

SuperGSeg: Open-Vocabulary 3D Segmentation with Structured Super-Gaussians

Siyun Liang^{1,2*} Sen Wang^{1,4*} Kunyi Li^{1,4} Michael Niemeyer³ Stefano Gasperini^{1,4,5}
Hendrik P.A. Lensch² Nassir Navab¹ Federico Tombari^{1,3}
¹Technical University of Munich ²University of Tübingen
³Google ⁴Munich Center for Machine Learning ⁵VisualAIs

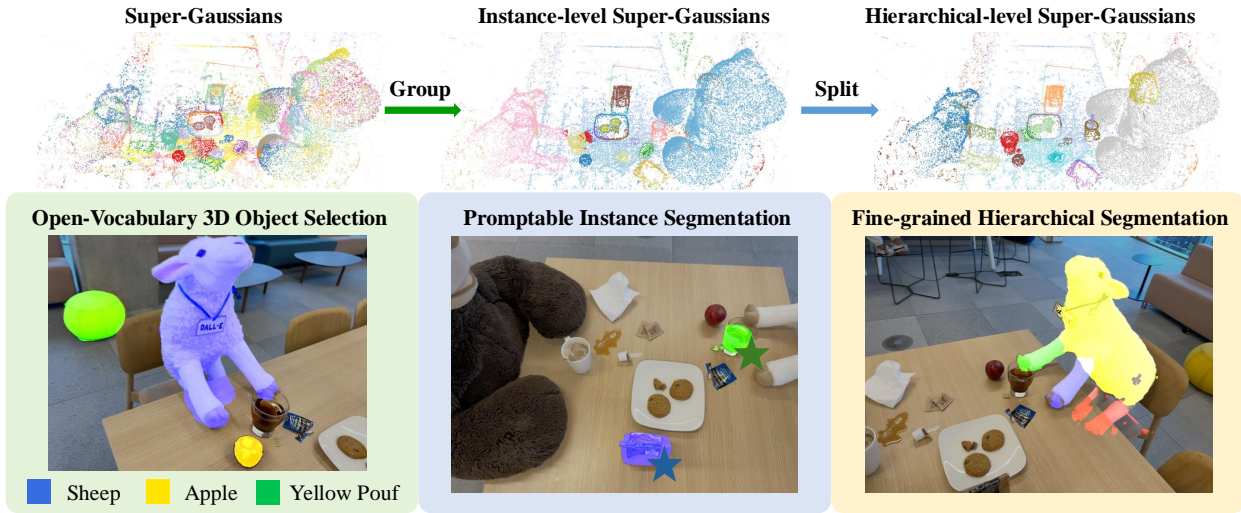


Figure 1. We present **SuperGSeg**, a novel method that clusters similar Gaussians into superpoint-like representations, termed Super-Gaussians (SuperGs). SuperGSeg enables efficient integration of diverse feature fields for comprehensive 3D scene understanding. **Left**: Querying SuperGs’ language features enables open-vocabulary 3D object selection, producing consistent 3D masks that extend beyond 2D visible surfaces, e.g., the leg of the sheep under the table. **Middle**: Grouping SuperGs by instance features enables promptable instance segmentation. **Right**: Further splitting instances via hierarchical features enables fine-grained hierarchical segmentation.

Abstract

3D Gaussian Splatting has recently gained traction for its efficient training and real-time rendering. While its vanilla representation is mainly designed for view synthesis, recent works extended it to scene understanding with language features. However, storing additional high-dimensional features per Gaussian for semantic information is memory-intensive, which limits their ability to segment and interpret challenging scenes. To this end, we introduce SuperGSeg, a novel approach that fosters cohesive, context-aware hierarchical scene representation by disentangling segmentation and language field distillation. SuperGSeg first employs neural 3D Gaussians to learn geometry, instance and hierarchical segmentation features from multi-view images

with the aid of off-the-shelf 2D masks. These features are then leveraged to create a sparse set of Super-Gaussians. Super-Gaussians facilitate the lifting and distillation of 2D language features into 3D space. They enable hierarchical scene understanding with high-dimensional language feature rendering at moderate GPU memory costs. Extensive experiments demonstrate that SuperGSeg achieves remarkable performance on both open-vocabulary object selection and semantic segmentation tasks. More results at supergseg.github.io.

1. Introduction

3D Gaussian Splatting (3DGS) [1] has rapidly emerged as a compelling alternative to NeRF [2] for its efficient training, real-time rendering, and explicit 3D representation. These advantages make 3DGS well-suited for a broad range of ap-

* Equal contribution.

plications, including 3D reconstruction [3–5], content generation [6], and scene understanding [7–12]. A particularly promising direction involves extending 3DGS frameworks to open-vocabulary understanding, enabling flexible, language-driven interaction with 3D scenes [13, 14].

Several recent methods aim to enable such open-vocabulary capabilities in 3DGS by distilling language features from both 2D [7, 9, 15, 16] and 3D [11, 12] perspectives. In 2D-based methods, language features extracted from images are lifted into 3D by exploiting the multi-view consistency inherent in 3DGS rendering. To reduce the substantial memory and computation overhead of storing and processing high-dimensional language features for each Gaussian, these methods employ dimensionality reduction techniques [7, 9]. However, this compression inevitably discards fine-grained semantic information. Another limitation is their inability to recognize partially occluded objects, which is often necessary in 3D understanding tasks. Text queries are performed on rendered pixels, which only capture the visible surface along each viewing ray. Consequently, objects that are partially or fully hidden cannot be retrieved. In contrast, 3D methods [11, 12] perform text queries directly in 3D space at the point level, which enables the retrieval of occluded objects by rendering the queried Gaussians into masks (see Figure 5), but also introduces new limitations. By directly associating language features with individual Gaussians and decoupling alpha blending, they cannot render consistent language feature maps in pixel space, which in turn makes them unsuitable for tasks such as pixel-wise dense semantic segmentation in 2D.

To address the aforementioned issues, we introduce a novel approach that: (1) preserves high-dimensional language feature embeddings without information loss, (2) handles occlusions by operating directly in 3D space, and (3) supports multi-granular segmentation, ultimately enabling open-vocabulary queries in both 2D and 3D, as shown in Figure 1. Inspired by superpoints [17] in point cloud analysis, our method clusters millions of Gaussians into a compact set of Super-Gaussians (SuperGs). However, due to the inherent noise in Gaussian point clouds, clustering solely based on Gaussian positions often produces suboptimal groupings. Instead, we leverage instance and hierarchical features extracted from grouped SAM masks [18] to guide clustering via an adaptive online clustering network [19]. For open-vocabulary scene understanding, we further distill 2D CLIP features [13] onto SuperGs that integrate both spatial and semantic information. This compact representation allows language features to be assigned at the SuperG level rather than to each individual Gaussian [7–9], thereby reducing the number of learnable language features from millions to only thousands, significantly lowering memory usage while retaining the full descriptive power of the original high-dimensional features.

Extensive experiments on the LERF-OVS [7] and ScanNet [20] datasets show that our method achieves remarkable performance in open-vocabulary 3D object retrieval and scene-level semantic segmentation, demonstrating superior capability in producing complete and consistent masks for 3D object retrieval and capturing fine-grained scene details for 2D dense pixel-wise segmentation. We summarize the main contributions as follows:

- We introduce SuperGSeg, a novel 3D scene understanding framework built on Super-Gaussian representations, enabling effective high-dimensional language feature distillation without information loss.
- We propose a novel neural Gaussian rasterization pipeline that distills instance and hierarchical feature fields, facilitating Super-Gaussian clustering and supporting multi-granular scene understanding.
- We design an online clustering network that adaptively fuses geometric, semantic, and appearance cues to generate Super-Gaussians, thus improving clustering quality.

2. Related Work

3D Open-Vocabulary Understanding. Advancements in universal 2D scene understanding, driven by foundation models such as CLIP [13] and SAM [21], have motivated the integration of language-aligned features into 3D scene representations. Early efforts incorporated these 2D features [13, 22] into NeRF-based representations [23, 24], enabling open-vocabulary queries in 3D scenes but at the cost of slow rendering and high memory usage. More recently, the emergence of 3DGS as a high-quality, real-time alternative for novel view synthesis has inspired extensions toward 3D scene understanding. For example, LangSplat [7] employs a scene-specific language autoencoder to compress high-dimensional CLIP features, providing clear object boundaries in rendered feature images while reducing memory usage. Feature3DGS [9] introduces a parallel Gaussian rasterizer with a lightweight convolutional decoder to distill high-dimensional features for tasks like scene editing and segmentation. However, these dimensionality reduction techniques inevitably discard fine-grained semantic information. OpenGaussian [11] instead directly associates uncompressed, lossless CLIP features with 3D Gaussians, preserving complete semantics and enabling the retrieval of visually occluded objects by performing queries directly in 3D space. Nevertheless, its decoupled language codebook design makes per-pixel 2D language feature rendering infeasible, thereby limiting performance on dense, pixel-wise semantic prediction tasks.

Despite notable progress, most existing methods focus primarily on instance-level knowledge while neglecting fine-grained part-level semantics [10, 11], or require separate models for different semantic granularities [7]. While recent methods [18, 25] explore hierarchical 3D understand-

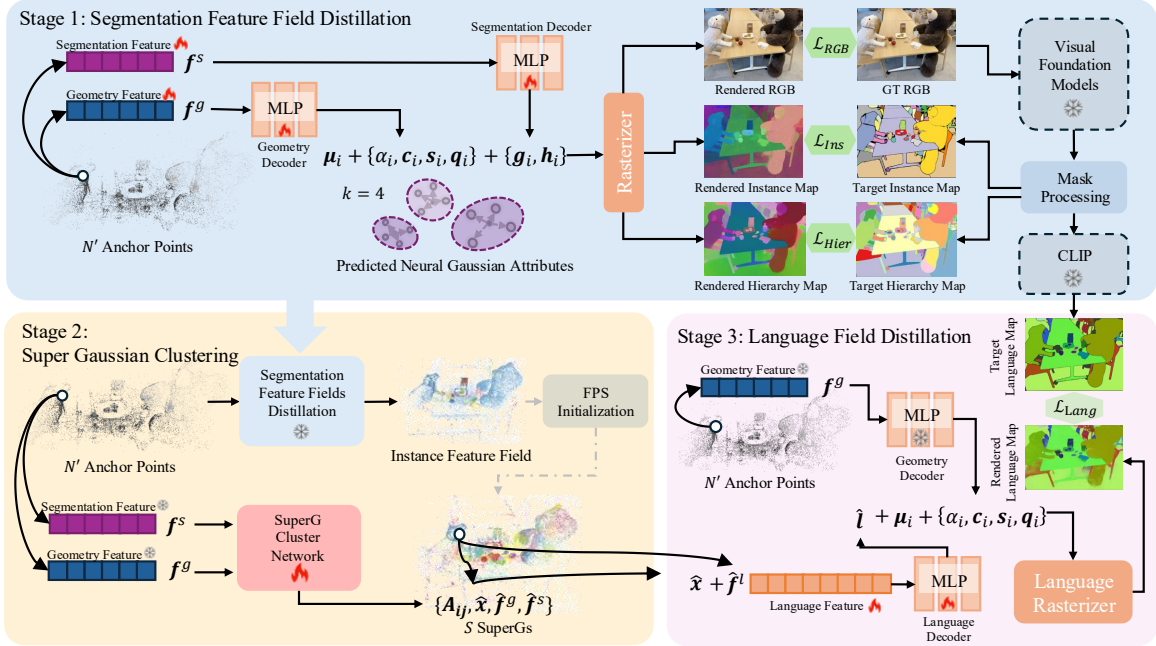


Figure 2. **SuperGSeg Overview.** We initialize the 3D Gaussians from a sparse set of anchor points, each generating k Gaussians with corresponding attributes. First, we train the appearance and segmentation features using RGB images and segmentation masks generated by SAM [21]. Next, we use the segmentation features and their spatial positions to produce a sparse set of Super-Gaussians, each carrying a 512-dimensional language feature. Finally, we train this high-dimensional language feature using a 2D feature map from CLIP [13].

ing at the part-level, they lack support for open-vocabulary language queries, leaving the joint modeling of multi-granularity 3D representation with language feature largely unexplored. In contrast, our method integrates both instance and hierarchical features from 2D inputs, and introduces a Super-Gaussian based language field that fuses segmentation information with the spatial distribution of 3D Gaussians, thereby enabling open-vocabulary, multi-granularity, and occlusion-robust 3D segmentation.

Superpoints. Superpoints have long served as fundamental primitives for various point cloud understanding tasks [19, 26–31]. Early approaches, such as Voxel Cloud Connectivity Segmentation (VCCS) [32], segment a voxelized 3D grid into spatially coherent regions using region-growing variants of K -means clustering. More recent works leverage learned point cloud representations [33, 34] to infer superpoints directly from 3D scans [19, 27, 35]. Superpoints have also been adopted for open-vocabulary 3D segmentation [17], demonstrating robustness in complex scenes. However, directly applying superpoint methods to 3DGS is challenging due to noisy Gaussian geometry. To address this, we leverage instance- and part-level cues from 2D foundation models to guide superpoint formation, effectively bridging high-quality 2D features with noisy 3D Gaussian representations.

3. Method

Given a set of posed RGB images, our goal is to reconstruct a 3D scene with a compact language feature field that supports open-vocabulary querying of arbitrary concepts. To achieve this, we propose a three-stage training paradigm, as shown in Figure 2. In the first stage, we train a neural variant of 3DGS [36] to reconstruct scene geometry using N' anchor points, each having a geometry feature f^g and a segmentation feature f^s . Anchor points are then spawned into a set of neural Gaussians and optimized. In the second stage, a learnable cluster network groups the anchors into S SuperGs using f^g , f^s , and anchor position \mathbf{x} , ensuring geometric and semantic consistency. Since $S \ll N'$, this yields a far more compact representation. In the third stage, we learn a language feature \hat{f}^l for each SuperG, enabling open-vocabulary queries on just S SuperGs rather than millions of individual Gaussians.

3.1. Preliminaries: Neural Gaussian Splatting

We begin with Stage 1 of our pipeline: modeling the scene geometry with Scaffold-GS [36] structure. Vanilla 3DGS represents a scene with N Gaussians, each parameterized by a center μ , opacity α , color \mathbf{c} , scale \mathbf{s} and quaternion \mathbf{q} . These Gaussians are projected onto the image plane [37] and rendered into RGB images via α -blending. While achiev-

ing leading rendering quality and speed, optimizing each Gaussian independently often leads to overfitting, redundancy, and degraded robustness in challenging regions such as texture-less surfaces. Scaffold-GS addresses these issues by voxelizing the scene into N' anchor points, each at position \mathbf{x} . From each anchor, k neural Gaussians are derived, where centers are computed as \mathbf{x} plus learnable offsets, and the remaining attributes $(\alpha, \mathbf{c}, \mathbf{s}, \mathbf{q})$ are produced on the fly from the anchor’s geometry feature \mathbf{f}^g via dedicated MLPs. By tying Gaussians to anchors, Scaffold-GS constrains their spatial distribution to the scene structure, preventing uncontrolled growth and improving robustness.

Training in 3DGS typically relies on a photometric loss \mathcal{L}_{RGB} , where rendered RGB images are supervised against ground-truth views. Unlike vanilla 3DGS that optimizes $(\mu, \alpha, \mathbf{c}, \mathbf{s}, \mathbf{q})_N$, with N often reaching millions for complex scenes, Scaffold-GS optimizes only $(\mathbf{f}^g)_{N'}$, the Gaussian offsets, and MLP weights, which significantly reduces parameters. This anchor-based formulation naturally yields a coarse partition of the Gaussian space, providing a strong basis for our subsequent clustering into SuperGs.

3.2. Segmentation Feature Field Distillation

Given N' anchor points representing the scene geometry, the next step is to group them into S superpoints, each forming a SuperG through its derived neural Gaussians. Ideally, each SuperG should align with a single semantic entity in the scene. However, clustering anchors solely by their geometry features \mathbf{f}^g or positions \mathbf{x} is suboptimal, since anchors from distinct objects can be spatially adjacent or geometrically similar. To overcome this limitation, we introduce an additional segmentation feature \mathbf{f}^s , distilled from 2D SAM masks, which encodes both instance- and part-level semantic cues to guide the SuperG clustering.

Hierarchical Partitioning of SAM Masks. Given an input RGB image, SAM [21] generates a set of 2D segmentation masks. These masks can, however, overlap with each other, leading to pixels belonging to multiple masks and thus obscuring the inherent part-instance hierarchy. Prior works either train separate models for each mask level [7, 38, 39], which is less efficient, or rely only on coarse instance-level masks [11, 40], discarding the finer part-instance relations. To overcome this, we adopt a hierarchical representation [18] that restructures the masks into non-overlapping instance-level masks \mathcal{M} for whole objects and part-level patches \mathcal{P} for finer components, which together provide supervision for learning both object-level semantics and intra-object details in the segmentation feature field. Implementation details and example mask visualizations are provided in Appendix B.

Instance and Hierarchical Feature Field. As shown in Figure 2, we assign each anchor point a segmentation feature \mathbf{f}^s . We pass \mathbf{f}^s together with the anchor position \mathbf{x}

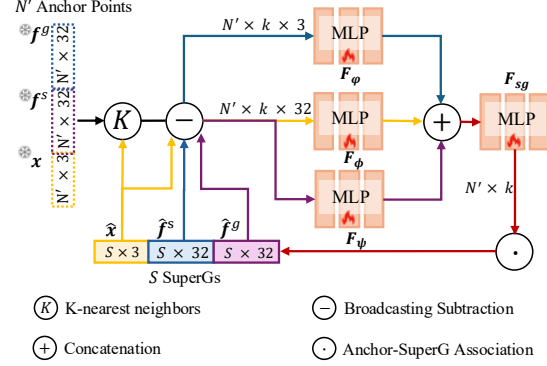


Figure 3. The architecture of the SuperG Cluster Network.

to a segmentation decoder to get the instance feature \mathbf{g} and hierarchical feature \mathbf{h} for each neural Gaussian. Through the vanilla Gaussian Splatting pipeline, we rasterize \mathbf{g} and \mathbf{h} to generate the 2D instance feature map $\hat{\mathbf{G}} \in \mathbb{R}^{D_g \times H \times W}$ and the 2D hierarchical feature map $\hat{\mathbf{H}} \in \mathbb{R}^{D_h \times H \times W}$.

To train the segmentation features, we leverage a contrastive learning objective [18, 41] to enforce cross-view consistency, encouraging features from the same mask to be similar while pushing apart those from different masks. Specifically, we represent the set of SAM-generated instance-level masks as $\mathcal{M} = \{\mathbf{m}^p \in \mathbb{R}^{H \times W} \mid p = 1, \dots, |\mathcal{M}|\}$. Given an instance mask \mathbf{m}^p , we collect all rendered instance features whose pixels fall inside the mask, and denote this set as $\hat{\mathbf{g}}^p = \{\hat{\mathbf{g}}_t^p \in \hat{\mathbf{G}} \mid t = 1, \dots, |\hat{\mathbf{g}}^p|\}$. We compute the mean instance feature value within \mathbf{m}^p as $\bar{\mathbf{g}}^p$ and the contrastive instance feature loss \mathcal{L}_{Ins} is:

$$\mathcal{L}_{Ins} = -\frac{1}{|\mathcal{M}|} \sum_{p=1}^{|\mathcal{M}|} \sum_{t=1}^{|\hat{\mathbf{g}}^p|} \log \frac{\exp(\hat{\mathbf{g}}_t^p \cdot \bar{\mathbf{g}}^p / \tau_p)}{\sum_{q=1}^{|\mathcal{M}|} \exp(\hat{\mathbf{g}}_t^p \cdot \bar{\mathbf{g}}^q / \tau_q)}, \quad (1)$$

where τ is the cluster temperature. We adopt a similar hierarchical feature loss \mathcal{L}_{Hier} from Omniseg3D [18], but applied to part-level patches \mathcal{P} to supervise our hierarchical feature \mathbf{h} . We refer to Appendix B for more details. Combined with the reconstruction loss introduced in Section 3.1, these objectives define the overall training loss for Stage 1:

$$\mathcal{L}_{stage1} = \mathcal{L}_{RGB} + \lambda_{Ins} \mathcal{L}_{Ins} + \lambda_{Hier} \mathcal{L}_{Hier}. \quad (2)$$

3.3. Super-Gaussian Clustering

After learning anchor-level geometry and segmentation features, we proceed to Stage 2, where anchors are grouped into semantically meaningful SuperGs to form a compact representation. However, contrastive learning struggles to separate objects that never co-occur in training [25], potentially grouping too distant Gaussians. To ensure spatial compactness and semantic consistency, we incorporate the anchor positions \mathbf{x} alongside segmentation features \mathbf{f}^s , while geometric features \mathbf{f}^g provide appearance cues for refine-

ment. A straightforward baseline is to apply K -means clustering [12] to the concatenated feature space of $\{\mathbf{x}, \mathbf{f}^g, \mathbf{f}^s\}$. Yet, this approach fails when appearance cues misalign with semantics (e.g., diverse textures within an object). Moreover, K -means assumes equal importance across concatenated features, without the flexibility to adapt their relative relevance during clustering. To improve the clustering quality, we instead propose a learnable SuperG clustering network (see Figure 3), inspired by [19]. It follows two steps: initialization and iterative refinement.

Super-Gaussian Initialization. We apply the Farthest Point Sampling algorithm [42] on anchor points to initialize SuperGs, averaging each a position $\hat{\mathbf{x}}$. Each SuperG has a geometry feature $\hat{\mathbf{f}}^g$ and segmentation feature $\hat{\mathbf{f}}^s$, which are initialized as the mean value of the corresponding anchors' features $\{\mathbf{f}^g, \mathbf{f}^s\}$.

Super-Gaussian Update. We denote the nearest k SuperGs to the i -th anchor as \mathcal{N}_i . The association probability matrix $\mathbf{A} \in \mathbb{R}^{N' \times k}$ [19, 43] is used to weight the contribution of each SuperG to its corresponding anchor, where N' is the number of anchors and k is the number of nearest SuperGs. Specifically, the association probability between the j -th SuperG ($j \in \mathcal{N}_i$) and the i -th anchor is:

$$\mathbf{A}_{ij} = F_{sg} \left(F_\phi(\mathbf{x}_i, \hat{\mathbf{x}}_j), F_\varphi(\mathbf{f}_i^s, \hat{\mathbf{f}}_j^s), F_\psi(\mathbf{f}_i^g, \hat{\mathbf{f}}_j^g) \right), \quad (3)$$

where F_ϕ , F_φ , and F_ψ are lightweight MLP decoders that output relevance weights in terms of spatial, semantic, and geometric information, respectively. The concatenated weights are then passed to the prediction decoder F_{sg} for the normalized association probability matrix prediction. Unlike K -means, this design dynamically adjusts the contribution of each SuperG to its corresponding anchor.

We iteratively update SuperGs through the association matrix \mathbf{A} . At iteration $t + 1$, each SuperG's position and features are updated with its corresponding anchors:

$$\hat{\mathbf{e}}_j^{t+1} = \frac{1}{\sum_{i=1}^{N'} \mathbb{I}(j \in \mathcal{N}_i) \mathbf{A}_{ij}^t} \sum_{i=1}^{N'} \mathbb{I}(j \in \mathcal{N}_i) \mathbf{A}_{ij}^t \mathbf{e}_i, \quad (4)$$

where \mathbb{I} denotes the indicator function, $\mathbf{e} \in \{\mathbf{x}, \mathbf{f}^g, \mathbf{f}^s\}$ are the anchor's attributes and $\hat{\mathbf{e}} \in \{\hat{\mathbf{x}}, \hat{\mathbf{f}}^g, \hat{\mathbf{f}}^s\}$ are SuperG's.

We optimize the SuperG clustering network to learn the association matrix \mathbf{A} , ensuring that the derived SuperG attributes $\hat{\mathbf{e}}$ accurately reconstruct the anchor attributes \mathbf{e} . Note that \mathbf{e} from Stage 1 (Section 3.2) are now frozen:

$$\mathcal{L}_{recon, \mathbf{e}} = \frac{1}{N'} \sum_{i=1}^{N'} \|\mathbf{e}_i - \sum_{j \in \mathcal{N}_i} \mathbf{A}_{ij} \hat{\mathbf{e}}_j\|. \quad (5)$$

However, anchors within the same SuperG may be semantically similar yet spatially distant, especially when contrastive learning fails to optimize instances that never co-occur in the same view. To enforce spatial coherence, we introduce a compactness objective:

$$\mathcal{L}_{compact, \mathcal{X}} = \frac{1}{S} \sum_{j=1}^S \sum_{\mathbf{x} \in \mathcal{X}_j} \|\mathbf{x} - \hat{\mathbf{x}}_j\|, \quad (6)$$

where \mathcal{X}_j is the set of anchors' position assigned to the j -th SuperG. This loss encourages assigned anchors to cluster around their SuperG center and avoid fragmentation.

3.4. Language Field Distillation

Building on the clustering from Stage 2 (Section 3.3), in Stage 3, we distill 2D CLIP features into our compact set of S SuperGs, rather than into millions of individual 3D Gaussians to enable open-vocabulary 3D scene understanding. This design ensures consistent, robust, and high-dimensional language representations, while avoiding the feature degradation typically caused by the lossy compression used in Gaussian-based distillation approaches.

Since all Gaussians within a SuperG are expected to share the same semantics, we assign each SuperG a learnable latent language feature $\hat{\mathbf{f}}^l$. As shown in Figure 2, this latent feature, together with the SuperG position $\hat{\mathbf{x}}$, is decoded by a language feature MLP F_L to produce a CLIP-aligned feature: $\hat{\mathbf{l}} = F_L(\hat{\mathbf{f}}^l, \hat{\mathbf{x}})$. We then modify the rasterizer to render a language feature map $\hat{\mathbf{L}}$, using $\hat{\mathbf{l}}$ and the anchor-SuperG association map \mathbf{A} . For supervision, instance masks obtained in Section 3.2 are encoded using the CLIP image encoder to produce target 2D CLIP features \mathbf{L} . The latent features $\hat{\mathbf{f}}^l$ and the decoder F_L are jointly optimized using a cosine similarity loss:

$$\mathcal{L}_{Lang} = 1 - \cos(\hat{\mathbf{L}}, \mathbf{L}). \quad (7)$$

4. Experiments

4.1. Experimental Setup

Datasets. We evaluate our method on the **open-vocabulary novel view semantic segmentation and object selection tasks** using the ScanNet v2 [20] and LERF-OVS [7] datasets. ScanNet v2 [20] includes posed RGB images and 2D semantic labels of indoor scenes. We randomly select 8 scenes from the dataset. These include a variety of indoor environments, e.g., living rooms, bedrooms, kitchens, and offices. For each scene, we split the data into a training set (composed of every 20th image from the original sequence) and a test set (derived from the intermediate images between the training set samples). For semantic segmentation, we specifically use the 20 object categories. LERF-OVS [7] consists of complex in-the-wild scenes captured with consumer-level devices, annotated with ground truth masks of textual queries to enable evaluation for open-vocabulary object selection tasks.

Baselines and Metrics. We compare our method with representative NeRF-based and 3DGS-based baselines, including LERF [23], LangSplat [7], LEGaussian [8], and

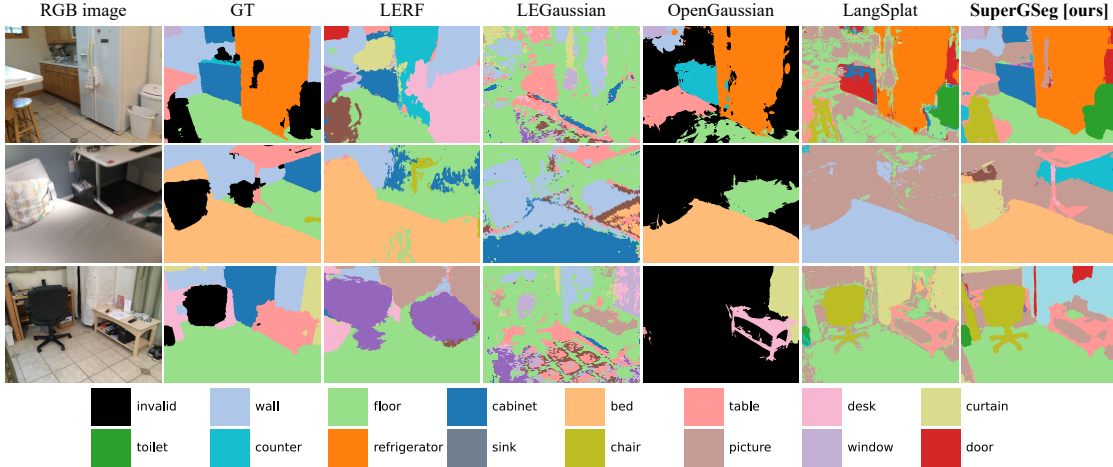


Figure 4. Qualitative comparison of semantic segmentation predictions on the ScanNet v2 dataset [20].

Method	mean		wall		floor		cabinet		chair		refrigerator		curtain	
	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
LERF [23]	38.5	60.4	35.2	82.8	60.1	68.8	52.0	82.7	10.9	10.9	69.9	90.2	70.2	77.8
LEGaussians [8]	8.7	33.2	17.9	53.1	14.6	20.6	2.7	18.6	0.4	28.7	9.0	74.3	1.9	10.4
OpenGaussian [11]	24.1	68.7	13.4	96.6	31.2	74.4	0.3	22.9	36.5	83.4	88.0	98.3	17.7	79.2
LangSplat [7]	27.6	48.3	45.3	72.6	43.3	45.6	24.8	56.7	18.0	48.5	0.7	33.3	46.8	66.5
SuperGSeg [ours]	54.7	74.7	58.8	92.9	53.6	86.5	69.8	83.8	80.4	83.8	79.4	80.2	61.8	64.5

Table 1. Comparison on the ScanNet v2 dataset [20]. We report the mean result and detailed scores for the most common object categories, following the evaluation protocol of [44]. Results for more categories are provided in the Appendix E.

OpenGaussian [11]. For the open-vocabulary semantic segmentation task, CLIP-encoded text features are compared with rendered 2D language feature maps via cosine similarity to produce per-pixel semantic predictions [9], evaluated with mean Intersection over Union (mIoU) and mean Accuracy (mAcc). For the open-vocabulary object selection task, we perform text queries directly in 3D space [11], retrieving the most relevant SuperGs and rendering them into 2D for evaluation with mIoU and mAcc. Since NeRF is an implicit representation without explicit 3D positions, LERF cannot be applied to this task. We also report inference-time efficiency, measuring both runtime and memory consumption for text queries on trained 3D scenes. Specifically, we perform multiple queries from different viewpoints and report the average query time. We consider this metric particularly important for assessing the feasibility of deploying models on resource-constrained devices and enabling real-time querying in practical scenarios.

Implementation Details. The training process is divided into 3 stages. In the first stage, we train the Scaffold-GS [36] with instance and hierarchical features for 30k iterations. In the second stage, we freeze the geometry and segmentation features from stage one and train only the SuperG clustering network for another 30k iterations. In the last stage,

we freeze all other parameters and optimize the language features for each SuperG for 10k iterations. For more implementation details, we refer to Appendix B.

4.2. Open-Vocabulary Semantic Segmentation

Quantitative Results. As shown in Table 1, SuperGSeg achieves the best overall scores in both mIoU and mAcc among the compared methods, demonstrating its effectiveness in capturing the open-set information of the scene, yielding remarkable performance in a variety of object categories. In comparison, LEGaussian [16] shows lower performance on both metrics, suggesting limited generalization across multiple object categories. LangSplat [7] performs better than LEGaussian but still shows reduced accuracy in more diverse categories. OpenGaussian [11] obtains competitive results on certain large structures such as wall and floor, but its overall scene-level performance remains below ours. LERF [23] achieves the second-highest mIoU, though its relatively low mAcc suggests difficulties in producing clear segmentation boundaries.

Qualitative Results. As shown in Figure 4, our method produces sharper and more semantically consistent masks than the compared methods. While OpenGaussian [11] demonstrates competitive performance in 3D object-level

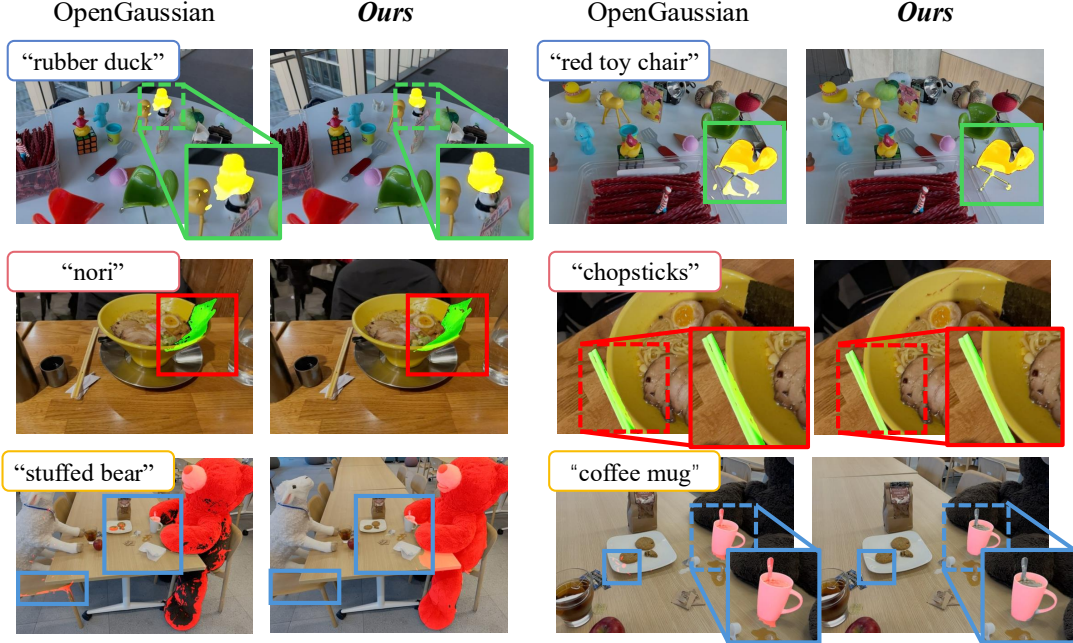


Figure 5. Qualitative comparison on the LERF-OVS dataset [23] for the open-vocabulary 3D object selection task. Text queries for each scene are displayed in quotation marks. SuperGSeg delivers more precise and less noisy segmentation masks.

Method	Inference		mean		<i>figurines</i>		<i>teatime</i>		<i>ramen</i>		<i>waldo_kitchen</i>	
	Time	Mem.	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
LangSplat [7]	3.28s	18GB	9.66	12.41	10.16	8.93	11.38	20.34	7.92	11.27	9.18	9.09
LEGaussians [8]	4.42s	5GB	16.21	23.82	17.99	23.21	19.27	27.12	15.79	26.76	11.78	18.18
OpenGaussian [11]	5.55s	9GB	38.36	51.43	39.29	55.36	60.44	76.27	31.01	42.25	22.70	31.82
SuperGSeg [ours]	0.50s	4GB	35.94	52.02	43.68	60.71	55.31	77.97	18.07	23.94	26.71	45.45

Table 2. Open-vocabulary 3D object selection comparison on the LERF-OVS dataset [7]. LERF [23] is not applicable for this task. We report the mIoU and mAcc of compared methods as provided in [11], and measure inference cost using their official implementations.

semantic segmentation (Section 4.3), it struggles in dense pixel-wise semantic segmentation. This is evident with occlusions due to projections onto 2D-pixel space. Without alpha blending, the occluded Gaussians cannot be effectively distinguished from one another. Instead, LangSplat [7] produces fine border segmentation but often includes incorrect semantic labels and noisy predictions, likely due to the lossy encoding of language information. LERF [23] presents accurate semantic prediction but with imprecise boundaries, limiting its applicability in fine-grained segmentation tasks.

4.3. Open-Vocabulary Object Selection

Quantitative Results. SuperGSeg improves over baseline methods that assign and optimize language features per Gaussian [7, 9, 16]. As shown in Table 2, clustering Gaussians into SuperGs enhances both spatial and semantic accuracy over per-Gaussian methods. We further compare SuperGSeg to OpenGaussian [11], another method exploring 3D Gaussian clustering. OpenGaussian’s direct 2D

CLIP feature association yields a slightly higher mIoU by avoiding alpha-blending artifacts, but it underperforms in 2D semantic segmentation on ScanNet (Section 4.2). In contrast, SuperGSeg maintains competitive mIoU for 3D object selection while surpassing OpenGaussian in 2D semantic segmentation, enhancing its versatility across real-world applications. Our higher mAcc, especially in complex LERF-OVS scenes such as *figurines* and *waldo kitchen*, reflects the precision of Super-Gaussian clustering and instance grouping. By accurately segmenting Gaussians in 3D, SuperGSeg renders more complete 2D masks with sharper boundaries, improving semantic consistency in challenging settings. In addition, SuperGSeg reduces inference latency to around 0.5s per query and decreases memory usage to 4GB, more than 50% lower than the next best baseline at 9GB. These improvements, enabled by SuperGs, demonstrate the potential for real-time querying on resource-constrained devices.

Qualitative Results. For visualization, we query language features in 3D space and render the resulting 3D masks

to 2D. As shown in Figure 5, SuperGSeg delivers precise 3D object selection without spurious outliers and produces clearer boundaries. Thanks to the 3D understanding capability, our SuperGSeg allows for effective localization of occluded regions (e.g., the *stuffed bear* leg under a table). Notably, its high-quality features distinguish the coffee mug from its contents and spoon, showcasing the efficacy of distilling fine-grained features into SuperGs.

Ablation Study. We conduct ablation studies on various components of our method to validate the necessity of SuperGs, as summarized in Table 3. The baseline without SuperG (case a) trains the language feature field by directly optimizing per-anchor features, which results in limited semantic consistency. To analyze how different feature types affect SuperG formation, we evaluate grouping based solely on anchor coordinates and geometric features (case b), instance features (case c), and hierarchical features (case d). The results indicate that grouping Gaussians into SuperG improves semantic consistency compared to per-anchor optimization, but relying only on coordinates and geometry remains suboptimal. Both instance and hierarchical features contribute substantially to accurate SuperG assignments, and the best performance is achieved with our full model (case f), which combines both. We further compare *K*-means clustering for Gaussian grouping (case e) with our learnable SuperG assignment (case f). By dynamically adapting to variations in the feature space, our learnable predictor produces higher-quality SuperGs, yielding consistently higher mIoU and improved mAcc. Additional ablation studies on components of the SuperG clustering network are provided in Appendix D.

#	w/ Learned SuperG	w/ ins	w/ hier	mIoU \uparrow	mAcc. \uparrow
a)				10.12	14.49
b)	✓			12.08	16.95
c)	✓	✓		53.91	64.41
d)	✓		✓	49.04	66.10
e)		✓	✓	53.77	67.80
f)	✓	✓	✓	55.31	77.97

Table 3. SuperG ablation study, *teatime* scene of LERF-OVS.

4.4. Application

Beyond language-based querying, SuperGs serve as a multi-granularity representation of 3D scenes by integrating instance- and part-level knowledge, readily applicable to tasks such as cross-frame segmentation and hierarchical instance decomposition, without requiring task-specific re-training. For example, a click on a reference image retrieves SuperGs with matching hierarchical features, allowing the selected part to be consistently rendered across views. In addition to cross-view querying, SuperGSeg enables cross-level queries: clicking on a part retrieves its parent object

using instance features, while clicking on an object reveals its constituent parts, which supports seamless navigation from parts to instances and vice versa, as illustrated in Figure 6. Furthermore, the granularity of instance-to-part segmentation can be adjusted by varying the threshold on hierarchical feature similarity, as shown in Figure 7. Additional implementation details are provided in Appendix A.

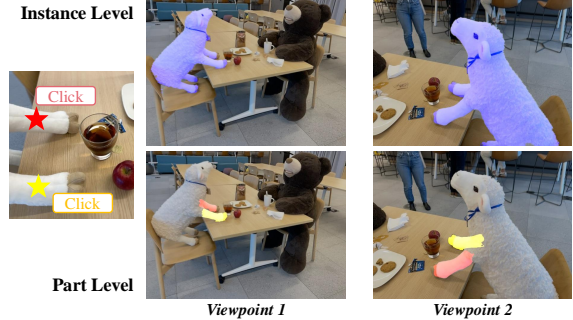


Figure 6. Cross-level and cross-frame segmentation visualization.

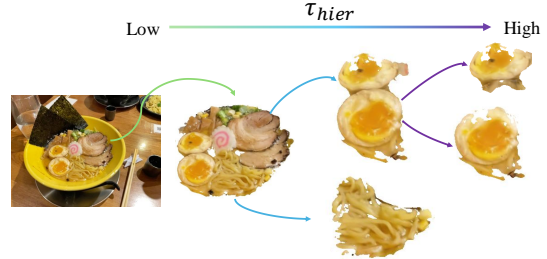


Figure 7. Visualization of intra-object hierarchy definition.

5. Conclusion

We present SuperGSeg, a novel framework for 3D scene understanding that represents scenes using compact Super-Gaussians, ensuring semantic and appearance consistency. By leveraging neural Gaussians, our method captures instance- and part-level segmentation features, guiding Super-Gaussian clustering through an adaptive online learning algorithm. Experiments show that integrating high-dimensional language features significantly improves open-set 3D language querying, demonstrating the framework’s remarkable performance. Furthermore, the Super-Gaussian representation is readily adaptable to a wide range of 3D scene understanding tasks.

6. Acknowledgement

This work was supported by the German Research Foundation (DFG): SFB 1233, Robust Vision: Inference Principles and Neural Mechanisms, TP A1, project number: 276693517, and the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Siyun Liang.

References

- [1] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D Gaussian Splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):139–1, 2023. [1](#)
- [2] 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 *Proceedings of the European Conference on Computer Vision*, pages 405–421, 2020. [1](#)
- [3] Zehao Yu, Torsten Sattler, and Andreas Geiger. Gaussian opacity fields: Efficient and compact surface reconstruction in unbounded scenes. *arXiv preprint arXiv:2404.10772*, 2024. [2](#)
- [4] Pinxuan Dai, Jiamin Xu, Wenxiang Xie, Xinguo Liu, Huamin Wang, and Weiwei Xu. High-quality surface reconstruction using gaussian surfels. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024.
- [5] 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. [2](#)
- [6] Taoran Yi, Jiemin Fang, Zanwei Zhou, Junjie Wang, Guan-jun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Xinggang Wang, and Qi Tian. Gaussiandreamerpro: Text to manipulable 3D gaussians with highly enhanced quality. *arXiv preprint arXiv:2406.18462*, 2024. [2](#)
- [7] 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. [2](#), [4](#), [5](#), [6](#), [7](#), [3](#)
- [8] 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. [5](#), [6](#), [7](#), [4](#)
- [9] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suyu 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. [2](#), [6](#), [7](#)
- [10] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3D scenes. In *Proceedings of the European Conference on Computer Vision*, pages 162–179, 2024. [2](#)
- [11] 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. [2](#), [4](#), [6](#), [7](#), [3](#)
- [12] Haijie Li, Yanmin Wu, Jiarui Meng, Qiankun Gao, Zhiyao Zhang, Ronggang Wang, and Jian Zhang. Instancegaussian: Appearance-semantic joint gaussian representation for 3d instance-level perception. *arXiv preprint arXiv:2411.19235*, 2024. [2](#), [5](#)
- [13] 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 *Proceedings of the International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [2](#), [3](#)
- [14] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022. [2](#)
- [15] Xingxing Zuo, Pouya Samangouei, Yunwen Zhou, Yan Di, and Mingyang Li. Fmgs: Foundation model embedded 3D gaussian splatting for holistic 3D scene understanding. *International Journal of Computer Vision*, pages 1–17, 2024. [2](#)
- [16] Yanyan Li, Chenyu Lyu, Yan Di, Guangyao Zhai, Gim Hee Lee, and Federico Tombari. Geogaussian: Geometry-aware gaussian splatting for scene rendering. In *European Conference on Computer Vision*, pages 441–457. Springer, 2024. [2](#), [6](#), [7](#)
- [17] Phuc Nguyen, Tuan Duc Ngo, Evangelos Kalogerakis, Chuang Gan, Anh Tran, Cuong Pham, and Khoi Nguyen. Open3Dis: Open-vocabulary 3D instance segmentation with 2d mask guidance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4018–4028, 2024. [2](#), [3](#)
- [18] Haiyang Ying, Yixuan Yin, Jinzhi Zhang, Fan Wang, Tao Yu, Ruqi Huang, and Lu Fang. Omnise3D: Omniversal 3D segmentation via hierarchical contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20612–20622, 2024. [2](#), [4](#)
- [19] Le Hui, Jia Yuan, Mingmei Cheng, Jin Xie, and Jian Yang. Superpoint network for point cloud oversegmentation. In *ICCV*, 2021. [2](#), [3](#), [5](#), [1](#)
- [20] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2017. [2](#), [5](#), [6](#), [4](#)
- [21] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. [2](#), [3](#), [4](#)
- [22] 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*, 2021. [2](#)
- [23] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. [2](#), [5](#), [6](#), [7](#), [4](#)

- [24] Francis Engelmann, Fabian Manhardt, Michael Niemeyer, Keisuke Tateno, Marc Pollefeys, and Federico Tombari. Opennerf: open set 3d neural scene segmentation with pixel-wise features and rendered novel views. *arXiv preprint arXiv:2404.03650*, 2024. 2
- [25] Chung Min Kim, Mingxuan Wu, Justin Kerr, Ken Goldberg, Matthew Tancik, and Angjoo Kanazawa. Garfield: Group anything with radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21530–21539, 2024. 2, 4
- [26] 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. 3
- [27] Damien Robert, Hugo Raguét, and Loic Landrieu. Efficient 3D semantic segmentation with superpoint transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17195–17204, 2023. 3
- [28] Damien Robert, Hugo Raguét, and Loic Landrieu. Scalable 3D panoptic segmentation as superpoint graph clustering. In *2024 International Conference on 3D Vision (3DV)*, pages 179–189. IEEE, 2024.
- [29] Loic Landrieu and Guillaume Obozinski. Cut pursuit: Fast algorithms to learn piecewise constant functions on general weighted graphs. *SIAM Journal on Imaging Sciences*, 10(4):1724–1766, 2017.
- [30] Yun Zhu, Le Hui, Yaqi Shen, and Jin Xie. SPGroup3D: Superpoint grouping network for indoor 3D object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 7811–7819, 2024.
- [31] Mingmei Cheng, Le Hui, Jin Xie, and Jian Yang. Sspcnet: Semi-supervised semantic 3D point cloud segmentation network. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 1140–1147, 2021. 3
- [32] Jeremie Papon, Alexey Abramov, Markus Schoeler, and Florentin Worgötter. Voxel cloud connectivity segmentation-supervoxels for point clouds. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2027–2034, 2013. 3
- [33] 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. 3
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [35] Loic Landrieu and Mohamed Boussaha. Point cloud over-segmentation with graph-structured deep metric learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7440–7449, 2019. 3
- [36] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20654–20664, 2024. 3, 6, 2
- [37] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa splatting. *IEEE Transactions on Visualization and Computer Graphics*, 8(3):223–238, 2002. 3
- [38] Jiahuan Cheng, Jan-Nico Zaech, Luc Van Gool, and Danda Pani Paudel. Occam’s lgs: A simple approach for language gaussian splatting. *arXiv preprint arXiv:2412.01807*, 2024. 4
- [39] Yansong Qu, Shaohui Dai, Xinyang Li, Jianghang Lin, Liujuan Cao, Shengchuan Zhang, and Rongrong Ji. Goi: Find 3d gaussians of interest with an optimizable open-vocabulary semantic-space hyperplane. *arXiv preprint arXiv:2405.17596*, 2024. 4
- [40] Kim Jun-Seong, Kim GeonU, Kim Yu-Ji, Yu-Chiang Frank Wang, Jaesung Choe, and Tae-Hyun Oh. Dr. splat: Directly referring 3d gaussian splatting via direct language embedding registration. In *CVPR*, 2025. 4
- [41] Myrna C Silva, Mahtab Dahaghin, Matteo Toso, and Alessio Del Bue. Contrastive gaussian clustering: Weakly supervised 3d scene segmentation. *arXiv preprint arXiv:2404.12784*, 2024. 4
- [42] 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. 5
- [43] Diwen Wan, Ruijie Lu, and Gang Zeng. Superpoint gaussian splatting for real-time high-fidelity dynamic scene reconstruction. In *Proceedings of the 41st International Conference on Machine Learning*, pages 49957–49972, 2024. 5
- [44] Haoran Chen, Kenneth Blomqvist, Francesco Milano, and Roland Siegwart. Panoptic vision-language feature fields. *IEEE Robotics and Automation Letters*, 9(3):2144–2151, 2024. 6
- [45] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 2
- [46] 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

SuperGSeg: Open-Vocabulary 3D Segmentation with Structured Super-Gaussians

Supplementary Material

This supplementary document provides additional details about our method. Section A elaborates on the design of the Super-Gaussian and demonstrates its application to downstream tasks. Section B provides further implementation details, including the MLP architectures for neural Gaussian feature decoding and the adaptation of OpenGaussian for 2D open-vocabulary semantic segmentation comparison. Section C reports detailed efficiency analysis, including training time, inference speed, and memory consumption. Section D reports extended ablation studies on SuperG hyperparameters and module variants. Section E presents additional quantitative and qualitative results. Lastly, Section F discusses the limitations of our approach and outlines potential directions for future work.

A. Super-Gaussian Details

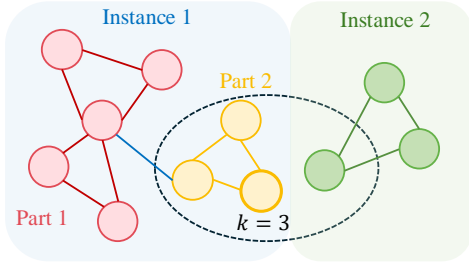


Figure 8. Example of Super-Gaussian graph. Each node represents a SuperG, connected to its k -nearest neighbors based on similarity in instance features. Through connected component analysis, the SuperG nodes are divided into two distinct instances. Within Instance 1, SuperG nodes are interconnected by similar hierarchical features, further splitting into two parts.

Module Design. As shown in Figure 3, our SuperG clustering network consists of four learnable MLPs. Inspired by SPNet [19], we design three attribute-specific learnable functions F_ϕ , F_ψ , and F_χ , each implemented as a single-hidden-layer MLP with ReLU activation. These functions independently encode the differences between an anchor and its k -nearest SuperGs in coordinates, segmentation features, and geometry, respectively, producing embeddings that reflect the relevance of each attribute for SuperG assignment. A final MLP, F_{sg} , then concatenates these embeddings and integrates spatial, semantic, and geometric cues into a probabilistic assignment.

Grouping Super-Gaussians for Instance and Hierarchical Segmentation. After training the SuperG association modules, we obtain a soft association map $\mathbf{A} \in \mathbb{R}^{N' \times k}$. During inference, each anchor point is assigned to one of its k -nearest neighbors with the highest probability, leading to a hard SuperG assignment $\hat{\mathbf{A}} \in \mathbb{R}^{N' \times 1}$. The attributes of each SuperG are then computed by averaging the attributes of its assigned anchors.

These SuperGs serve as the fundamental units for representing and interpreting the 3D scene. Specifically, as shown in Figure 8, we further construct a graph where nodes correspond to SuperGs. For instance segmentation, a node is connected to nodes within its k -nearest neighbors if their instance feature similarity exceeds a threshold τ_{Ins} . Instances are then obtained via connected component analysis on this restricted graph. Similarly, part segmentation is achieved by building a SuperG graph within each instance and identifying connected components based on hierarchical feature similarity with threshold τ_{Hier} . In practice, we set $k = 3$, $\tau_{Ins} = 0.8$, and $\tau_{Hier} = 0.9$.

B. Additional Implementation Details

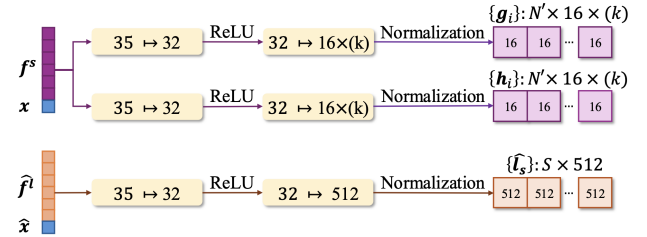


Figure 9. MLP structures for decoding different features.

Decoding Neural Gaussians from MLPs. We employ MLPs to decode latent features, as shown in Figure 9. Each MLP contains a single hidden layer of dimension 32. Their decoding targets, however, differ: instance feature decoder F_I and hierarchical feature decoder F_H decode anchor-level features, while language feature decoder F_L decodes SuperG-level features. Specifically, F_I and F_H take the anchor segmentation feature f^s and anchor position x as input, and predict the instance feature g and hierarchical feature h of the neural Gaussians spawned per anchor. In contrast, F_L predicts the CLIP-aligned feature \hat{t} for each SuperG, conditioned on its latent language feature f^l and its center \hat{x} .

Additional Information on the Contrastive Losses. The contrastive instance feature loss \mathcal{L}_{Ins} , which is computed on the set of SAM-generated instance-level masks \mathcal{M} and the rendered 2D instance feature map \hat{G} , has been discussed in Section 3.2. Inspired by [18], we elaborate a similar yet more complicated contrastive hierarchical feature loss \mathcal{L}_{Hier} . This loss is defined on the part-level patches $\mathcal{P} = \{\mathbf{p}^p \in \mathbb{R}^{H \times W} \mid p = 1, \dots, |\mathcal{P}|\}$ and the rendered 2D hierarchical feature map $\hat{H} \in \mathbb{R}^{D_h \times H \times W}$.

Given a patch \mathbf{p}^p , we collect the rendered hierarchical features from \hat{H} at each pixel, forming a set of feature vectors \mathbf{h}^p . The mean feature of the patch is then defined as $\bar{\mathbf{h}}^p$. We define the contrastive hierarchical feature loss at the minimum unit, on a pixel t with feature $\mathbf{h}_t^p \in \mathbf{h}^p$, as:

$$\mathcal{L}^{p,t}(r) = -\log \frac{\exp(\mathbf{h}_t^p \cdot \bar{\mathbf{h}}^r / \tau_r)}{\sum_{q=1}^{|\mathcal{P}|} \exp(\mathbf{h}_t^p \cdot \bar{\mathbf{h}}^q / \tau_q)}, \quad (8)$$

where τ is the temperature of the contrastive loss, and p, r, q are indices of patches. Subsequently, the hierarchical feature loss [18] can be written as:

$$\mathcal{L}_{Hier} = \sum_{p=1}^{|\mathcal{P}|} \sum_{d=1}^{d_{\max}^p} \mathcal{L}_{p,d}, \quad (9)$$

$$\mathcal{L}_{p,d} = \frac{\lambda^{d-1}}{|\mathcal{R}_d^p|} \sum_{t=1}^{|\{\mathbf{h}^p\}|} \sum_{r \in \mathcal{R}_d^p} \max(\mathcal{L}^{p,t}(r), \mathcal{L}_{\max}^{p,t}(d-1)), \quad (10)$$

where λ^{d-1} is a hyperparameter, \mathcal{R}_d^p denotes the index set of patches at hierarchy level d of patch p , and $r \in \mathcal{R}_d^p$ refers to a patch at level d . The maximum loss at level d ensures that the contrastive loss between the pixel feature t and patches with higher correlation (lower d) is always smaller than for patches with lower correlation:

$$\mathcal{L}_{\max}^{p,t}(d) = \max_{r \in \mathcal{R}_d^p} \mathcal{L}^{p,t}(r). \quad (11)$$

Additional Technical Details. We use the SAM ViT-H model [21] to generate 2D masks from the input images and then extract language features for each instance mask using the OpenCLIP ViT-B/16 model following [7]. The training process is divided into three stages. In the first stage, we train the Scaffold-GS [36] with instance and hierarchical feature attributes for 30k iterations. In the second stage, we freeze the geometry and multi-granularity features network from stage one and train only the SuperG clustering network for another 30k iterations. Finally, in the last stage, we freeze all other parameters and optimize the language features for each SuperG for 10k iterations. The embedding dimensions for \mathbf{f}^g and \mathbf{f}^s are set to 32 [36], while instance and hierarchical features are 16-dimensional [18]. For optimization, we use the Adam [45] optimizer for the MLPs with an initial learning rate of 0.01 and an exponential annealing schedule

of 0.001 as in [46].

OpenGaussian Implementation. OpenGaussian [11] assigns language features to instance-level Gaussians, enabling direct language queries on 3D point clouds. However, this approach does not natively support 2D pixel-level semantic segmentation, making direct evaluation on ScanNet more challenging. To enable a fair comparison, we first identify category-relevant 3D Gaussians by iterating over all text prompts to predict language feature maps for open-vocabulary semantic segmentation. For each instance-level Gaussian cluster, we determine the corresponding text prompt ID and store these IDs in a label map, which is then used to generate the final semantic segmentation. By following this approach, occlusions at the instance-level Gaussians are not explicitly handled, leading to the occlusion artifacts observed in Figure 4.

C. Training and Inference Efficiency

In Table 5, we report training time, inference time and memory consumption for LangSplat [7] and OpenGaussian [11] on the LERF-OVS dataset.

For LangSplat [7], the training consists of two stages: in S1, the 3DGS is pretrained without any additional feature fields for 30k iterations. In S2, the pretrained 3DGS is frozen, and a language feature field is optimized for another 30k iterations. For OpenGaussian [11], S1 corresponds to pretraining the 3DGS jointly with an instance feature field for 40k iterations. In S2, Gaussian clustering is performed in a coarse-to-fine manner, requiring an additional 30k iterations. Finally, in S3, the 2D language features are directly associated with the 3D Gaussian clusters.

When comparing our method to LangSplat and OpenGaussian across different training stages S1, S2, and S3, we find that our approach, while requiring longer training times in S1 and S2, achieves comparable efficiency in S3, indicating a trade-off due to the joint learning of instance and hierarchical features. GPU memory usage varies, with our method consuming more resources in S1 and S2 but significantly less in S3, showcasing the compact memory footprint enabled by our proposed SuperGs. Conversely, LangSplat exhibits consistent memory efficiency in S1 and S2, while OpenGaussian’s memory demands vary across different training stages. Most notably, our method consistently outperforms both baselines in inference speed across all scenarios, a critical advantage for real-time applications. This blend of rapid inference and flexible resource utilization highlights our method’s robustness for practical deployment in diverse computer vision tasks. In addition, our method supports multi-granularity scene understanding. Specifically, it enables semantic-level queries to retrieve groups of objects sharing the same language description, instance-level segmentation of a specific object, and further decom-

Ablation	mean		<i>figurines</i>		<i>teatime</i>		<i>ramen</i>		<i>waldo_kitchen</i>	
	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
$S = 250$	25.75	39.42	22.70	35.71	35.93	50.85	15.44	21.13	28.92	50.00
$S = 500$	34.05	50.03	27.44	48.21	57.42	77.97	20.42	23.94	30.93	50.00
$S = 1000$	35.94	52.02	43.68	60.71	55.31	77.97	18.07	23.94	26.71	45.45
$S = 2000$	27.61	40.60	27.64	46.43	47.58	62.71	12.46	16.90	22.76	36.36
$k = 3$	35.94	52.02	43.68	60.71	55.31	77.97	18.07	23.94	26.71	45.45
$k = 5$	33.70	46.76	21.59	35.71	68.75	84.75	16.96	21.13	27.48	45.45
$k = 10$	33.81	46.88	43.56	62.50	41.56	55.93	21.82	28.17	28.31	40.91

Table 4. Additional ablation studies on the LERF-OVS dataset [7] about the parameter choices for the SuperG Clustering Network. We use $s = 1000$ SuperGs and $k = 3$ for k -nearest neighbor in our implementation by default.

	LangSplat		OpenGaussian			Ours		
	S1	S2	S1	S2	S3	S1	S2	S3
Train	20m	45m	50m	20m	10m	90m	85m	30m
Memory	8G	6G	14G	17G	22G	18G	14G	4G
Inference	3.28s		5.55s			0.56s		
Memory	18GB		9GB			4GB		

Table 5. Comparison of the time and GPU memory requirements during training (top rows) and inference (bottom rows).

position of this instance into fine-grained parts. In contrast, LangSplat [7] only supports semantic-level queries. OpenGaussian [11] extends to instance-level understanding but, similar to LangSplat, does not support finer-grained part segmentation. This demonstrates that our method provides a more comprehensive representation for 3D scene understanding.

D. Additional Ablation Studies

In Section 4, we perform ablation studies to evaluate the necessity of SuperGs and analyze the performance of the instance feature field and hierarchical feature field. In this section, we further investigate the necessity of our proposed SuperG clustering network and measure the impact of its individual components.

Super-Gaussian Clustering Approaches. Given a pre-trained scene, our objective is to group anchors into meaningful SuperGs using their coordinates, segmentation features, and geometric properties. We explore two alternative approaches. First, we evaluate a simple K -means clustering algorithm by concatenating the aforementioned attributes and clustering them into $k = 1000$ SuperGs. Second, we experiment with a traditional supervoxel generation approach. Specifically, we map the anchors to a point cloud, using the concatenated features as normals. We then apply the Voxel Cloud Connectivity Segmentation (VCCS) algorithm [32] to compute the SuperGs. Finally, we compare these two non-

learning-based approaches with our learning-based method for Anchor-to-SuperG association.

We observed that K -means fails to prevent the overlap of resulting SuperGs across instances. Meanwhile, VCCS [32], originally designed for dense point clouds, struggles with the sparse structure of Gaussians. Its region-growing mechanism incorrectly clusters a large number of anchors together, which hinders the learning of the language feature field. The results in Table 6 show that our method is better suited for grouping Gaussians, achieving better performance.

Method	mIoU \uparrow	mAcc. \uparrow
K -means	53.77	67.80
VCCS[32]	0.45	0.00
Ours	55.31	77.97

Table 6. Ablation study of SuperG clustering approaches on the *teatime* scene of LERF-OVS.

w/ F_ϕ	w/ F_φ	w/ F_ψ	mIoU \uparrow	mAcc. \uparrow
\checkmark	\checkmark	\checkmark	32.41	40.68
			48.29	67.80
			58.07	75.66
\checkmark	\checkmark	\checkmark	37.12	62.71
		\checkmark	55.31	77.97

Table 7. Ablation study on the SuperG clustering network and its components on the *teatime* scene of LERF-OVS.

Super-Gaussian Clustering Network. As introduced in Section A, we employ three MLPs F_ϕ , F_φ , and F_ψ to capture the coordinate, segmentation, and geometric relationships between anchors and their k -nearest neighbors. To evaluate the contributions of these MLPs, we conduct an ablation study. Notably, in experiments where none of these MLPs are used, we directly concatenate the differences of the attributes as input to F_{sg} for predicting the association matrix.

Method	mIoU mean	mAcc	mIoU wall	mAcc	mIoU floor	mAcc	mIoU cabinet	mAcc	mIoU table	mAcc	mIoU desk	mAcc	mIoU curtain	mAcc
LERF [23]	38.5	60.4	35.2	82.8	60.1	68.8	52.0	82.7	10.2	80.1	14.4	16.1	70.2	77.8
LEGaussians [8]	8.7	33.2	17.9	53.1	14.6	20.6	2.7	18.6	0.0	0.0	0.5	13.5	1.9	10.4
OpenGaussian [11]	24.1	68.7	13.4	96.6	31.2	74.4	0.3	22.9	0.1	1.0	30.6	35.6	17.7	79.2
LangSplat [7]	27.6	48.3	45.3	72.6	43.3	45.6	24.8	56.7	21.9	87.4	0.1	6.4	46.8	66.5
SuperGSeg [ours]	54.7	74.7	58.8	92.9	53.6	86.5	69.8	83.8	35.7	54.8	15.0	16.7	61.8	64.5
	toilet		counter		refrigerator		chair		sink		window		door	
LERF [23]	25.2	25.2	24.4	42.8	69.9	90.2	10.9	10.9	25.8	37.1	11.5	11.5	64.5	67.5
LEGaussians [8]	13.7	16.3	10.7	27.0	9.0	74.3	0.4	28.7	0.3	0.4	0.0	44.4	1.4	4.7
OpenGaussian [11]	73.0	98.4	3.0	9.3	88.0	98.3	36.5	83.4	3.0	3.7	75.0	88.8	75.4	97.0
LangSplat [7]	0.1	5.4	10.7	34.7	0.7	33.3	18.0	48.5	0.0	0.0	0.0	0.1	55.6	66.3
SuperGSeg [ours]	26.9	26.9	14.0	59.1	79.4	80.2	80.4	83.8	11.7	12.0	54.7	77.0	58.2	58.3

Table 8. Comparison of mIoU and mAcc for various methods on each class of the ScanNet v2 dataset [20].

The results presented in Table 7 demonstrate that each MLP contributes to improving the SuperG assignments. The MLP F_ϕ , which accounts for the segmentation feature differences between the anchor and the SuperG, has the most significant impact. In particular, using only F_ϕ yields a relatively high mIoU, emphasizing its effectiveness in aligning semantic features. However, our full setup that integrates the F_ϕ and F_ψ for coordinate and geometric feature information further enhances mAcc. This suggests that incorporating additional spatial and geometric context refines the SuperG assignments, leading to a more precise understanding of the scene.

Parameters in Super-Gaussian Clustering Network.

We conduct ablation studies on the parameters involved in generating SuperGs using the SuperG clustering network. One crucial parameter is the total number of SuperGs predefined, denoted as S . Too few SuperGs fail to distinguish all instances, causing a single SuperG to span multiple instances, which undermines semantic accuracy. Conversely, too many SuperGs may introduce additional noise. Another parameter is the number of neighboring SuperGs k considered for each anchor when computing the association matrix between anchors and SuperGs.

As shown in Table 4, these two parameters are highly scene-specific, with the optimal number of SuperGs S and neighbors k varying across different scenes. Notably, for fair comparisons, we use the same parameter values, $S = 1000$ and $k = 3$, for all scenes. This parameter choice achieves optimal performance on average.

E. Additional Results

Qualitative Results in Occlusion Cases. As illustrated in Figure 10, our method queries objects directly in 3D space, effectively mitigating occlusion issues (e.g., the bear leg under the table can be retrieved). Moreover, the queried

objects exhibit multi-view consistency, enabling comprehensive scene understanding in 3D. For additional qualitative results, please refer to the accompanying videos.

Additional Quantitative Results on ScanNet. We report the results on more categories in the ScanNet dataset in Table 8.



Figure 10. The language-queried 3D masks rendering to arbitrary viewpoints remain multi-view consistent. Benefit from the 3D understanding, we enable render regions that are originally occluded and invisible in 2D.

F. Limitations

Despite the advancements achieved by our method, certain limitations remain. First, our approach inherits biases from the original visual foundation models, which may constrain performance and limit generalization to diverse or unseen scenarios. Second, our method is tailored for scene-specific language representation, requiring significant modeling time for each scene. This limits its applicability in tasks that demand rapid adaptation or broad generalization, such as in-the-wild scene understanding. Future work could focus on mitigating inherited biases and optimizing training pipelines to enhance scalability and generalization.