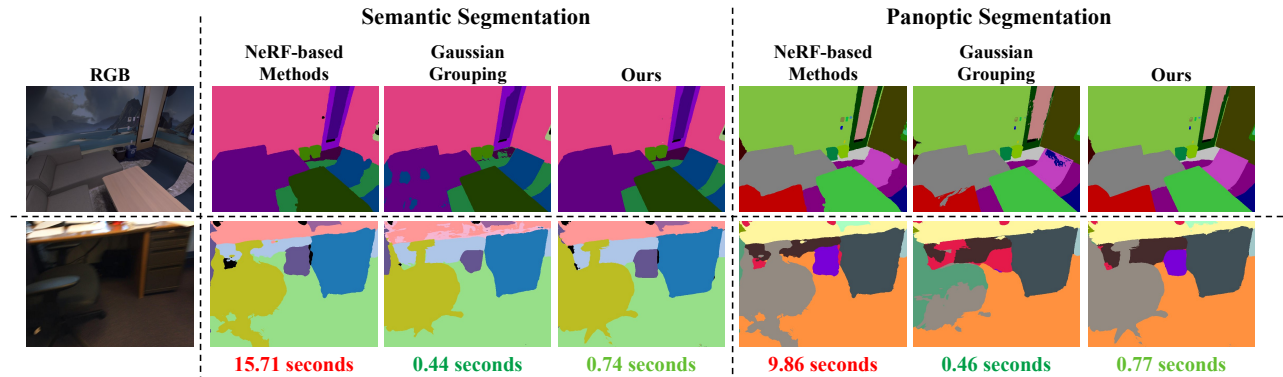{abc991227db, puzzling229, maguire, astjkl12345}@stu.xjtu.edu.cn, yuan.ze.jian@xjtu.edu.cn,

Figure 1: Our method can achieve compact and swift scene segmentation with only RGB input. Once trained, our model can render segmentation maps over 10× than NeRF-based segmentation(over 20× than Semantic-NeRF and 12× than Panoptic-Lifting which is accelerated using TensoRF) due to our selected 3D Gaussian Splatting representation. Compared to Gaussian Grouping, though sacrificing a bit speed, our method achieves much more compact and reliable zero-shot segmentation due to our designed decoder upon multi-scale fused features.

## Abstract

We propose **Co**mpact and **S**wift **Seg**menting 3D **Gaussians**(**CoSSegGaussians**), a method for compact 3D-consistent scene segmentation at fast rendering speed with only RGB images input. Previous NeRF-based 3D segmentation methods have relied on implicit or voxel neural scene representation and ray-marching volume rendering which are time consuming. Recent 3D Gaussian Splatting significantly improves the rendering speed, however, existing Gaussians-based segmentation methods(eg: Gaussian Grouping) fail to provide compact segmentation masks especially in zero-shot segmentation, which is mainly caused by the lack of robustness and compactness for straightforwardly assigning learnable parameters to each Gaussian when encountering inconsistent 2D machine-generated labels. Our method aims to achieve compact and reliable zero-shot scene segmentation swiftly by mapping fused spatial and semantically meaningful features for each Gaussian point with a shallow decoding network. Specifically, our method firstly optimizes Gaussian points' position, convariance and color attributes under the supervision of RGB images. After Gaussian Locat-

ing, we distill multi-scale DINO features extracted from images through unprojection to each Gaussian, which is then incorporated with spatial features from the fast point features processing network, i.e. RandLA-Net. Then the shallow decoding MLP is applied to the multi-scale fused features to obtain compact segmentation. Experimental results show that our model can perform high-quality zero-shot scene segmentation, as our model outperforms other segmentation methods on both semantic and panoptic segmentation task, meanwhile consumes approximately only 10% segmenting time compared to NeRF-based segmentation. Code and more results will be available at https://David-Dou.github.io/CoSSegGaussians

## 1 Introduction

In recent years, computer vision and computer graphics have achieved notable advancements, particularly in the neural rendering area. Neural Radiance Fields (NeRF)[Mildenhall *et al.*, 2021] and its subsequent methods have propelled the development of neural scene representations, which have shown significant capabilities for tasks such as novel view synthesis.

Besides learning a radiance field, some methods extended

from NeRF can achieve 3D-consistent segmentation for scene understanding. Early segmenting models rely on manually annotated 2D segmentation labels as supervision and employ ray-marching volume rendering which are time consuming. After the release of large visual models like DINO-Vit[Caron *et al.*, 2021] and SAM[Kirillov *et al.*, 2023], NeRF-based models can achieve few-shot or zero-shot scene segmentation with the introduction of semantically meaningful knowledge or machine-generated labels. Recently, benefiting from the emergence of 3D Gaussian Splatting[Kerbl *et al.*, 2023] which utilizes point-like Gaussians as scene representation and employs rapid rasterization method for rendering, some methods have applied 3D Gaussians as scene representation for further swift segmentation, such as Gaussian Grouping[Ye *et al.*, 2023]. Nevertheless, as previous Gaussians-based segmentation models commonly assign learnable parameters for each Gaussian which are updated according to the zero-shot segmented masks, 3D-inconsistency of the pseudo labels(though Gaussian Grouping utilizes mask association technique, inconsistency still occurs especially for complex scenes) will lead incompactness for scene segmentation. Fig. 2 shows that Gaussian Grouping fails to segment the *ground* and *cabinet* compactly and reliably mainly due to the inconsistency of machine-generated labels, which may further lead to failure for downstream tasks, like scene editing.
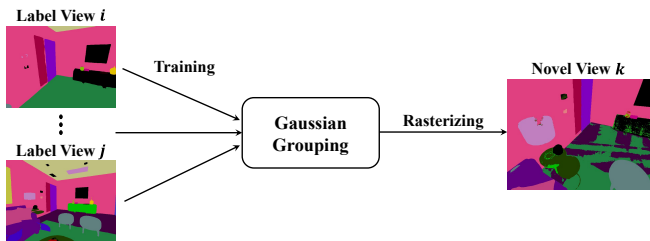


Figure 2: Reason for incompact zero-shot segmentation in previous Gaussians-based segmentation method.

To improve the segmenting compactness especially for zero-shot situation as well as enhance model's training and inference speed, we aim to design a compact scene segmentation model based on 3D Gaussians. We conjecture that replacing attached learnable attributes on each Gaussian with a decoder upon segmenting features for each Gaussian can enhance the model's robustness to inconsistency if the feature acts as a excellent prior. Motivated by the DINO distillation in interactive scene segmentation and spatial feature extraction point cloud segmentation, we think fusion of them can be used as suitable segmening feature.

Therefore in this paper, we propose CoSSegGaussians, which utilizes multi-scale fused features as segmenting prior to be further neurally decoded to compact segmenting features. Specifically, the method firstly optimizes 3D Gaussian parameters to locate Gaussians. Then considering the explicitness of Gaussians, multi-scale DINO features are distilled through explicit unprojection. For further segmentation, encoding layers of RandLA-Net(a swift point cloud segmentation model) is employed to extract Gaussians' spatial features which is concatenated with distilled DINO features and de-

coded as segmentation logits.

In summary, our work makes the following contributions:

• We propose a network decoding method to represent segmenting 3D Gaussians instead of learnable parameters on each Gaussian point, which will enhance the segmenting compactness and robustness especially in zero-shot segmentation.

• We present an explicit method to distill extracted 2D DINO features into 3D Gaussians by utilizing unprojection.

• We design a multi-scale spatial and DINO features fusion module, which assists to achieve compact and swift segmentation.

## 2 Related Work

### 2.1 Neural Scene Representation

Neural scene representation has seen significant advancements in recent years, particularly in the realm of the novel view synthesis task. NeRF[Mildenhall *et al.*, 2021] utilizes neural networks to implicitly encode the geometry and appearance of 3D space, combining differentiable volume rendering to achieve high-quality novel view synthesis. Subsequent works have introduced numerous improvements, with some aiming to further enhance rendering quality. Mip-NeRF[Barron *et al.*, 2021] has improved the way light is modeled, while NeRF++[Zhang *et al.*, 2020] and Mip-NeRF 360[Barron *et al.*, 2022] have enhanced the sampling method. Another category of works focus on accelerating the neural scene optimization and inference. Plenoxel[Fridovich-Keil *et al.*, 2022], TensoRF[Chen *et al.*, 2022], and Instant NPG[Müller *et al.*, 2022] use explicit voxels or hybrid representations to improve upon purely implicit MLP representations. Recently, 3D Gaussian Splatting[Kerbl *et al.*, 2023] has been proposed to represent the scene using 3D Gaussians, which is able to achieve high-quality and real-time rendering and garners widespread attention.

### 2.2 NeRF-based Scene Segmentation

Numerous approaches have been developed to extend NeRF to segmentation field for 3D scene understanding, using merely 2D labels as supervision. Semantic-NeRF[Zhi *et al.*, 2021] facilitates concurrent optimization of the 3D radiance and semantic field by amalgamating semantic, geometric and appearance encoding for rendering semantically consistent segmentation masks. DM-NeRF[Wang *et al.*, 2022] achieves 3D object-level segmentation by enhancing semantic field training with the incorporation of the designed object field. However, these methods rely on dense-view annotated labels to achieve high-quality segmentation results. Benefiting from the development of zero-shot 2D segmentation models pretrained on large-scale datasets, some methods extend them for segmenting scenes without manually labeled annotations. Panoptic Lifting[Siddiqui *et al.*, 2023] builds scene representation on TensoRF and achieves zero-shot 3D-consistent panoptic segmentation supervised by Mask2Former's[Cheng *et al.*, 2022] segmentation. Gaussian Grouping[Ye *et al.*, 2023] extends 3D Gaussian representation by jointly reconstructing and segmenting anything in full open-world 3D scenes with instance and stuff level modeling and achieves

rapid scene segmentation and editing. Beyond those semantic or panoptic scene segmentation tasks, some interactive 3D segmentation methods are also proposed which aim to segment required object in the scene according to the user's prompt on single image. They often distill large-model extracted 2D knowledge into neural scene representation via rendering and optimization to transfer user-provided information on single-view. For example, DFF[Kobayashi *et al.*, 2022] and ISRF[Goel *et al.*, 2023] distills knowledge from DINO for click-prompted and language-prompted segmentation, LERF[Kerr *et al.*, 2023] lift CLIP[Radford *et al.*, 2021] features into NeRF to achieve 3D open-vocabulary segmentation and SAGA[Cen *et al.*, 2023] applies features from SAM[Kirillov *et al.*, 2023] to 3D Gaussians representation to realize fast and compact interactive segmentation.

## 2.3 Point Cloud Segmentation

Considering that 3D Gaussian is a point-like representation, we conjecture that modules in scene point cloud segmentation can be extended to Gaussians-based scene segmentation. PointNet[Qi *et al.*, 2017a] directly learns a spatial encoding of each point and PointNet++[Qi *et al.*, 2017b] extends it with local feature extractor based on Farthest Point Sampling(*FPS*) and is trained with hierarchical feature learning architecture. However, *FPS* is computationally expensive for large-scale point clouds and its splitting strategy for scene segmentation even increases the computational burden. SPG[Landrieu and Simonovsky, 2018] preprocesses the large point clouds as superpoint graphs to learn per superpoint semantics which is based on a more time-consuming graph partitioning. To improve model efficiency especially for large-scale point cloud segmentation, RandLANet[Hu *et al.*, 2020] employ the remarkably efficient random sampling combined with Local Feature Aggregation(Local Spatial Encoding+Attentive Pooling). Trading-off the efficiency and performance, we try to combine 3D Gaussians with the swift RandLA-Net for segmenting improvement.

## 3 Method

Given only posed RGB images of a 3D scene, our method aims to build an expressive representation to capture appearance, geometry as well as compact segmenting identity of the scene. Our proposed method, CoSSegGaussians, enables compact novel-view 3D-consistent scene segmentation, while consuming much less rendering time compared to NeRF-based segmentation methods. Fig. 3 provides an overview of our model's architecture. As we design our method based on the recent 3D Gaussian Splatting which is able to represent scenes and render novel-view images swiftly, we first review 3D Gaussian Splatting in 3.1. Then we introduce our method to capture compact segmenting features in 3.2 and 3.3. In addition, we present the training startegy of our model in 3.4

## 3.1 Preliminaries: 3D Gaussian Splatting

3D Gaussian splatting is a recent effective scene representation, which is able to achieve remarkable quality of 3D scene reconstruction and exhibits much higher inference speed compared to previous Neural Radiance Field representation.

It represents the scene explicitly with 3D Gaussians sharing great amount of similarity with point cloud, as each Gaussian is parameterized by its centroid $x \in R^3$, the scale $s \in R^3$ and rotational quaternion $q \in R^4$($s$ and $q$ are used jointly to represent the Gaussian's 3D covariance), opacity $\alpha \in R$ and color $c$ as the three degrees of spherical harmonics(SH) coefficients. To supervise aforementioned learnable attributes of 3D Gaussians, Gaussian Splatting projects them onto 2D image plane for given viewpoints to render RGB image by $\alpha$-blending, which is a differentiable rasterization method and is implemented cuda-friendly. For each pixel, It can be formulated as

$$C = \sum_{i \in \mathcal{N}} c_i \alpha'_i \prod_{j=1}^{i-1}(1 - \alpha'_j) \tag{1}$$

where $c_i$ represents the $i$-th Gaussian's color, $\alpha'_i$ represents the $i$-th Gaussian's influence factor which is obtained by multiplying the projected 2D covariance with the earned per-point opacity $\alpha_i$,

## 3.2 Multi-scale DINO Feature Unprojection

DINO features[Caron *et al.*, 2021] extracted from RGB images have been employed as semantically meaningful correspondences in recent segmentation works due to their semantic consistency within the same image and across image collections. Previous neural radiance field methods[Kobayashi *et al.*, 2022] distill DINO features by assigning learnable attributes on scene representation and supervise them by comparing the rendered feature map with the 2D extracted feature. However, considering that 3D Gaussians utilize explicit point clouds to represent the scene, here we choose to use an explicit method to distill 2D DINO features into 3D scene representation. Inspired by semantic tracing method in [Chen *et al.*, 2023] which achieves inverse rendering multi-view segmentation masks of certain object by weighted summing 2D masks(the weight is calculated by the affected Gaussian's opacity and transmittance), we conjecture that such weighted-sum solution can also be used for unprojecting semantic correspondences. As more comprehensive and extensive prior knowledge of the scene will enhance the segmenting compactness, here we utilize multi-scale DINO features for distillation. Thus we extend the inverse rendering for DINO feature distillation, as the $i$-th Gaussian's distilled feature for scale $n$ can be expressed as:

$$\boldsymbol{f}_{n,i} = \sum \alpha'_i(\boldsymbol{p}) * \prod_{j=1}^{i-1}(1 - \alpha'_j(\boldsymbol{p})) * \boldsymbol{f}_n(\boldsymbol{p}), n = 1 \cdots 4 \tag{2}$$

where $\boldsymbol{f}(\boldsymbol{p})$ represents DINO feature vector on pixel $\boldsymbol{p}$.

As the feature inverse rendering method can also record the counter of affected pixels for each Gaussian, we manually set threshold to prune points according to the counter, which will improve efficiency for further feature processing experimentally. We provide the unprojected DINO features in the **Supplementary Material** to show the effectiveness of features distillation through unprojection.

## 3.3 RandLA-based Feature Fusion

Previous Gaussians-based segmentation methods lack compactness especially with machine-generated segmentation la-
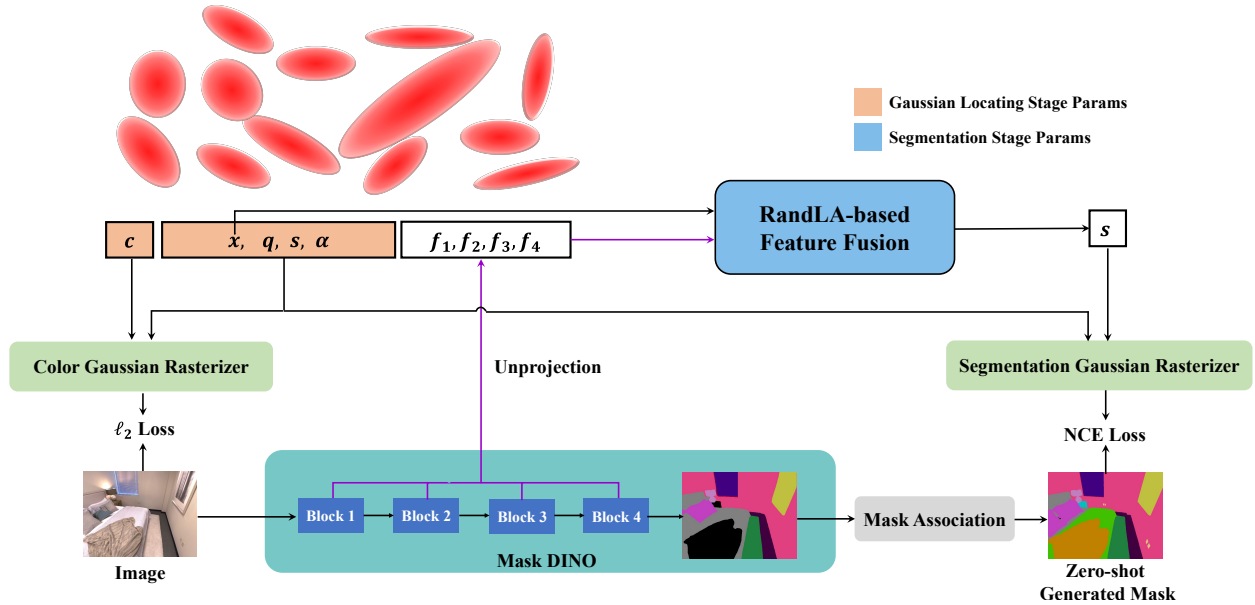
Figure 3: Overview of our method. Our model performs segmentation by capturing a scene based on 3D Gaussian Splatting. Gaussians are firstly located by optimizing parameters via differentiable rasterizer, supervised by $\ell_2$ photometric loss. Then multi-scale extracted 2D DINO features are distilled into Gaussians by explicit unprojection to be further incorporated with spatial features derived from RandLA-Net. Learnable decoder for the fused features is used for predict segmenting logits on Gaussians and is supervised by the zero-shot segmented and associated masks, with NCE loss.

bels, which may be caused by inconsistent 2D segmentation labels according to our hypothesis as shown in Fig.2. To alleviate this problem, we utilize shallow MLP network to encode fused spatial and distilled DINO features which are more robust to such zero-shot issue compared to directly assigning learnable attributes as segmentation features.

Inspired by previous point cloud segmentation methods, here we employ pretrained RandLA-Net[Hu *et al.*, 2020], an efficient segmentation method for large-scale and complex scenes, as Gaussians' spatial feature extractor. As rasterization only depends on Gaussians inside the viewing cone and RandLA-Net doesn't need points from complete objects due to its designed attentive pooling module(not as *PointNet++*-based methods which split the scene into spatial blocks and process points inside each block), we employ only $N$ visible Gaussians for each view as the input for efficiency. Then we treat these Gaussians as points and utilize pretrained Dilated Residual Block consisting of Local Spatial Encoding, Attentive Pooling together with Random Sampling on $N$ Gaussians' centroids to obtain spatial features. For the following decoding modules, we incorporate the spatial features with unprojected DINO prior by extending the previous skip connection. Specifically, we concatenate each scale's distilled DINO features with the skipped Dilated Residual Block encoding features and the decoding features, from local to global. Decoding layer's MLP are updated in the segmentation training to decode compact segmentation features for each Gaussian and $(N, K)$ features are obtained($K$ represents segmentation classes in the scene).

## 3.4 Training Strategy

As our compact segmentation method is based on the position of 3D Gaussians, the training is conducted in a two-stage strategy: Gaussian Locating Stage and Segmentation Stage.

Gaussian Locating Stage is similar to 3D Gaussians Splatting training, as Gaussians' centroids $x$, scales $s$, rotational quaternions $q$, opacities $\alpha$ and colors $c$ are supervised with $\ell_1$ loss and $\ell_{D-SSIM}$ by comparing the rendering with the ground-truth RGB image.

In Segmentation Stage, multi-scale DINO features are firstly unprojected to all located Gaussians and can be further queried with the sampling point id in the concatenating operation for RandLA-based feature fusion. After feature fusion, we extend the aforementioned Gaussian Rasterizer from color space to segmentation logits space, which can be expressed as:

$$S = \sum_{i \in \mathcal{N}} s_i \alpha_i' \prod_{j=1}^{i-1} (1 - \alpha_j') \qquad (3)$$

As for the supervisory labels, we adapt the mask generation method in [Ye *et al.*, 2023]. We first employ a DINO-based 2D scene segmentation method(i.e. MaskDINO) to generate masks for all RGB images. Then multi-view images of a scene is treated as a video sequence so a well-trained zero-shot tracker[Cheng *et al.*, 2023] is employed to propagate and associate the generated masks for cross-view consistency enhancement. The segmentation parameters are optimized with the NCE loss by comparing the rendering mask with the gen-
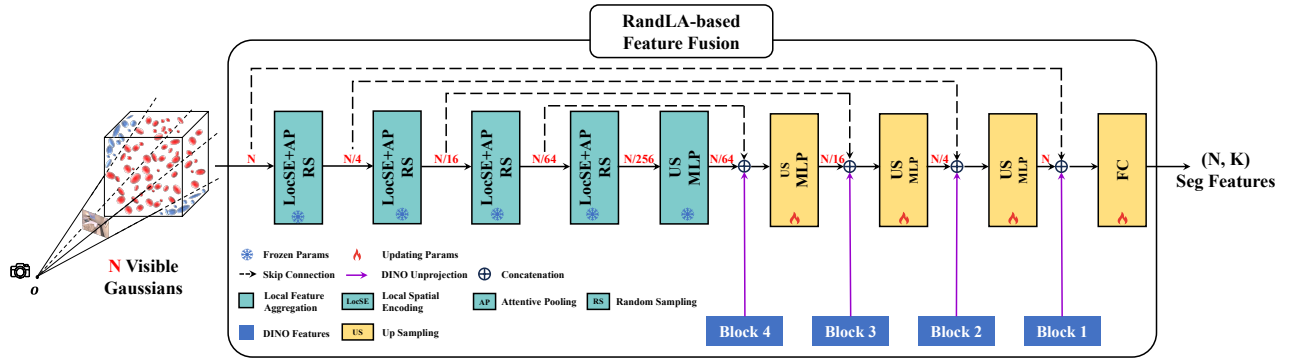
Figure 4: Our designed RandLA-based feature fusion. N visible Gaussians for each viewpoint are utilized as input to the fusion module. RandLA-Net's encoding layers(Local Features Aggregation consisting of Local Spatial Encoding & Attentive Pooling) are frozen to extract Gaussians' spatial features, which are then skipped and concatenated with multi-scale distilled DINO features. MLP decoder is used to predict segmentation logits for each Gaussian by decoding fused features and requires gradient optimization during training.

erated mask, which can be formulated as:

$$\ell_s = -\sum_p \sum_{k \in K} \kappa_r[k] \log(\mathbf{S}(p)[k]) \quad (4)$$

Therefore, the objective for our method is:

$$\ell_{GL} = \lambda_{rgb}\ell_1 + (1 - \lambda_{rgb})\ell_{D-SSIM} \\ \ell_{Seg} = \lambda_s \ell_s \quad (5)$$

## 4 Experiments

We evaluate our method on the novel view scene semantic and panoptic segmentation task, and further demonstrate the potential for scene manipulation as an application of our method. More experimental setups and results can be found in **Supplementary Material**.

### 4.1 Experimental Setup

**Datasets**
We show experimental results on scenes from two public indoor scene datasets: Replica[Straub *et al.*, 2019] and Scan-Net[Dai *et al.*, 2017]. Replica is a 3D dataset consisting of high-fidelity room or office scenes and each scene consists of RGB images with corresponding 2D segmentation masks, rendered by the authors of [Zhi *et al.*, 2021]. ScanNet is a large-scale challenging real-world dataset and each scene comprises images accompanied by 2D segmenting masks and camera poses. For both datasets we select 7 scenes for training and evaluation(similar to [Wang *et al.*, 2022]). To generate training labels, we utilize MaskDINO[Li *et al.*, 2023] which is pretrained on ADE20K[Zhou *et al.*, 2019] for zero-shot semantic and panoptic segmentation on indoor scene images. The evaluation is conducted between all poses' segmentation predictions and corresponding targets in the testset. It should be noticed that the annotated segmentation masks are only available in evaluation as ground-truth labels and are not used during training.

**Evaluation Metrics**
For semantic segmentation, we employ the mean intersection over union(**mIoU**) to evaluate the quality of predicted

segmentation maps. For panoptic segmentation, aside from **mIoU**, mean average precision (**mAP**) on masks is utilized, with IoU thresholds set at 0.50. In addition, for better evaluation of panoptic segmentation, we introduce scene-level panoptic quality(**PQ^scene**) derived from [Siddiqui *et al.*, 2023] for every scene, as we compare all pairs of *thing* or *stuff* segmentation subsets for all renderings in each scene instead of one rendered image, and record a match when IoU>0.50.

In addition to the above performance metrics, we also use **FPS**(Frames Per Second) and learnable parameters in training to evaluate model efficiency.

**Multi-scale Features Setup**
For spatial features, we employ RandLANet pretrained on the large-scale indoor dataset S3DIS[Armeni *et al.*, 2017] and its encoding layers are frozen during training. In reference to related feature distillation work[Tschernezki *et al.*, 2022; Kobayashi *et al.*, 2022], we use DINO[Caron *et al.*, 2021] as the teacher network to extract semantically meaningful prior from RGB images. For multi-scale setting, we select outputs of 4 transformer blocks inspired by [Sharma *et al.*, 2023]. Before unprojection to 3D Gaussians, these features are concatenated with the expanded global feature tensor then reduced to 64 dimensions by PCA and $L_2$ normalized.

**Implementation Details**
We implement our method based on Gaussian Splatting[Kerbl *et al.*, 2023] and RandLANet[Hu *et al.*, 2020](Pytorch implementation). We add segmentation logits, which is the processed feature from RandLA-based decoder, as a feature of each Gaussian and implement forward and backward rasterization in reference to [Ye *et al.*, 2023]. In training, $\lambda_{rgb} = 0.2$ and $\lambda_s = 0.1$. We use Adam optimizer for updating both Gaussians parameters and point processing network, with Gaussians' learning rates identical to Gaussian Splatting and 0.001 for MLP decoder. All datasets are trained and evaluated for 10K iterations on one NVIDIA A800 GPU.

### 4.2 Comparisons

We conduct a comparative analysis on two kinds of segmentation tasks: semantic and panoptic segmentation, as

the first one assigns semantic labels for all pixels, and the second attach labels for not only foreground objects(defined as *thing*) but background(defined as *stuff*) while provides object-level labels only for *thing*. In training, we employ semantic and panoptic mode of MaskDINO to generate zero-shot 2D segmentation labels for each task respectively. We compare our method with NeRF-based 3D segmentation methods(Semantic-NeRF[Zhi *et al.*, 2021] in semantic segmentation and Panoptic-Lifting[Siddiqui *et al.*, 2023], a segmenting method built on TensoRF for acceleration in panoptic segmentation) and the recent Gaussians-based segmentation, i.e. Gaussian Grouping[Ye *et al.*, 2023], which are all trained with the same generated labels and evaluated on all test viewpoints.

### Semantic Segmentation

| Model | Replica | | | ScanNet | | |
|---|---|---|---|---|---|---|
| | mIoU(%) | FPS | Learnable Params(MB) | mIoU(%) | FPS | Learnable Params(MB) |
| Semantic-NeRF | 69.81 | ≈ 5 | 9.7 | 69.67 | ≈ 5 | 10.4 |
| Gaussian Grouping | 66.84 | ≈140 | 838.0 | 67.99 | ≈150 | 838.6 |
| Ours | **79.34** | ≈ 90 | 680.3 | **75.40** | ≈100 | 684.5 |

Table 1: Comparison on zero-shot novel-view semantic segmentation. Our model outperforms baselines by a wide margin and achieves much faster speed compared to the NeRF-based method, showcasing our efficient tackling for compact and swift zero-shot scene segmentation.

Tab. 1 quantitatively shows segmentation results under machine-generated semantic label supervision. It can be observed that our method obtains remarkably higher segmenting performance(+9.53%/5.73% mIoU on Replica/ScanNet than Semantic-NeRF, +12.5%/7.41% mIoU on Replica/ScanNet than Gaussian Grouping). Though sacrificing rendering time compared to Gaussian Grouping, our method also achieves significantly faster rendering speed than NeRF-based semantic segmentation model(about 10 times) and performance improvement from Gaussian Grouping proves the meaning of the little more time consumption.
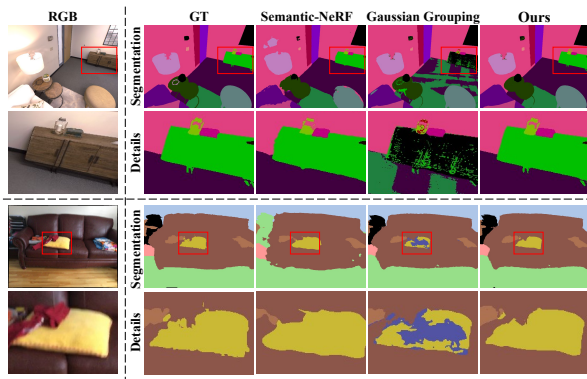


Figure 5: Qualitative comparison on semantic segmentation.

Similar segmenting performance improvement can also be observed qualitatively in Fig. 5. We can also find that Gaussian Grouping fail to provide compact semantic masks especially for some large-scale objects such as ground and cabinet, which may be due to the lack of grouping constraints for large objects as Gaussian Grouping's 3D regularization loss only restricts spatial neighbor point features. In contrast, the introduction of spatial and DINO features in our method can improve the segmenting compactness of these objects and lead to the performance enhancement.
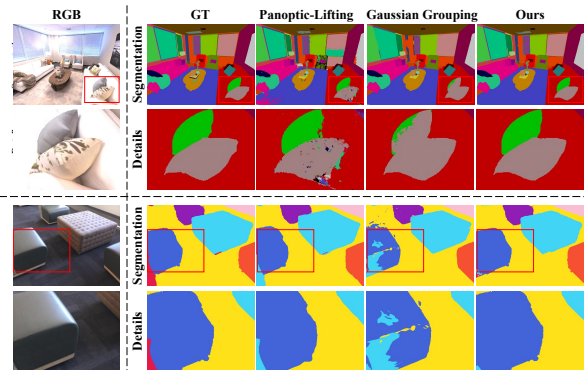
### Panoptic Segmentation



Figure 6: Qualitative comparison on panoptic segmentation.

Quantitative panoptic segmentation results is presented in Tab. 2, which also shows that our model outperforms baselines in the zero-shot panoptic segmentation task, leading to approximately 10%/9%/10% improvement in mIoU/mAP$^{0.50}$/PQ$^{scene}$ compared to Panoptic-Lifting and 7%/10%/8% improvement in mIoU/mAP$^{0.50}$/PQ$^{scene}$ compared to Gaussian Grouping. It can be noticed that Gaussian Grouping can achieve better panoptic than semantic segmentation mainly due to its KNN grouping for segmenting scenes finer-grainedly. However, our method still outperforms it with a bit more time fusing and processing features as a more tight and compact spatial regularization way.

As depicted in Fig 6, Panoptic-Lifting and Gaussian Grouping fail to segment some objects compactly such as the pillow and the seat. Conversely our model can alleviate the problem and provide a more compact segmentation mask.

### 4.3 Abaltion Studies

To discover each component's contribution to zero-shot segmenting improvement, we conduct a series of ablation experiments on semantic segmentation task for the Replica dataset using the same generated training labels. As a first baseline (Tab 3 row 1), we disable both spatial and unprojected DINO features and train the segmenting features on Gaussians from scratch. Though achieving the faster speed, the baseline will suffer from a significant performance decrease(-17.55% mIoU) compared to the final model. Visualized results for segmentation rendering and segmented 3D Gaussians in Fig 7 (a) and (e) also demonstrates the performance decline.

### DINO Feature Distillation

To verify the effect of DINO feature distillation through unprojection, we conduct a comparison on the model with-

| | Replica | | | | | ScanNet | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Model | mIoU(%) | mAP$^{0.50}$(%) | PQ$^{scene}$(%) | FPS | Learnable Params(MB) | mIoU(%) | mAP$^{0.50}$(%) | PQ$^{scene}$(%) | FPS | Learnable Params(MB) |
| Panoptic-Lifting | 66.22 | 72.25 | 64.34 | ≈ 10 | 100.1 | 67.01 | 66.95 | 60.74 | ≈ 10 | 102.5 |
| Gaussian Grouping | 71.15 | 69.98 | 66.52 | ≈140 | 840.2 | 68.70 | 67.03 | 61.83 | ≈150 | 841.6 |
| Ours | **79.49** | **80.97** | **75.60** | ≈ 90 | 681.1 | **74.37** | **76.14** | **68.68** | ≈ 90 | 685.4 |

Table 2: Comparison on zero-shot novel-view panoptic segmentation. The leading performance advantage similar to semantic segmentation also occurs in the task.
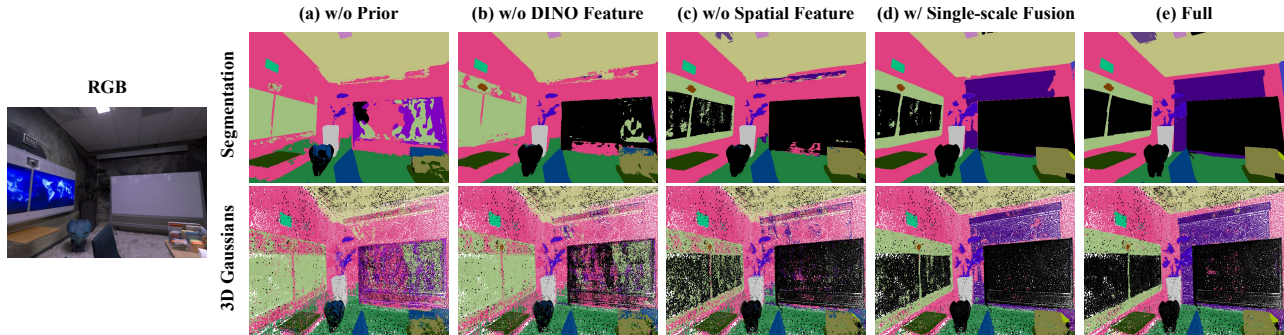


Figure 7: Visualized ablation results. Both segmentation renderings and 3D Gaussians colored with segmentation results are presented.

| Spatial Feature | DINO Feature | mIoU(%) | FPS |
|---|---|---|---|
| × | × | 61.79 | ≈160 |
| √ | × | 67.53 | ≈110 |
| × | √ | 72.33 | ≈150 |
| √ | √ (single-scale fusion) | 75.69 | ≈100 |
| √ | √ | **79.34** | ≈ 90 |

Table 3: Quantitative ablation results as segmenting performance and speed are reported.

out DINO distillation(using only spatial information to predict segmentation logits in segmentation stage) with the full model. Upon comparing Tab. 3, row 2 with row 7, we observe that the addition of distillation lead to remarkable improvement in segmentation(mIoU +11.81%). As depicted in Fig 7 (b) and (e), the full model exhibits inferior segmentation performance and distilled DINO features can enhance segmenting compactness for large objects, such as the board, due to the semantically meaningful correspondences in the 2D DINO feature map. These findings clearly indicate the crucial role played by DINO feature distillation in the zero-shot segmentation.

**Gaussians' Spatial Features**
The effect of spatial features on Gaussian points extracted from RandLA-Net is also evaluated as we employ MLP network to directly decode the distilled DINO features on Gaussians to obtain segmentation logits. As shown in Tab. 3, row 3, the absence of spatial information results in a performance decline, specifically -7.01% mIoU, and also shown in Fig. 7 (c).

Furthermore, FPS comparison on Tab. 3, row 1 row 3 and row 5 indicates that extracting and processing Gaussian's spatial features cause the time consumption, leading to a 60 FPS

decrease approximately. However, such time cost is not so severe in practice due to the swift Gaussian rendering and the segmenting performance enhancement proves the importance of spatial feature introduction.

**Multi-scale Feature Fusion**

Additionally, we evaluate the effect of our designed multi-scale spatial and DINO features fusion module by only incorporating distilled features from DINO block 1 and the global spatial feature as single-scale feature fusion. As shown in Tab. 3, row 5, the absence of multi-scale feature fusion results in a segmenting performance decrease, a specific decline of 3.65% in mIoU. By comparing Fig. 7 (e) with (d), we can find that multi-scale fusion can further enhance the segmenting compactness for objects, like the screens.

## 4.4 Application

Once our model is trained under panoptic segmentation mode, in addition to generating novel-view segmentation masks, it can also be used for further 3D object-level applications based on segmentation masks, such as scene manipulation(i.e. novel-view synthesis of a scene with 3D object removed, duplicated or manipulated under affine transformations). Benefiting from the 3D Gaussians representation' explicit nature, we can conveniently apply manipulations on the trained scene.

As shown in Fig. 8, with scene manipulation operations, our method can achieve more realistic results compared to Gaussian Grouping due to the segmenting compactness. More applications such as the language-based 3D object segmentation are available in **Supplementary Material**.
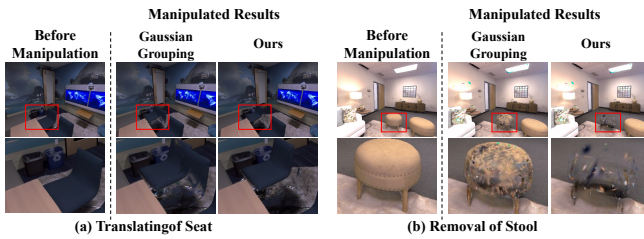
Figure 8: Scene manipulation results based on the trained segmentation model. Compared to Gaussian Grouping? our model can get higher fidelity results.

# 5  Conclusion

We propose a method CoSSegGaussians, which achieves compact and swift 3D-consistent scene segmentation with only posed RGB images. Our method unprojects DINO features extracted from images as prior and concatenate the multi-scale distilled features with the spatial information from RandLA-Net to obtain fused features, which is then passed to the MLP-decoder to get the compact segmentation for each Gaussian. Results illustrate that our model can accomplish the zero-shot segmentation task reliably and efficiently. Furthermore, we present object-level scene manipulation results as downstream application of our model.

## References

Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017.

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.

Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.

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.

Jiazhong Cen, Jiemin Fang, Chen Yang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Segment any 3d gaussians. *arXiv preprint arXiv:2312.00860*, 2023.

Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer Nature Switzerland Cham, 2022.

Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521*, 2023.

Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022.

Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1316–1326, 2023.

Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017.

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.

Rahul Goel, Dhawal Sirikonda, Saurabh Saini, and PJ Narayanan. Interactive segmentation of radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4201–4211, 2023.

Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11108–11117, 2020.

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.

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.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems*, 35:23311–23330, 2022.

Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4558–4567, 2018.

Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3041–3050, 2023.

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. *Communications of the ACM*, 65(1):99–106, 2021.

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.

Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.

Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

Prafull Sharma, Julien Philip, Michaël Gharbi, Bill Freeman, Fredo Durand, and Valentin Deschaintre. Materialistic: Selecting similar materials in images. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023.

Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kontschieder. Panoptic lifting for 3d scene understanding with neural fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9043–9052, 2023.

Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.

Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In *2022 International Conference on 3D Vision (3DV)*, pages 443–453. IEEE, 2022.

Bing Wang, Lu Chen, and Bo Yang. Dm-nerf: 3d scene geometry decomposition and manipulation from 2d images. *arXiv preprint arXiv:2208.07227*, 2022.

Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. *arXiv preprint arXiv:2312.00732*, 2023.

Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.

Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021.

Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127(3):302–321, 2019.