# Training-Free Hierarchical Scene Understanding for Gaussian Splatting with Superpoint Graphs

Shaohui Dai*    Yansong Qu*    Zheyan Li    Xinyang Li    Shengchuan Zhang    Liujuan Cao†

daish@stu.xmu.edu.cn,caoliujuan@xmu.edu.cn

Key Laboratory of Multimedia Trusted Perception and Efficient Computing,

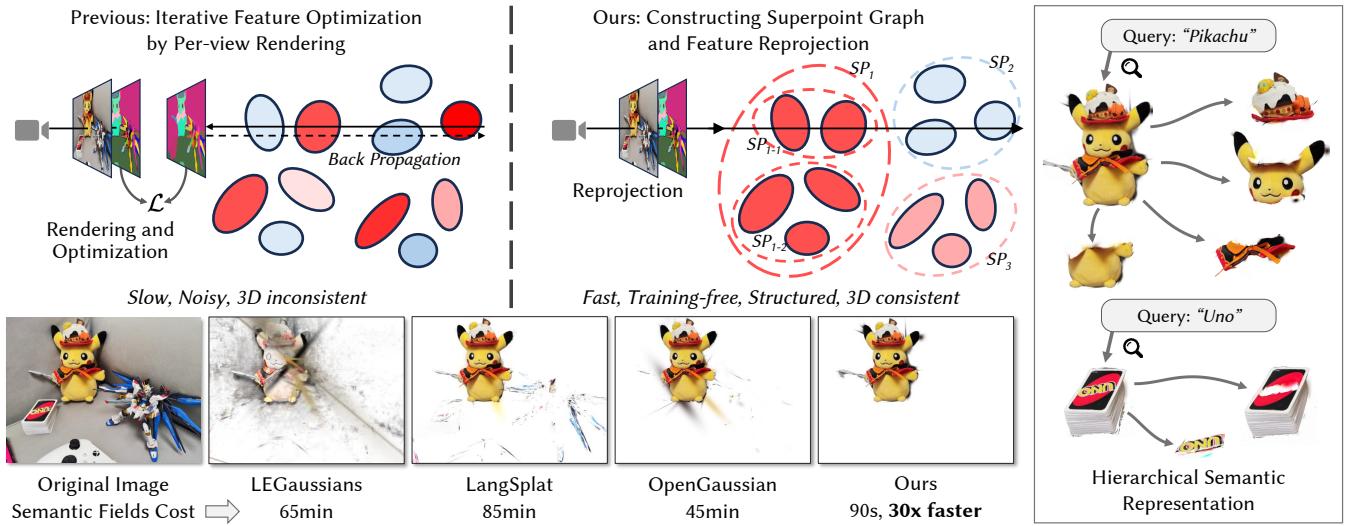Ministry of Education of China, Xiamen University, China

**Figure 1: We propose a method for open-vocabulary 3D scene understanding by integrating Gaussian Splatting with a superpoint graph. The top row shows that existing methods iteratively optimize features over Gaussian primitives, resulting in slow processing and inconsistent 3D semantics. In contrast, our approach clusters Gaussians into superpoints (denoted as *SP*) and uses feature reprojection to build a semantic field. The bottom row highlights 3D consistency and speed improvements. The examples on the right demonstrate the capability of our work for open-vocabulary and hierarchical scene understanding.**

## Abstract

Bridging natural language and 3D geometry is a crucial step toward flexible, language-driven scene understanding. While recent advances in 3D Gaussian Splatting (3DGS) have enabled fast and high-quality scene reconstruction, research has also explored incorporating open-vocabulary understanding into 3DGS. However, most existing methods require iterative optimization over per-view 2D semantic feature maps, which not only results in inefficiencies but also leads to inconsistent 3D semantics across views. To address these limitations, we introduce a training-free framework that constructs a superpoint graph directly from Gaussian primitives. The superpoint graph partitions the scene into spatially compact and semantically coherent regions, forming view-consistent 3D entities and providing a structured foundation for open-vocabulary understanding. Based on the graph structure, we design an efficient reprojection strategy that lifts 2D semantic features onto the superpoints, avoiding costly multi-view iterative training. The resulting representation ensures strong 3D semantic coherence and naturally supports hierarchical understanding, enabling both coarse- and fine-grained open-vocabulary perception within a unified semantic

field. Extensive experiments demonstrate that our method achieves state-of-the-art open-vocabulary segmentation performance, with semantic field reconstruction completed over 30× faster. Our code will be available at https://github.com/Atrovast/THGS.

## CCS Concepts

• **Computing methodologies** → **Scene understanding**.

## Keywords

Open-vocabulary, Scene Understanding, Gaussian Splatting, Superpoint

## 1 Introduction

In recent years, computer vision has made significant progress, particularly in developing systems that perceive and interact with the three-dimensional world. One emerging challenge in this context is 3D open-vocabulary scene understanding, where machines are expected to identify and localize arbitrary regions within a 3D environment based on free-form natural language input. This capability plays an important role in applications such as augmented reality and robotics, enabling users to refer to objects or regions in a scene

---

*Equal Contribution.

†Corresponding author.

using language descriptions rather than predefined labels. As 3D perception continues to scale in both complexity and scene size, the demand for accurate, efficient, and semantically structured 3D understanding is becoming increasingly critical.

Given the scarcity of large-scale 3D datasets with language annotations, many approaches [6, 9, 11, 20, 27, 33] leverage vision-language models (VLMs), such as CLIP [28] and LSeg [16], to distill open-vocabulary semantics into 3D scene representations. These methods extract semantic features from 2D images and lift them into 3D, enabling language-driven interaction with reconstructed scenes, without requiring dense 3D supervision. Among various 3D representations, 3D Gaussian Splatting (3DGS) [8] has recently emerged as a promising alternative for scene reconstruction due to its efficient optimization and real-time rendering. These advantages make it well-suited for embedding semantic information. Building on this, several recent works have explored combining image-derived language features with 3DGS to construct semantic fields under open-vocabulary settings [9, 19, 24, 32, 37, 42, 44].

Despite recent progress, existing methods still suffer from two key limitations. First, approaches like LangSplat [24] and LEGaussians [32] iteratively optimize semantic features on each view by rendering high-dimensional embeddings and backpropagation. However, they treat each Gaussian primitive in isolation and do not enforce consistency in 3D space, weakening semantic continuity across neighboring primitives. As a result, these methods often produce noisy semantic features and suffer from misalignment between 2D observations and the underlying 3D structure. Second, semantic field reconstruction is often time-consuming, even surpassing the cost of appearance modeling. Semantic information is generally higher-level and lower-frequency than visual details and is expected to exhibit spatial smoothness in 3D. Yet, existing methods optimize dense high-dimensional embeddings across all primitives without spatial regularization, resulting in an ill-posed process that converges slowly.

To overcome the inefficiencies and semantic inconsistencies of existing methods, we rethink how semantic information should be represented and transferred in 3D, as illustrated in Figure 1. Instead of optimizing dense per-view embeddings for individual primitives, we observe that high-level semantics are typically low-frequency, spatially coherent, and better captured through structured abstraction. This insight motivates our training-free framework, which builds hierarchical semantic fields by grouping Gaussian primitives into superpoints—compact clusters that are semantically homogeneous and geometrically coherent. We begin with SAM-guided segmentation to partition the scene into superpoints that align well with object boundaries across views, thereby enforcing 3D semantic consistency from the start. To further capture the multi-scale nature of real-world semantics, we progressively merge these superpoints into a multi-level graph, guided by segmentation masks of progressively coarser granularity. This hierarchy naturally supports open-vocabulary understanding at both object and part levels. Rather than relying on slow, view-by-view optimization to distill semantics, we directly reproject 2D semantic features onto the superpoint graph using a simple and efficient aggregation mechanism. This training-free, one-pass design eliminates costly iterative updates and enables semantic field reconstruction within minutes—achieving over 30×

speedup compared to existing approaches. Together, these innovations result in a scalable, hierarchical, and view-consistent solution for 3D open-vocabulary scene perception.

Our main contributions are summarized as follows:

- We propose a training-free framework for constructing 3D open-vocabulary semantic fields with strong spatial consistency and no iterative optimization.
- A contrastive Gaussian partitioning strategy is incorporated to improve boundary precision and ensure semantic coherence.
- We introduce a hierarchical merging and reprojection strategy that produces structured superpoint representations for multi-level semantic understanding.
- Extensive experiments demonstrate that our method achieves state-of-the-art open-vocabulary scene understanding, while reducing semantic field reconstruction time by over 30×.

## 2 Related Works

### 2.1 Gaussian Splatting

3D Gaussian Splatting (3DGS) has emerged as a powerful and efficient representation for 3D scenes [5, 8, 18, 25, 31]. By modeling scenes as collections of explicit Gaussian primitives and rasterizing them for image synthesis, 3DGS achieves real-time rendering and fast reconstruction. Its point-based structure also provides flexibility for applications such as editing and manipulation. Subsequent works have aimed to improve 3DGS along several dimensions, including geometric fidelity, scalability, and robustness. To enhance surface accuracy, SuGaR [4] introduces geometric constraints to Gaussian primitives. 2DGS [5] and PGSR [2] replace 3D ellipsoids with 2D disks for tighter surface alignment. Scaffold-GS [22] proposes a hierarchical anchor-based structure to reduce redundancy and improve scalability. Other efforts focus on robustness under sparse-view conditions [17, 43], anti-aliasing performance [38, 41], and adaptation to in-the-wild images [12, 30, 35, 36].

### 2.2 Open-vocabulary Scene Understanding with Radiance Fields

Vision-language models (VLMs) such as CLIP [28] have enabled open-vocabulary reasoning in visual domains, inspiring a series of methods for 3D semantic understanding. Early works [9, 11, 21, 33] lift CLIP features from 2D views into NeRF scenes, allowing text-based localization. However, the implicit volumetric nature of NeRF limits their training efficiency and real-time utility. Recent methods based on Gaussian Splatting offer faster reconstruction and explicit point-based representations. LangSplat [24] compresses CLIP features using an autoencoder and incorporates SAM masks for hierarchical supervision. LEGaussians [32] proposes a quantized and smoothed embedding strategy to preserve rendering quality. Gaussian Grouping [39] uses SAM [10] priors to improve boundary accuracy and region separation. GOI [26] introduces an optimizable semantic-space Hyperplane to achieve more accurate open-vocabulary semantic perception at test time. OpenGaussian [37] applies contrastive learning with SAM supervision to construct object-level semantic fields. SuperGSeg [19] shares a similar motivation to ours, drawing from superpoint-based methods. It performs
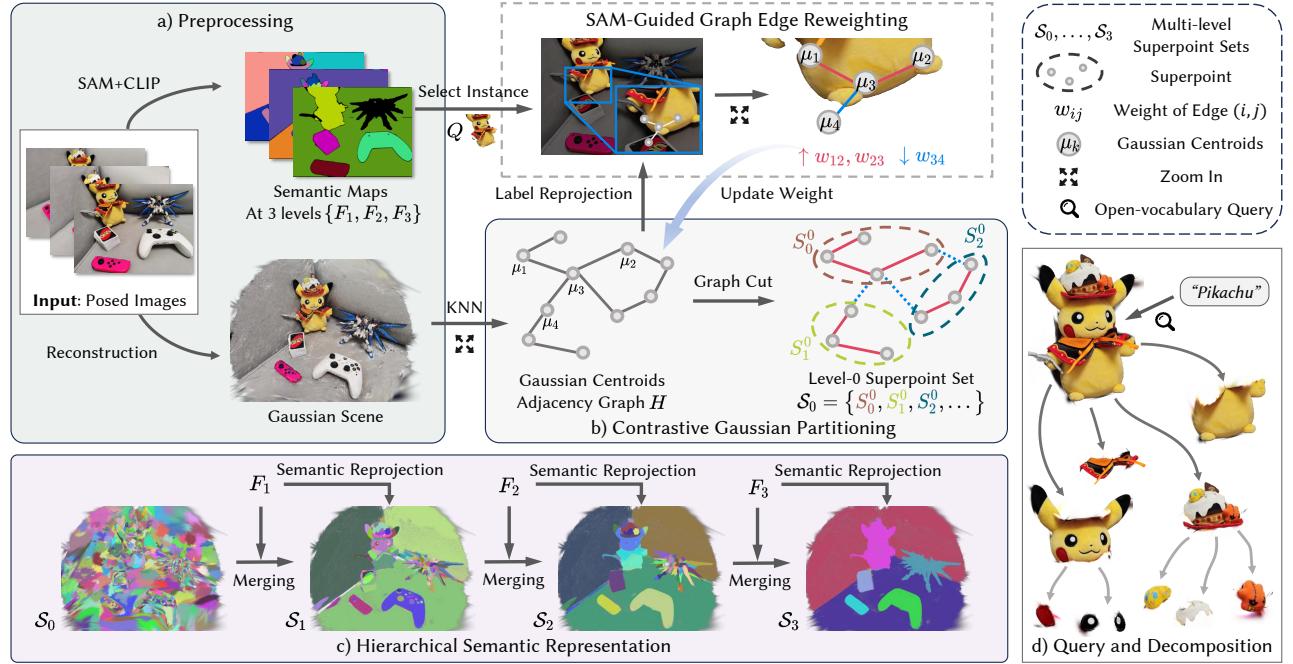
**Figure 2: Framework Overview. a) Preprocessing: Scene reconstruction and extraction of 2D semantic maps. b) Contrastive Gaussian Partitioning: A Gaussian adjacency graph is created, and its edge weights are adjusted using SAM-guided contrastive cues. The scene is then partitioned into superpoints. c) Hierarchical Semantic Representation: Superpoints are progressively merged to form a multi-level superpoint graph, while semantic features are reprojected onto each level. d) Query and Decomposition: The resulting hierarchical graph enables open-vocabulary query and part-based decomposition of scene objects.**

multi-view instance and hierarchical segmentation on Scaffold-GS by leveraging Super-Gaussians to enhance interpretability.

## 2.3 Superpoint-based Point Cloud Segmentation

Superpoint-based methods have gained increasing attention in 3D point cloud segmentation due to their efficiency and ability to preserve local geometric structure [7, 13, 15, 23, 29, 40]. Inspired by superpixels in 2D image segmentation [1], these approaches group neighboring points into compact and geometrically consistent regions called superpoints, which serve as mid-level representations for downstream semantic reasoning. SPG [15] and SPT [29] propose building superpoint graphs to model spatial and geometric relations between regions, and further refine superpoint features using neural networks. SAI3D [40] leverages SAM to iteratively merge superpoints into segments aligned with multi-view 2D masks, producing consistent instance-level groupings without requiring semantic supervision. Open3DIS [23] combines superpoints with a 3D instance segmentation network and embeds CLIP features at the instance level to support open-vocabulary segmentation. These approaches demonstrate the effectiveness of superpoint for scalable and structured 3D semantic understanding.

## 3 Methods

Given a set of posed images, we first reconstruct a 3D Gaussian scene $G$ using Gaussian Splatting [8]. Our goal is to extend $G$ with

open-vocabulary semantics, enabling users to query and segment objects in both 2D and 3D using natural language descriptions.

Unlike previous methods that rely on per-view optimization and lack mechanisms for global consistency, we construct a unified 3D semantic representation directly on a superpoint graph built from Gaussian primitives (GPs). This design enables fast, training-free reconstruction of a view-consistent semantic field, as illustrated in Figure 2. We begin by preprocessing each image with frozen vision-language models to extract 2D feature maps and perform scene reconstruction. We treat the Gaussian primitives in $G$ as a point cloud and perform a contrastive Gaussian partitioning to segment the scene into simple, semantically homogeneous clusters, referred to as *superpoints*. For the over-segmented superpoints, we adopt a hierarchical merging strategy that groups adjacent superpoints into semantically distinct larger clusters, resulting in a multi-level *superpoint graph* aligned with the levels of SAM masks. During this process, semantic features are reprojected onto the multi-level superpoints, enabling efficient construction of a semantic field.

## 3.1 Preliminaries: Gaussian Splatting

Gaussian Splatting represents a 3D scene as a collection of Gaussian primitives, which extend traditional point clouds with additional shape, opacity, and appearance parameters. Given known camera poses, the scene is rendered by splatting these primitives onto the image plane and rasterizing the resulting 2D ellipses in a differentiable manner.

Each Gaussian primitive is parameterized by a centroid $\mu \in \mathbb{R}^3$, a 3D covariance matrix $\Sigma$, opacity $\alpha$, and spherical harmonics coefficients $c$ for view-dependent color. To ensure that $\Sigma$ remains valid and interpretable, it is factorized as:

$$\Sigma = RSS^T R^T, \tag{1}$$

where $R$ is a rotation matrix and $S$ is a scaling matrix. The full set of learnable parameters for the $i$-th Gaussian primitive is given by: $\theta_i = \{\mu_i, c_i, \alpha_i, R_i, S_i\}$.

Rendering is performed via volumetric alpha compositing. For each pixel, the color $C$ is computed as:

$$C = \sum_{i \in G'} c_i \alpha_i T_i, \tag{2}$$

where $G'$ is the depth-sorted set of Gaussian primitives contributing to the pixel, and $T_i = \prod_{j=1}^{i-1}(1 - \alpha_j)$ denotes the accumulated transmittance up to the $i$-th primitive.

To achieve high-fidelity surface reconstruction, we adopt 2D Gaussian Splatting (2DGS) [5], which improves geometric accuracy by modeling each surface element as a 2D disk rather than a 3D ellipsoid. In addition to higher geometric precision, 2DGS provides surface normals $\mathbf{n}_i$ for each primitive, which can be utilized in the following processing.

## 3.2 Contrastive Gaussian Partitioning

We begin by partitioning the Gaussian primitives in the scene into a set of spatially compact and semantically coherent clusters, referred to as *superpoints*. To obtain these superpoints, we treat the GPs as a 3D point cloud and construct an adjacency graph based on spatial proximity. A graph partitioning algorithm is then applied to segment this graph into geometrically simple clusters.

However, unlike traditional point-based partitioning, GPs possess spatial extent and may span multiple semantic regions when projected onto the image plane. As a result, naive clustering can group semantically inconsistent content into a single superpoint. To mitigate semantic ambiguity, we use SAM-generated 2D masks as guidance to reweight the graph edges, encouraging connections within the same region and suppressing those across boundaries. This reweighting introduces a contrastive constraint into the partitioning process, encouraging the formation of semantically homogeneous superpoints aligned with object boundaries and consistent across views.

### 3.2.1 *Graph-based Gaussian Partitioning.* To adapt point cloud segmentation techniques to Gaussian primitives, we ignore the volumetric nature of each GP and represent it solely by its geometric centroid $\mu$. The set of centroids is denoted as $G_\mu$. We construct an undirected $K$-nearest neighbor graph by connecting each GP to its $K$ closest neighbors in Euclidean space. The resulting adjacency graph $H$ is defined as the combination of the vertex set $G_\mu$ and the edge set $E$:

$$H = (G_\mu, E), \quad E = \{(i, j) \mid \mu_i, \mu_j \in G_\mu\}. \tag{3}$$

For each node $\mu_i$, we construct a feature vector $p_i$ by concatenating its position and additional GP features, such as color or normal:

$$p_i = \text{concat}\{\mu_i, c_i, \mathbf{n}_i, \ldots\}, \quad \mu_i \in G_\mu. \tag{4}$$

And the edge weight $w_{ij}$ is computed inversely proportional to the normalized feature distance between $p_i$ and $p_j$:

$$w_{ij} = \frac{1}{1 + \text{dis}(\mu_i, \mu_j)}, \ \text{dis}(\mu_i, \mu_j) = \frac{\|p_i - p_j\|_2}{\frac{1}{|E|}\sum_{(k,l) \in E}\|p_k - p_l\|_2}. \tag{5}$$

This feature construction approach is adapted from [15, 29], with minor modifications tailored to Gaussian primitives.

With the graph weights defined, we apply the Cut Pursuit algorithm [14] to partition the graph into superpoints. Edges with lower weights are more likely to be cut, encouraging the separation of geometrically or semantically dissimilar regions. The output of the algorithm yields a set of constant connected components in the graph, which we define as level-0 superpoint set $\mathcal{S}_0$.

### 3.2.2 *SAM-guided Graph Edge Reweighting.* To encourage semantic consistency during superpoint partitioning, we refine the edge weights in the adjacency graph using SAM-predicted segmentation masks before applying Cut Pursuit. While the initial graph captures geometric proximity, it does not account for semantic alignment. Since Gaussian primitives have non-zero spatial extent, they may overlap multiple regions or objects, leading to semantic ambiguity.

To mitigate this, we introduce a mask-guided edge reweighting strategy. For each view, SAM masks are indexed with integer labels and reprojected onto the GPs (see Sec. 3.3). For any connected pair of nodes $\mu_i$ and $\mu_j$, if their mask labels differ ($l_i \neq l_j$), the connection is down-weighted; otherwise, it is strengthened.

We update the edge weights $w_{ij}$ in the adjacency graph accordingly:

$$w'_{ij} = w_{ij} + \delta_+ \cdot [l_i = l_j] - \delta_- \cdot [l_i \neq l_j], \tag{6}$$

where $[\cdot]$ is the Iverson bracket, $\delta_+$ and $\delta_-$ control the influence of SAM masks on the graph structure. In practice, we ensure that $w'_{ij} > 0$ to preserve meaningful graph connectivity. To account for reduced reliability in distant projections [40], $\delta_+$ and $\delta_-$ are scaled by a depth-aware decay function, diminishing the effect of masks when the corresponding GPs are far from the camera.

The updated graph is then used in the Cut Pursuit algorithm to generate superpoints. This reweighting step helps align superpoint boundaries more closely with semantic object boundaries, promoting intra-region consistency and inter-region separation.

## 3.3 Feature Reprojection for Semantic Guidance

Our framework requires transferring 2D segmentation information from SAM masks to 3D Gaussian primitives to support key components such as edge reweighting, hierarchical merging, and semantic assignment. We propose a training-free feature reprojection mechanism that robustly associates 2D mask labels with 3D Gaussian primitives in the presence of boundary ambiguity.

This process consists of two steps: (1) encoding each SAM mask label into a latent semantic embedding, and (2) reprojecting these embeddings onto Gaussian primitives through a rendering-guided aggregation scheme. This enables direct and reliable assignment of semantic features from 2D to 3D without iterative optimization.

### 3.3.1 *Latent label encoding.* For each SAM-generated mask $m_t$, where $t \in \{1, 2, ..., M\}$, we assign a unique integer label $t$ to identify the mask. However, directly reprojecting raw integer labels onto 3D Gaussian primitives often leads to semantic ambiguity, especially

near region boundaries. One-hot encodings are also suboptimal due to their variable length and high dimensionality.

To overcome these limitations, we map each label $t$ to a fixed-dimensional latent vector $Y_t \in \mathbb{R}^D$, sampled from a standard normal distribution and normalized to unit length:

$$Y_t = \text{normalize}(\text{rand}(D)), \quad t \in \{1, 2, ..., M\}. \tag{7}$$

Here, $M$ is the total number of SAM-predicted masks. Inspired by spline positional encoding [34], this approach distributes the label embeddings uniformly on the hypersphere, making them maximally discriminative under cosine similarity. In practice, we assign the latent label embedding $Y_t$ to every pixel belonging to mask $m_t$, resulting in a pixel-aligned semantic feature map $\widetilde{Y}$ in image space.

*3.3.2 Rendering-guided reprojection.* For each view, we rasterize the 3D scene using Gaussian Splatting and compute the transmittance $T_k^{\mathbf{x}}$ of each Gaussian primitive $k$ along the pixel ray $\mathbf{x}$, following the volumetric rendering formulation (Equation 2). Since transmittance reflects the contribution of each GP to the final pixel color, we use it as a soft weighting factor to aggregate 2D semantic information onto the GPs.

The latent feature $y_k$ of Gaussian primitive $k$ is computed by aggregating label embeddings $\widetilde{Y}^{\mathbf{x}}$ from all contributing pixels, weighted by their transmittance values $T_k^{\mathbf{x}}$:

$$y_k = \frac{\sum_{\mathbf{x}} T_k^{\mathbf{x}} \widetilde{Y}^{\mathbf{x}}}{\sum_{\mathbf{x}} T_k^{\mathbf{x}}}. \tag{8}$$

We then assign a discrete semantic label $l_k$ by matching latent feature $y_k$ to the closest entry $Y_t$ in the label embedding set $\{Y_t\}_{t=1}^M$ via cosine similarity:

$$l_k = \underset{t \in \{1,2,...,M\}}{\arg\max} \cos(y_k, Y_t). \tag{9}$$

This feature reprojection mechanism effectively preserves semantic boundaries and reduces ambiguity near object edges.

## 3.4 Hierarchical Semantic Representation

To build structured and interpretable scene semantics beyond the over-segmented superpoints $\mathcal{S}_0$ (from Sec. 3.2), we adopt a bottom-up hierarchical merging strategy inspired by SAI3D [40]. While SAI3D is initially developed for point-level instance grouping, we extend this approach to operate on Gaussian primitives and generalize it to construct a multi-level representation.

We use multi-level SAM masks $\{F_1, F_2, F_3\}$ to guide hierarchical merging. Each corresponds to a specific level of segmentation granularity. At each level, adjacent superpoints are aggregated based on their affinity score, producing progressively larger clusters. This process yields a multi-level superpoint graph $\mathcal{S}_1$, $\mathcal{S}_2$, and $\mathcal{S}_3$, which represent sub-parts, parts, and complete objects, respectively. Each level builds upon the previous one.

Semantic features are assigned to superpoints at all levels via the feature reprojection scheme, resulting in a hierarchical semantic field. This field enables open-vocabulary understanding at multiple levels and supports object decomposition.

*3.4.1 Superpoint Merging based on Affinity Score.* We adopt a hierarchical superpoint merging strategy inspired by SAI3D [40],

which merges superpoints based on their affinity scores. The affinity scores reflect the likelihood that two regions belong to the same object instance. We compute the score using 2D SAM-generated masks as guidance.

At level $q$, for each pair of superpoints $S_u^q$ and $S_v^q$, we reproject the 2D masks $m_t$ ($t \in \{1, 2, ..., M\}$) onto the scene using the rendering-guided method described in Sec. 3.3. For $S_u^q$ and $S_v^q$, we compute a histogram feature $\mathbf{h}_u^q$ and $\mathbf{h}_v^q$ that captures the distribution of its constituent Gaussian primitives across the reprojected masks.

The affinity score between $S_u^q$ and $S_v^q$ is then defined as the cosine similarity between their histogram features:

$$A_{u,v}^q = \cos\left(\mathbf{h}_u^q, \mathbf{h}_v^q\right). \tag{10}$$

The affinity scores are averaged across views. Adjacent superpoints with affinity scores exceeding a predefined threshold are merged to form coarser superpoints.

*3.4.2 Semantic Feature Assignment.* The resulting multi-level superpoints are spatially coherent and significantly sparser than the original Gaussian primitives (typically by two orders of magnitude), making semantic field construction more efficient. Moreover, since both superpoint partitioning and hierarchical merging are guided by SAM masks, the resulting superpoints exhibit strong alignment with object boundaries. This allows us to directly reproject semantic features to superpoints without any further optimization.

During the merging process, we establish the correspondence between each superpoint and the set of masks it overlaps with. We then assign a semantic feature to each superpoint by computing a weighted average of the features from the corresponding masks. Let $f_t$ denote the semantic feature of mask $m_t$. The semantic feature $f(S_k^q)$ of superpoint $S_k^q$ is computed as:

$$f(S_k^q) = \sum_{t=1}^{M} \omega_t \cdot f_t, \quad \omega_t = \frac{\text{NumVis}(S_k^q, m_t)}{\left|S_k^q\right|}. \tag{11}$$

Here, $\text{NumVis}(S_k^q, m_t)$ denotes the number of visible Gaussian primitives in $S_k^q$ that are assigned to mask $m_t$, and $\left|S_k^q\right|$ is the total number of GPs in the superpoint. The weight $\omega_t$ thus reflects the proportion of $S_k^q$ associated with mask $m_t$. We then aggregate multi-view features to obtain overall superpoint semantic features.

Semantic features derived from SAM masks of three levels are reprojected onto the corresponding superpoint levels $\mathcal{S}_1$, $\mathcal{S}_2$, and $\mathcal{S}_3$, representing sub-parts, parts, and whole objects respectively. These layers are nested within a hierarchical superpoint graph, forming a unified semantic field that supports multi-level, open-vocabulary scene understanding.

## 3.5 Evaluation

We follow the evaluation protocol proposed in LERF [9], adapted to our superpoint-based hierarchical representation. For each text query, we compute a relevance score between the query embedding $\phi_{\text{qry}}$ and each superpoint semantic feature $\phi_{\text{sp}}$ across one or multiple hierarchy levels.

For each superpoint, the relevance score is defined as:

$$\min_i \frac{\exp(\phi_{\text{sp}} \cdot \phi_{\text{qry}})}{\exp(\phi_{\text{sp}} \cdot \phi_{\text{qry}}) + \exp(\phi_{\text{sp}} \cdot \phi_{\text{canon}}^i)}, \tag{12}$$
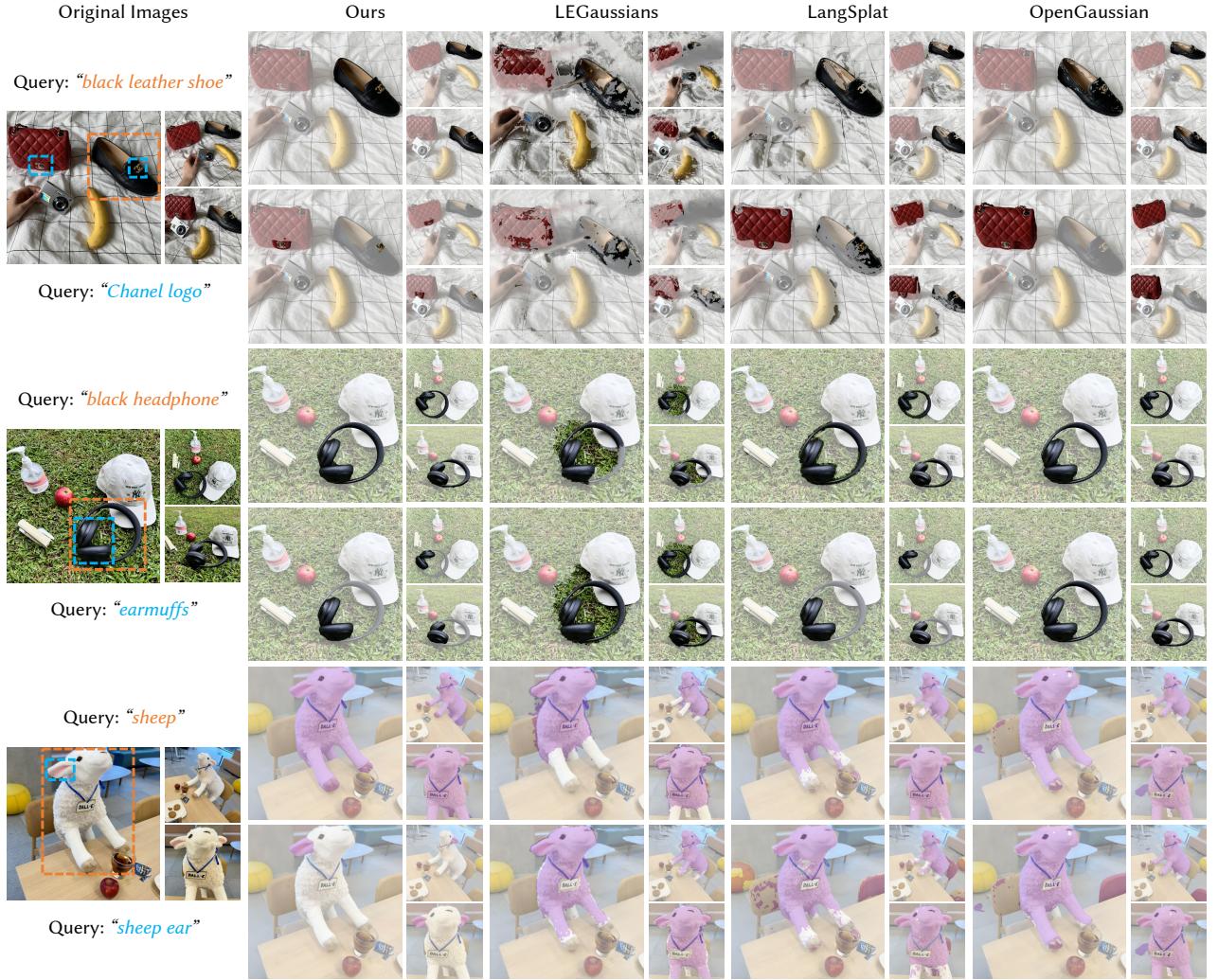
Figure 3: Qualitative comparisons of open-vocabulary segmentation on multi-view 2D images. We compare our method with LEGaussians [32], LangSplat [24], and OpenGaussian [37]. Each scene includes an object- and a part-level query. Query results are highlighted, and corresponding ground-truth regions are marked with orange (object) and blue (part) bounding boxes.

where $\phi_{\text{canon}}^i$ denotes the CLIP embedding of a predefined canonical concept (i.e., "object", "things", "stuff", "texture"). This contrastive scoring helps disambiguate the target object from generic or background categories.

After computing relevance scores, we select the top-ranked superpoints as query results, and retrieve their associated Gaussian primitives for further rendering or analysis. The query can be performed at a single hierarchical level or jointly across multiple levels for more robust localization.

To visualize the result in 2D, we do not render the full high-dimensional semantic field. Instead, we rasterize a binary presence mask by assigning a value of $b_k = 1$ to each selected GP and rendering it using standard volumetric accumulation:

$$B = \sum_{k \in G'} b_k \alpha_k T_k. \tag{13}$$

The resulting soft mask $B$ is thresholded at 0.5 to obtain a binary segmentation map.

## 4 Implementation Details

Our method is implemented on top of 2D Gaussian Splatting [5]. For semantic feature extraction, we follow LangSplat [24], using OpenCLIP ViT-B/16 as the vision-language encoder, alongside SAM ViT-H for 2D segmentation guidance. To support multi-level representation, we generate SAM masks at three different levels and obtain a CLIP embedding for each mask by feeding the masked region into the CLIP image encoder.

Our hierarchical semantic field construction is fully training-free and takes less than 2 minutes on average using a single RTX 3090 GPU. The runtime may vary depending on the number of input views and scene complexity.

Table 1: Quantitative comparison of 2D open-vocabulary segmentation on the LERF-OVS [9] and 3DOVS [21] datasets. "SR Time" refers to the Semantic Field Reconstruction Time. The best and second-best results are highlighted in red and orange.

| Methods | LERF-OVS mIoU (%) | | | | | LERF SR Time | 3DOVS mIoU (%) | | | | | | 3DOVS SR Time |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Figurines | Ramen | Teatime | Waldo | Overall | SR Time | Bed | Bench | Lawn | Room | Sofa | Overall | SR Time |
| LEGau. [9] | 27.60 | 13.75 | 45.21 | 23.71 | 27.57 | 65min | 11.21 | 62.66 | 58.61 | 64.03 | 9.02 | 41.11 | 75min |
| LangSplat [9] | 44.7 | 51.2 | 65.1 | 44.5 | 51.4 | 85min | 92.5 | 94.2 | 96.1 | 94.1 | 90.0 | 93.4 | 90min |
| GOI [26] | 36.89 | 35.09 | 66.56 | 45.23 | 45.94 | 15min | 97.70 | 73.24 | 96.26 | 73.91 | 85.16 | 85.26 | 14min |
| OpenGau. [37] | 69.73 | 21.11 | 63.35 | 34.91 | 47.28 | 45min | 57.87 | 74.94 | 63.95 | 40.16 | 70.56 | 61.49 | 55min |
| Ours | 57.30 | 43.46 | 68.33 | 50.65 | 54.94 | 90s | 96.01 | 95.14 | 96.88 | 92.71 | 93.89 | 94.93 | 25s |

## 5 Experiments

### 5.1 Experimental Setup

We evaluate our method on three datasets: LERF-OVS, 3DOVS, and ScanNet, covering both open-vocabulary 2D segmentation and 3D semantic understanding.

**Datasets.** The LERF-OVS benchmark is derived from the LERF dataset [9], which consists of 14 real-world indoor scenes with posed multi-view RGB images. LangSplat [24] selects four scenes (*figurines*, *ramen*, *teatime*, and *waldo kitchen*) from LERF and adds 2D mask annotations along with natural language descriptions for selected objects. 3DOVS dataset [21] includes 10 indoor and outdoor scenes featuring long-tail objects captured under diverse poses and backgrounds. Each scene provides 2D segmentation masks and text descriptions. We evaluate on five commonly used scenes: *bed*, *bench*, *lawn*, *sofa*, and *room*. ScanNet [3] is a large-scale RGB-D dataset with over 1,500 indoor scenes. It offers RGB-D sequences, point clouds from scans, and semantic annotations. Following OpenGaussian [37], we evaluate on 10 random scenes.

**Evaluation Protocol.** For all datasets, we reconstruct a 3D semantic field from posed multi-view images. Given a natural language query, the relevant region is retrieved via computing the relevance between text query and superpoint features, as described in Sec. 3.5. For LERF-OVS and 3DOVS, we reproject the retrieved 3D semantic regions onto 2D views and compare them with ground-truth masks. We report per-scene and overall mean Intersection-over-Union (mIoU), along with the average semantic field reconstruction time. For ScanNet, we evaluate directly in 3D by first assigning predicted instance-level labels to Gaussian primitives and then comparing them with the ground truth point-level semantic annotations. We report mean Intersection-over-Union (mIoU) and mean class accuracy (mAcc), which are computed over 19, 15, and 10-class subsets following the evaluation of OpenGaussian [37].

Table 2: Quantitative comparison for 3D semantic segmentation on ScanNet [3] dataset.

| Methods | 19 classes | | 15 classes | | 10 classes | |
| --- | --- | --- | --- | --- | --- | --- |
| | mIoU | mAcc | mIoU | mAcc | mIoU | mAcc |
| LangSplat [24] | 3.78 | 9.11 | 5.35 | 13.20 | 8.40 | 22.06 |
| LEGau. [32] | 3.84 | 10.87 | 9.01 | 22.22 | 12.82 | 28.62 |
| OpenGau. [37] | 24.73 | 41.54 | 30.13 | 48.25 | 38.29 | 55.19 |
| Ours | 34.39 | 50.74 | 39.61 | 57.07 | 46.38 | 64.74 |

### 5.2 Comparisons

We conduct a comparative evaluation of our approach in contrast with LEGaussians [32], LangSplat [24], GOI [26], and OpenGaussian [37].

**Qualitative Results.** We present the qualitative results generated by our method, along with comparisons to other approaches. Figures 3 and 4 provide a detailed overview of open-vocabulary scene understanding performance on the LERF-OVS and 3DOVS datasets, evaluated at both the 2D pixel level and the 3D Gaussian primitive level.

Figure 3 presents the multi-view segmentation results on 2D images. For each scene, we provide two natural language queries—at object and part levels—to demonstrate both cross-view consistency and hierarchical understanding. LEGaussians, trained only with 2D supervision and lacking geometric priors, often produces blurry and imprecise boundaries. Our method and LangSplat, both guided by SAM, show multi-level understanding and can distinguish object and its part (e.g., "headphones" vs. "earmuffs"). However, LangSplat struggles with finer categories like "sheep ear" and shows inconsistency in cases like the "Chanel logo", where our method remains accurate. OpenGaussian and our approach both retrieve in 3D and project to 2D, enabling consistent multi-view appearance. Yet, OpenGaussian's fixed-size codebook limits its part-level expressiveness, though it performs well at coarse object localization.

Figure 4 shows qualitative results rendered directly from querying on 3D Gaussian primitives. We compare our method against two representative baselines. LangSplat, which relies on per-view 2D optimization, produces spatially inconsistent and noisy 3D segmentations, with many scattered or misaligned points. Both our method and OpenGaussian yield relatively clean 3D segmentation results. However, OpenGaussian often lacks spatial coherence, mistakenly including scattered points from unrelated objects (e.g., "sheep" and "apple" in the bottom left scene). In contrast, our superpoint-based framework enforces structural regularity and spatial compactness, resulting in more coherent and reliable segmentations.

Overall, our method offers more coherent 3D segmentation, better view consistency, and supports multi-level perception in both 2D and 3D.

**Quantitative Results.** Tables 1 and 2 summarize the quantitative results across three benchmarks. Our method achieves state-of-the-art performance on both 2D and 3D open-vocabulary segmentation tasks. On LERF-OVS and 3DOVS (Table 1), we evaluate open-vocabulary on 2D images. Our method achieves strong per-scene performance and outperforms existing methods in overall

**Figure 4: Qualitative comparison of open-vocabulary 3D segmentation. We compare our method with OpenGaussian [37] and LangSplat [24] by visualizing the predicted 3D Gaussian primitives. The queried regions are annotated with colored bounding boxes on the original images to indicate object locations.**

mIoU. On ScanNet (Table 2), we directly evaluate in 3D against ground-truth semantic labels. Our method surpasses LangSplat and LEGaussian by a large margin and improves over OpenGaussian by approximately 9% in both mIoU and mAcc.

In addition to segmentation accuracy, Table 1 reports the semantic field reconstruction time, highlighting the efficiency of our framework. LEGaussians, LangSplat and OpenGaussian rely on iterative multi-view optimization, leading to slow convergence. In contrast, our method employs a training-free, forward-only pipeline based on a superpoint graph representation. This design yields up to 30× speedup on LERF-OVS and over 100× on 3DOVS, where fewer input views further amplify the advantage. While GOI improves efficiency through a feature clustering codebook, it remains within a training-based paradigm and exhibits limited segmentation quality. On the 3DOVS dataset, our method achieves over 30× faster reconstruction compared to GOI, demonstrating superior performance in both accuracy and speed.

## 5.3 Ablation Study

We present ablation results in Table 3, analyzing the impact of key components in our framework. Within the Contrastive Gaussian Partitioning module, removing edge reweighting results in less accurate superpoint boundaries, which hinders both hierarchical merging and semantic precision. Replacing the depth-aware decay with fixed coefficients $\delta_+$ and $\delta_-$ causes the model to over-reliance on reprojection results from regions far away from the camera, reducing boundary reliability.

For the hierarchical semantic representation, using only instance-level SAM masks limits the model's ability to capture fine-grained semantics. Without progressive merging, independently constructing all levels from $\mathcal{S}_0$ undermines hierarchical coherence and leads

**Table 3: Evaluation metrics for ablation studies on LERF-OVS [9] dataset. "SR Time" refers to the Semantic Field Reconstruction Time.**

| Method | mIoU | SR Time |
|---|---|---|
| w/o Edge Reweighting | 48.50 | 81s |
| w/o Depth Decay on $\delta_+$ and $\delta_-$ | 51.59 | 84s |
| Instance Level Only | 41.70 | 35s |
| w/o Progressive Merging | 44.26 | 120s |
| Ours Full | 54.94 | 90s |

to over- or under-segmentation. These results highlight the importance of both contrastive partitioning and hierarchical representation for achieving accurate 3D semantic understanding.

## 6 Conclusion

We present a novel framework for training-free, open-vocabulary 3D scene understanding by integrating Gaussian Splatting with a hierarchical superpoint graph. We begin with contrastive partitioning of Gaussian primitives and progressively merge superpoints under SAM-guided cues to form a multi-level graph structure. An efficient feature reprojection strategy is then used to construct a semantic field aligned with the multi-level superpoint graph. This training-free pipeline yields semantically consistent and hierarchically structured representations without iterative optimization. Extensive experiments on LERF-OVS, 3DOVS, and ScanNet demonstrate state-of-the-art performance in open-vocabulary segmentation while achieving significant speedup, highlighting the effectiveness of our approach for open-vocabulary 3D scene understanding.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. 2012. SLIC superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence* 34, 11 (2012), 2274–2282.

[2] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. 2024. PGSR: Planar-based Gaussian Splatting for Efficient and High-Fidelity Surface Reconstruction. *ArXiv preprint* abs/2406.06521 (2024). https://arxiv.org/abs/2406.06521

[3] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. 2017. ScanNet: Richly-Annotated 3D Reconstructions of Indoor Scenes. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2432–2443. doi:10.1109/CVPR.2017.261

[4] Antoine Guédon and Vincent Lepetit. 2024. SuGaR: Surface-Aligned Gaussian Splatting for Efficient 3D Mesh Reconstruction and High-Quality Mesh Rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 5354–5363. doi:10.1109/CVPR52733.2024.00512

[5] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2024. 2D Gaussian Splatting for Geometrically Accurate Radiance Fields. In *SIGGRAPH 2024 Conference Papers*. Association for Computing Machinery. doi:10.1145/3641519.3657428

[6] Chi Huang, Xinyang Li, Shengchuan Zhang, Liujuan Cao, and Rongrong Ji. 2024. NeRF-DetS: Enhancing Multi-View 3D Object Detection with Sampling-adaptive Network of Continuous NeRF-based Representation. *ArXiv preprint* abs/2404.13921 (2024). https://arxiv.org/abs/2404.13921

[7] Le Hui, Jia Yuan, Mingmei Cheng, Jin Xie, Xiaoya Zhang, and Jian Yang. 2021. Superpoint Network for Point Cloud Oversegmentation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 5490–5499. doi:10.1109/ICCV48922.2021.00546

[8] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* 42, 4 (2023), 1–14.

[9] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. 2023. LERF: Language Embedded Radiance Fields. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*. IEEE, 19672–19682. doi:10.1109/ICCV51070.2023.01807

[10] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. 2023. Segment Anything. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*. IEEE, 3992–4003. doi:10.1109/ICCV51070.2023.00371

[11] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. 2022. Decomposing NeRF for Editing via Feature Field Distillation. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). http://papers.nips.cc/paper_files/paper/2022/hash/93f250215e4889119807b6fac3a57aec-Abstract-Conference.html

[12] Jonas Kulhanek, Songyou Peng, Zuzana Kukelova, Marc Pollefeys, and Torsten Sattler. 2024. WildGaussians: 3D Gaussian Splatting In the Wild. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.). http://papers.nips.cc/paper_files/paper/2024/hash/25c0fe7b157821dd3140727dc07461da-Abstract-Conference.html

[13] Loïc Landrieu and Mohamed Boussaha. 2019. Point Cloud Oversegmentation With Graph-Structured Deep Metric Learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 7440–7449. doi:10.1109/CVPR.2019.00762

[14] Loic Landrieu and Guillaume Obozinski. 2017. Cut pursuit: Fast algorithms to learn piecewise constant functions on general weighted graphs. *SIAM Journal on Imaging Sciences* 10, 4 (2017), 1724–1766.

[15] Loïc Landrieu and Martin Simonovsky. 2018. Large-Scale Point Cloud Semantic Segmentation With Superpoint Graphs. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 4558–4567. doi:10.1109/CVPR.2018.00479

[16] Boyi Li, Kilian Q. Weinberger, Serge J. Belongie, Vladlen Koltun, and René Ranftl. 2022. Language-driven Semantic Segmentation. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net. https://openreview.net/forum?id=RriDjddCLN

[17] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. 2024. DNGaussian: Optimizing Sparse-View 3D Gaussian Radiance Fields with Global-Local Depth Normalization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 20775–20785. doi:10.1109/CVPR52733.2024.01963

[18] Xinyang Li, Zhangyu Lai, Linning Xu, Yansong Qu, Liujuan Cao, Shengchuan Zhang, Bo Dai, and Rongrong Ji. 2024. Director3D: Real-world Camera Trajectory and 3D Scene Generation from Text. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.). http://papers.nips.cc/paper_files/paper/2024/hash/89566c18d5b3e9836e8e16fde010b41d-Abstract-Conference.html

[19] Siyun Liang, Sen Wang, Kunyi Li, Michael Niemeyer, Stefano Gasperini, Nassir Navab, and Federico Tombari. 2024. SuperGSeg: Open-Vocabulary 3D Segmentation with Structured Super-Gaussians. *ArXiv preprint* abs/2412.10231 (2024). https://arxiv.org/abs/2412.10231

[20] Guibiao Liao, Kaichen Zhou, Zhenyu Bao, Kanglin Liu, and Qing Li. 2024. OV-NeRF: Open-vocabulary Neural Radiance Fields with Vision and Language Foundation Models for 3D Semantic Understanding. *ArXiv preprint* abs/2402.04648 (2024). https://arxiv.org/abs/2402.04648

[21] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El-Saddik, Christian Theobalt, Eric P. Xing, and Shijian Lu. 2023. Weakly Supervised 3D Open-vocabulary Segmentation. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). http://papers.nips.cc/paper_files/paper/2023/hash/a76b693f36916a5ed84d6e5b39a0dc03-Abstract-Conference.html

[22] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. 2024. Scaffold-GS: Structured 3D Gaussians for View-Adaptive Rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 20654–20664. doi:10.1109/CVPR52733.2024.01952

[23] Phuc D. A. Nguyen, Tuan Duc Ngo, Evangelos Kalogerakis, Chuang Gan, Anh Tuan Tran, Cuong Pham, and Khoi Nguyen. 2024. Open3DIS: Open-Vocabulary 3D Instance Segmentation with 2D Mask Guidance. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 4018–4028. doi:10.1109/CVPR52733.2024.00385

[24] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. 2024. LangSplat: 3D Language Gaussian Splatting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 20051–20060. doi:10.1109/CVPR52733.2024.01895

[25] Yansong Qu, Dian Chen, Xinyang Li, Xiaofan Li, Shengchuan Zhang, Liujuan Cao, and Rongrong Ji. 2025. Drag Your Gaussian: Effective Drag-Based Editing with Score Distillation for 3D Gaussian Splatting. *arXiv preprint arXiv:2501.18672* (2025).

[26] Yansong Qu, Shaohui Dai, Xinyang Li, Jianghang Lin, Liujuan Cao, Shengchuan Zhang, and Rongrong Ji. 2024. Goi: Find 3d gaussians of interest with an optimizable open-vocabulary semantic-space hyperplane. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 5328–5337.

[27] Yansong Qu, Yuze Wang, and Yue Qi. 2023. Sg-nerf: Semantic-guided point-based neural radiance fields. In *2023 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 570–575.

[28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. http://proceedings.mlr.press/v139/radford21a.html

[29] Damien Robert, Hugo Raguet, and Loïc Landrieu. 2023. Efficient 3D Semantic Segmentation with Superpoint Transformer. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*. IEEE, 17149–17158. doi:10.1109/ICCV51070.2023.01577

[30] Sara Sabour, Lily Goli, George Kopanas, Mark Matthews, Dmitry Lagun, Leonidas Guibas, Alec Jacobson, David J Fleet, and Andrea Tagliasacchi. 2024. Spotlessplats: Ignoring distractors in 3d gaussian splatting. *ArXiv preprint* abs/2406.20055 (2024). https://arxiv.org/abs/2406.20055

[31] You Shen, Zhipeng Zhang, Xinyang Li, Yansong Qu, Yu Lin, Shengchuan Zhang, and Liujuan Cao. 2025. Evolving High-Quality Rendering and Reconstruction in a Unified Framework with Contribution-Adaptive Regularization. *arXiv preprint arXiv:2503.00881* (2025).

[32] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. 2024. Language Embedded 3D Gaussians for Open-Vocabulary Scene Understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 5333–5343. doi:10.1109/CVPR52733.2024.00510

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

[34] Peng-Shuai Wang, Yang Liu, Yu-Qi Yang, and Xin Tong. 2021. Spline Positional Encoding for Learning 3D Implicit Signed Distance Fields. In *Proceedings of*

the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021, Zhi-Hua Zhou (Ed.). ijcai.org, 1091–1097. doi:10.24963/IJCAI.2021/151

[35] Yuze Wang, Junyi Wang, Ruicheng Gao, Yansong Qu, Wantong Duan, Shuo Yang, and Yue Qi. 2025. Look at the Sky: Sky-aware Efficient 3D Gaussian Splatting in the Wild. *IEEE Transactions on Visualization and Computer Graphics* (2025).

[36] Yuze Wang, Junyi Wang, and Yue Qi. 2024. WE-GS: An In-the-wild Efficient 3D Gaussian Representation for Unconstrained Photo Collections. *ArXiv preprint* abs/2406.02407 (2024). https://arxiv.org/abs/2406.02407

[37] Yanmin Wu, Jiarui Meng, Haijie Li, Chenming Wu, Yahao Shi, Xinhua Cheng, Chen Zhao, Haocheng Feng, Errui Ding, Jingdong Wang, and Jian Zhang. 2024. OpenGaussian: Towards Point-Level 3D Gaussian-based Open Vocabulary Understanding. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.). http://papers.nips.cc/paper_files/paper/2024/hash/21f7b745f73ce0d1f9bcea7f40b1388e-Abstract-Conference.html

[38] Zhiwen Yan, Weng Fei Low, Yu Chen, and Gim Hee Lee. 2024. Multi-Scale 3D Gaussian Splatting for Anti-Aliased Rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 20923–20931. doi:10.1109/CVPR52733.2024.01977

[39] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. 2023. Gaussian grouping: Segment and edit anything in 3d scenes. *ArXiv preprint* abs/2312.00732 (2023). https://arxiv.org/abs/2312.00732

[40] Yingda Yin, Yuzheng Liu, Yang Xiao, Daniel Cohen-Or, Jingwei Huang, and Baoquan Chen. 2024. SAI3D: Segment any Instance in 3D Scenes. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 3292–3302. doi:10.1109/CVPR52733.2024.00317

[41] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. 2024. Mip-Splatting: Alias-Free 3D Gaussian Splatting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 19447–19456. doi:10.1109/CVPR52733.2024.01839

[42] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya You, Zhangyang Wang, and Achuta Kadambi. 2024. Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 21676–21685. doi:10.1109/CVPR52733.2024.02048

[43] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. 2024. Fsgs: Real-time few-shot view synthesis using gaussian splatting. In *European conference on computer vision*. Springer, 145–163.

[44] Xingxing Zuo, Pouya Samangouei, Yunwen Zhou, Yan Di, and Mingyang Li. 2024. Fmgs: Foundation model embedded 3d gaussian splatting for holistic 3d scene understanding. *ArXiv preprint* abs/2401.01970 (2024). https://arxiv.org/abs/2401.01970

## A Evaluation Protocol on ScanNet

Following the experimental setup of OpenGaussian [37], we evaluate the performance of 3D semantic segmentation on the ScanNet [3] dataset. In the following, we detail the specific evaluation protocol on the ScanNet dataset used in OpenGaussian.

During the scene reconstruction stage, we use the raw scanned point clouds provided by the dataset as the initialization for Gaussian Splatting. The densification and position optimization are disabled, keeping all Gaussian primitives fixed to the input point cloud. Only other geometric and appearance attributes are optimized. For semantic field construction, each method utilizes the reconstructed Gaussian scenes in its own way.

During testing, we classify each GP based on its open-vocabulary features, and assess mIoU and mAcc by comparing the GP predictions with the ground-truth semantic labels of the ScanNet point cloud. This provides an evaluation of 3D point-level semantic understanding.

The predefined categories for classification are derived from the commonly occurring object classes in the ScanNet dataset. The 19 categories (as defined by ScanNet) used for text queries are: wall, floor, cabinet, bed, chair, sofa, table, door, window, bookshelf, picture, counter, desk, curtain, refrigerator, shower curtain, toilet, sink, and bathtub. The 15-category subset excludes picture, refrigerator, shower curtain, and bathtub, while the 10-category subset further excludes cabinet, counter, desk, curtain, and sink.

OpenGaussian evaluates on randomly selected 10 scenes from ScanNet: *scene0000_00, scene0062_00, scene0070_00, scene0097_00, scene0140_00, scene0200_00, scene0347_00, scene0400_00, scene0590_00,* and *scene0645_00.* To ensure a fair comparison, we conduct our evaluation on the same set of scenes.

## B Experimental Details

### B.1 Interactive Segmentation via Hierarchical Superpoint Graph

In addition to open-vocabulary querying, our method naturally enables interactive segmentation, owing to its structured and 3D-consistent representation. Unlike methods that rely on per-view 2D supervision, our approach explicitly models semantic entities in 3D space, allowing users to directly interact with the scene at the object or part level.

This capability is further enhanced by our multi-level superpoint graph, which organizes the scene into a semantic hierarchy—from fine-grained components to whole objects. As shown in Figure 5, a user can interactively select a superpoint corresponding to a queried concept (e.g., nose of the toy bear) and then navigate through the hierarchy to refine or expand the region (e.g., separating wider area of the mouse and the whole toy bear), enabling flexible and semantically coherent scene editing.

Together, these properties make our method well-suited for interactive applications that require accurate and hierarchical 3D scene understanding.

## B.2 Additional Comparisons

*B.2.1 Qualitative Results.* Figure 6 extends the main paper's comparison by presenting additional qualitative results for both 2D and
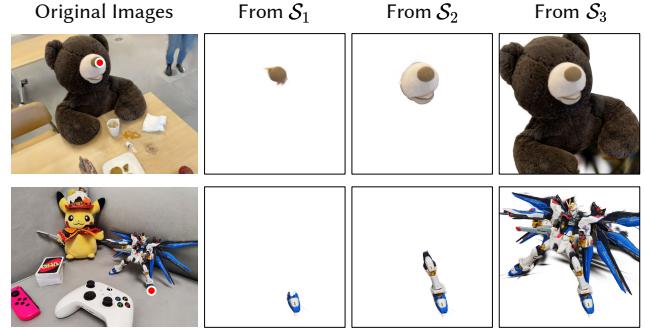


**Figure 5: Interactive Segmentation Results. Starting from a point prompt (red circle), our method enables coarse-to-fine object segmentation by traversing the hierarchical superpoint graph—retrieving regions from fine-grained components ($\mathcal{S}_1$) to complete objects ($\mathcal{S}_3$). This demonstrates the flexibility of our hierarchical representation for intuitive part-level interaction.**

3D open-vocabulary segmentation. We compare our method with LEGaussians, LangSplat, GOI, and OpenGaussian, across multiple scenes and query types, including both object-level and part-level descriptions.

Methods that rely on iterative 2D optimization without 3D constraints often produce inconsistent predictions between 2D observations and the underlying 3D scene structure. In contrast, our approach delivers cleaner segmentations with fewer artifacts and more accurate, consistent results. Moreover, it supports multi-level semantic understanding, enabling both coarse object-level and fine-grained part-level segmentation within a unified framework.

*B.2.2 Quantitative Results.* We evaluate the performance of open-vocabulary segmentation using multiple metrics. In addition to the mIoU results reported in the main paper, Table 4 presents the mean class accuracy (mAcc) across all evaluated scenes. This complementary metric further demonstrates the effectiveness and robustness of our framework. Our method consistently achieves strong performance across individual scenes and outperforms prior approaches in overall mAcc on both LERF-OVS [9] and 3DOVS [21] datasets, highlighting its capability for accurate and efficient open-vocabulary 3D understanding.

### B.3 Ablation Study: Effect of Gaussian Representation

To assess the impact of the underlying Gaussian representation, we compare our pipeline built on 2DGS [5] with a variant that uses the original 3DGS [8] for scene reconstruction.

While both 2DGS and 3DGS represent scenes using Gaussian primitives, 2DGS adopts surface-aligned disks that better capture object geometry. This alignment preserves object boundaries and reduces ambiguity between nearby objects. As a result, the constructed adjacency graph more accurately captures object-level structure, making it less likely for primitives from different objects to be erroneously grouped into the same superpoint.

**Table 4: Additional quantitative comparison of 2D open-vocabulary segmentation on the LERF-OVS [9] and 3DOVS [21] datasets. "SR Time" refers to the Semantic Field Reconstruction Time. The best and second-best results are highlighted in red and orange.**

| Methods | LERF-OVS mAcc (%) | | | | | LERF SR Time | 3DOVS mAcc (%) | | | | | | 3DOVS SR Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Figurines | Ramen | Teatime | Waldo | **Overall** | | Bed | Bench | Lawn | Room | Sofa | **Overall** | |
| LEGau. [32] | 96.65 | 76.86 | 92.32 | 79.65 | 86.37 | 65min | 69.72 | 95.60 | 96.02 | 96.58 | 71.59 | 85.91 | 75min |
| LangSplat [9] | 98.81 | 94.16 | 97.32 | 95.90 | 96.55 | 85min | 99.17 | 99.05 | 99.75 | 99.82 | 98.99 | 99.35 | 90min |
| GOI [26] | 96.16 | 92.11 | 97.68 | 91.95 | 94.48 | 15min | 99.78 | 88.42 | 99.80 | 99.17 | 94.02 | 96.24 | 14min |
| OpenGau. [37] | 99.66 | 95.11 | 98.28 | 93.12 | 96.54 | 45min | 75.21 | 88.15 | 73.30 | 72.01 | 94.98 | 80.73 | 55min |
| Ours | 99.30 | 96.39 | 99.13 | 97.24 | 98.02 | 90s | 99.53 | 99.58 | 99.72 | 99.59 | 99.36 | 99.56 | 25s |



**Figure 6: Qualitative comparisons of open-vocabulary on 2D images (top four rows) and 3D Gaussian primitives (bottom three rows). We show results from our method alongside LEGaussians [32], LangSplat [24], GOI [26], and OpenGaussian [37]. Our approach delivers coherent 3D understanding and effectively supports open-vocabulary querying at both object and part levels. The queried foreground regions are highlighted, and the prompts are shown on the left side of each row.**

In contrast, 3DGS relies on volumetric ellipsoids and often introduces a larger number of floaters, particularly in textureless or blurred regions. These floaters introduce noisy or misleading proximity relations in the graph. This increases the risk of grouping unrelated Gaussians and consequently damages the spatial and semantic coherence of superpoints.

**Table 5: Evaluation of different Gaussian representations on the LERF-OVS [9] dataset.**

| Method | mIoU | mAcc | SR Time |
|---|---|---|---|
| Ours (w/ 2DGS) | 54.94 | 98.02 | 90s |
| Ours (w/ 3DGS) | 42.50 | 96.88 | 102s |

Table 5 shows that switching to 3DGS results in a drop in both mIoU and mAcc due to degraded boundary precision and object separation. While the runtime remains comparable, the segmentation quality is noticeably reduced. These findings confirm the advantage of 2DGS as a geometry-aware representation that enables more accurate and spatially consistent semantic field construction.