# GOI: Find 3D Gaussians of Interest with an Optimizable Open-vocabulary Semantic-space Hyperplane

Yansong Qu[*], Shaohui Dai[*], Xinyang Li, Jianghang Lin,

Liujuan Cao[†], Shengchuan Zhang, Rongrong Ji

Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Ministry of Education of China,

Xiamen University, Fujian, China

{quyans,daish}@stu.xmu.edu.cn,imlixinyang@gmail.com,hunterjlin007@stu.xmu.edu.cn
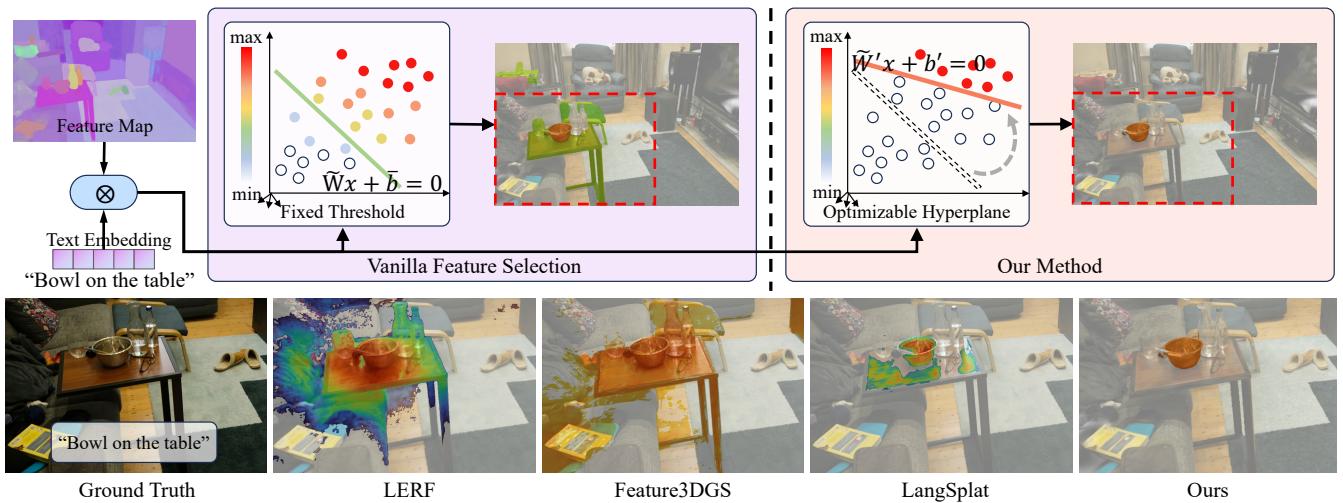
{caoliujuan,zsc_2016,rrj}@xmu.edu.cn

**Figure 1:** We propose GOI, an innovative approach to 3D open-vocabulary scene understanding based on 3D Gaussian Splatting [20]. In the top row, we emphasize our key contribution: the Optimizable Semantic-space Hyperplane (OSH). Instead of relying on a manually set, fixed empirical threshold for relative feature selection, which frequently lacks universal accuracy, OSH is fine-tuned for each query to accurately locate target regions in response to natural language prompts. The bottom row showcases our superior performance in open-vocabulary querying compared to other approaches.

## ABSTRACT

3D open-vocabulary scene understanding, crucial for advancing augmented reality and robotic applications, involves interpreting and locating specific regions within a 3D space as directed by natural language instructions. To this end, we introduce GOI, a framework that integrates semantic features from 2D vision-language foundation models into 3D Gaussian Splatting (3DGS) and identifies 3D Gaussians of Interest using an Optimizable Semantic-space Hyperplane. Our approach includes an efficient compression method that utilizes scene priors to condense noisy high-dimensional semantic features into compact low-dimensional vectors, which are subsequently embedded in 3DGS. During the open-vocabulary querying process, we adopt a distinct approach compared to existing methods, which depend on a manually set fixed empirical threshold to select regions based on their semantic feature distance to the query text embedding. This traditional approach often lacks universal accuracy, leading to challenges in precisely identifying specific target areas. Instead, our method treats the feature selection process as a hyperplane division within the feature space, retaining only those features that are highly relevant to the query. We

leverage off-the-shelf 2D Referring Expression Segmentation (RES) models to fine-tune the semantic-space hyperplane, enabling a more precise distinction between target regions and others. This fine-tuning substantially improves the accuracy of open-vocabulary queries, ensuring the precise localization of pertinent 3D Gaussians. Extensive experiments demonstrate GOI's superiority over previous state-of-the-art methods. Our project page is available at https://quyans.github.io/GOI-Hyperplane/ .

## CCS CONCEPTS

• **Computing methodologies → Scene understanding**.

## KEYWORDS

Open-vocabulary, 3D scene understanding, 3D Gaussian Splatting, Semantic Field, Hyperplane

---

[*] Equal Contribution.
[†] Corresponding Author.

# 1 INTRODUCTION

The field of computer vision has witnessed a remarkable evolution in recent years, driven by advancements in artificial intelligence and deep learning. A critical aspect of this progress is the enhanced ability of computer systems to interpret and interact with the three-dimensional world. The growing complexity in technology use has spurred a significant demand for advanced 3D visual understanding. This evolution brings to the fore the significance of the open-vocabulary querying task [5, 28, 34] — the capacity to process and respond to user queries formulated in natural language, enabling a more natural and flexible interaction between users and the digital world. Such advancements hold the potential to enhance how human navigate and manipulate complex three-dimensional data [14, 19, 47], bridging the gap between human cognitive abilities and computerized processing [7, 21].

Due to the scarcity of large-scale and diverse 3D scene datasets with language annotations, earlier Methods [21, 29] distill the open-vocabulary multimodal knowledge from off-the-shelf vision-language models, such as CLIP [42] and LSeg [24], into Neural Radiance Fields (NeRF) [35]. However, because of the implicit representation inherent in NeRF, these methods encounter impediments in terms of speed and accuracy, considerably limiting their practical application. Recently, the 3D Gaussian Splatting (3DGS) [20] has emerged as an effective representation of 3D scenes, and there have been explorations in constructing semantic fields [40, 49, 60]. This lifting approach requires pixel-aligned semantic features, whereas CLIP encodes the entire image into one global semantic feature. [21, 30, 49] utilize a multi-scale feature pyramid that incorporates CLIP embeddings from image crops. This approach, however, leads to blurred semantic boundaries, a problem that persists despite the introduction of DINO [3] constraints, resulting in unsatisfactory query results.

In this work, we introduce 3D **G**aussians **O**f **I**nterest (GOI). We utilize the vision-language foundation model APE [48] to extract pixel-aligned semantic features from multi-view images. GOI leverages these semantic features to reconstruct a 3D Gaussian semantic field. Given the explicit representational nature of 3DGS, directly embedding high-dimensional semantic features into each 3D Gaussian results in high computational demands. To mitigate this, we introduce the Trainable Feature Clustering Codebook (TFCC), which compresses noisy high-dimensional features based on scene priors, significantly reducing storage and rendering costs while maintaining each feature's informational capacity. Moreover, current open-vocabulary query strategies call for setting a fixed empirical threshold to ascertain features proximate to the query text. This, however, results in a failure to precisely query the targets. We introduce the Optimizable Semantic-space Hyperplane (OSH) to address this issue. OSH is fine-tuned by the Referring Expression Segmentation (RES) model, which aims to identify binary segmentation masks in 2D RGB images for text queries and is recognized for its robust spatial and localization capabilities. The OSH enhances GOI's spatial perception for more precise phrasal queries like "the table under the bowl", aligning query results more closely with target regions. Additionally, we have meticulously expanded and annotated a subset of the Mip-NeRF360 [1] dataset, tailored for the open-vocabulary query task. Owing to our method's proficient 3D

open-vocabulary scene understanding, it is practical for a range of downstream applications, notably scene manipulation and editing.

In summary, the main contributions of our work include:

- We propose GOI, an innovative framework based on 3D Gaussian Splatting for accurate 3D open-vocabulary semantic perception. The Trainable Feature Clustering Codebook (TFCC) is further introduced to efficiently condense noisy high-dimensional semantic features into compact, low-dimensional vectors, ensuring well-defined segmentation boundaries.

- We introduce the Optimizable Semantic-space Hyperplane (OSH), which eschews the fixed empirical threshold for relative feature selection due to its limited generalizability. Instead, OSH is fine-tuned for each text query with the off-the-shelf RES model to precisely locate target regions.

- Extensive experiments demonstrate that our method outperforms the state-of-the-art methods, achieving substantial improvements in mean Intersection over Union (mIoU) of 30% on the Mip-NeRF360 dataset [1] and 12% on the Replica dataset [50].

# 2 RELATED WORK

## 2.1 Neural Scene Representation

Recent methods in representing 3D scenes with neural networks have made substantial progress. Notably, Neural Radiance Fields (NeRF) [35] have excelled in novel view synthesis, producing highly realistic new viewpoints. However, NeRF's reliance on a neural network for complete implicit representation of scenes leads to tedious training and rendering times. Many subsequent methods [6, 12, 18, 36, 43, 44] have concentrated on improving its performance. In order to enhance the quality of surface reconstruction, [10, 13, 33, 52, 53] uses the signed distance function (SDF) for surface expression and uses a novel volume rendering scheme to learn an SDF representation. On the other hand, some approaches [8, 9, 39, 41, 55, 56] have explored the combination of implicit and explicit representations, utilizing traditional geometric structures, such as point clouds or mesh, to enhance NeRF's performance and to enable more downstream tasks. Kerbl et al. proposed 3D Gaussian Splatting (3DGS) [20], which greatly accelerates the rendering speed of novel view synthesis and achieves high-quality scene reconstruction. Unlike NeRF that represents a 3D scene implicitly with neural networks, 3DGS represent a scene as a set of 3D Gaussian ellipsoids, and accomplish efficient rendering by rasterizing the Gaussian ellipsoids into images. The technique adopted by 3DGS, which entails encoding scene information into a collection of Gaussian ellipsoids, provides distinct advantages [25, 26, 54]. It permits easy manipulation of specific parts in the reconstructed scene without significantly affecting other components. We have extended the 3DGS to achieve open-vocabulary 3D scene perception.

## 2.2 2D Visual Foundation Models

Foundation Models (FM) are becoming an impactful paradigm in the content of AI. They are typically pre-trained on vast amounts of data, possess numerous model parameters, and can be adapted to a wide range of downstream tasks [2]. The efficacy of 2D visual foundation models is evident in multiple visual tasks, such

as object localization [31] and image segmentation [15–17]. The incorporation of multimodal capabilities substantially amplifies the perceptual ability of these models. For instance, CLIP [42], by using contrastive learning, aligns the features of text encoders and image encoders into the unified feature space. Similarly, SAM [22] showcases immersive capabilities as a promptable segmentation model, delivering competitive, even superior zero-shot performance vis-à-vis earlier fully-supervised models. DINO [4, 37], a self-supervised Vision Transformer (ViT) model, is trained on vast unlabeled images. The model deciphers a semantic representation of images, encompassing components such as object boundaries and scene layouts.

Moreover, recent efforts are focused on leveraging existing pre-trained models, thereby pushing the limit of Foundation Models. Grounding DINO [32] represents an open-set object detector executing target detection based on textual descriptions. It utilizes CLIP and DINO as basic encoders, and proposes a tight fusion approach for better synthesizing of visual-language information. Grounded SAM [45] integrates Grounding DINO with SAM, facilitating the detection and segmentation for arbitrary queries. APE [48] is a universal visual perception model designed for diverse tasks like segmentation and grounding. Rigorously designed visual-language fusion and alignment modules enable APE to detect anything in an image swiftly without heavy cross-modal interactions.

## 2.3 3D Scene Understanding

Earlier works, such as Semantic NeRF [59] and Panoptic NeRF [11], introduced the transfer of 2D semantic or panoptic labels into 3D radiance fields for zero-shot scene comprehension. Following this, [23, 51] capitalized on pixel-aligned image semantic features, which they lifted to 3D, rather than relying on pre-defined semantic labels. Vision-language models like CLIP exhibited impressive performance in zero-shot image understanding tasks. A subsequent body of work [21, 23, 30] proposed leveraging CLIP and CLIP-based visual encoders to extract dense semantic features from images, with the aim of integrating them into NeRF scenes.

The recently proposed 3D Gaussian Splatting has achieved leading benchmarks in areas of novel view synthesis and reconstruction speed. This advancement has made the integration of 3D scenes with feature fields more efficient. LangSplat [40], LEGaussians [49], Feature 3DGS [60], Gaussian Grouping [57] explored the integration of pixel-aligned feature vectors from 2D models like LSeg, CLIP, DINO and SAM into 3D Gaussian frameworks so as to enabling 3D open-vocabulary query and localization of scene areas.

## 3 METHODS

### 3.1 Problem Definition and Method Overview

Given a set of posed images $I = \{I_1, I_2, \ldots, I_K\}$, a 3D Gaussian scene $S$ can be reconstructed using the standard 3D Gaussian Splatting technique [20] based on $I$. Our method expands $S$ with open-vocabulary semantics, enabling us to precisely locate the Gaussians of interest based on a natural language query.

We begin by recapping the vanilla 3D Gaussian Splatting (Sec. 3.2). Figure 2 illustrates the overview pipeline of our method. Initially, we utilize an frozen image encoder, well-aligned with the language

space, to process each image $I_k$ and derive the 2D semantic feature maps $V = \{V_1, V_2, \ldots, V_K\}$ (Sec. 3.3). To integrate these 2D high-dimensional feature maps into 3DGS, while ensuring minimal storage and optimal computational performance, Trainable Feature Clustering Codebook (TFCC) is proposed (Sec. 3.4). We expand 3DGS to reconstruct 3D Gaussian Semantic Field (Sec. 3.5). Following this, we explain how to utilize the RES model to optimize the Semantic-space Hyperplane, thereby achieving accurate open-ended language queries in 3D Gaussians (Sec. 3.6).

### 3.2 Vanilla 3D Gaussian Splatting

3D Gaussian Splatting utilizes a set of 3D Gaussians, essentially Gaussian ellipsoids, which bears a significant resemblance to point clouds, to model the scene and accomplish fast rendering by efficiently rasterizing Gaussians into images, given cameras poses. Specifically, each 3D Gaussian is parameterized by its centroid $x \in \mathbb{R}^3$, a 3D anisotropic covariance matrix $\Sigma$ in world coordinates, an opacity value $\alpha$, and spherical harmonics (SH) $c$. In the rendering process, 3D Gaussians are projected on to the 2D image plane, which transforms 3D Gaussian ellipsoids into 2D ellipses. $\Sigma$ is transformed to $\Sigma'$ in camera coordinates:

$$\Sigma' = JW\Sigma W^T J^T, \tag{1}$$

where $W$ denotes the world-to-camera tranformation matrix and $J$ is the Jacobian matrix for the projective transformation. In practical, $\Sigma$ is decomposed into a rotation matrix $R$ and a scaling matrix $S$:

$$\Sigma = RSS^T R^T. \tag{2}$$

This decomposition is to ensure that $\Sigma$ is physically meaningful during the optimization. To summarize, the learnable parameters of the $i$-th 3D Gaussian are represented by $\theta_i = \{x_i, R_i, S_i, \alpha_i, c_i\}$.

A volumetric rendering process, similar to NeRF, is then employed in the rasterization to compute the color $C$ of each pixel.

$$C = \sum_{i \in G} c_i \alpha_i T_i, \tag{3}$$

where $G$ denotes a set of 3D Gaussians sorted by their depth, and $T_i$ represents the transmittance, defined as the cumulative product of the opacity values of Gaussians that superimpose on the same pixel, computed through $T_i = \prod_{j=1}^{i-1}(1 - \alpha_j)$.

### 3.3 Pixel-level Semantic Feature Extraction

Prior research has broadly employed CLIP for feature lifting in the 3D radiance field, owing to its superior capability in managing open-vocabulary queries. [23, 60] use LSeg [24] to extract pixel-aligned CLIP features. However, LSeg proves inadequate in recognizing long-tail objects. To compensate for CLIP's limitation for yielding only image-level features, methodologies such as [21, 40, 49] adopt a feature pyramid approach, using cropped image encoding to represent local features. These methods extract pixel-level features from the CLIP model, but the generated feature maps lack geometric boundaries and correspondence to the scene objects. As such, pixel-aligned DINO features are introduced and predicted simultaneously with the CLIP features, thus bounding CLIP with the object geometry. Leveraging the success of SAM, [27, 40] utilizes SAM explicitly to constrain the object-level boundaries of the features.
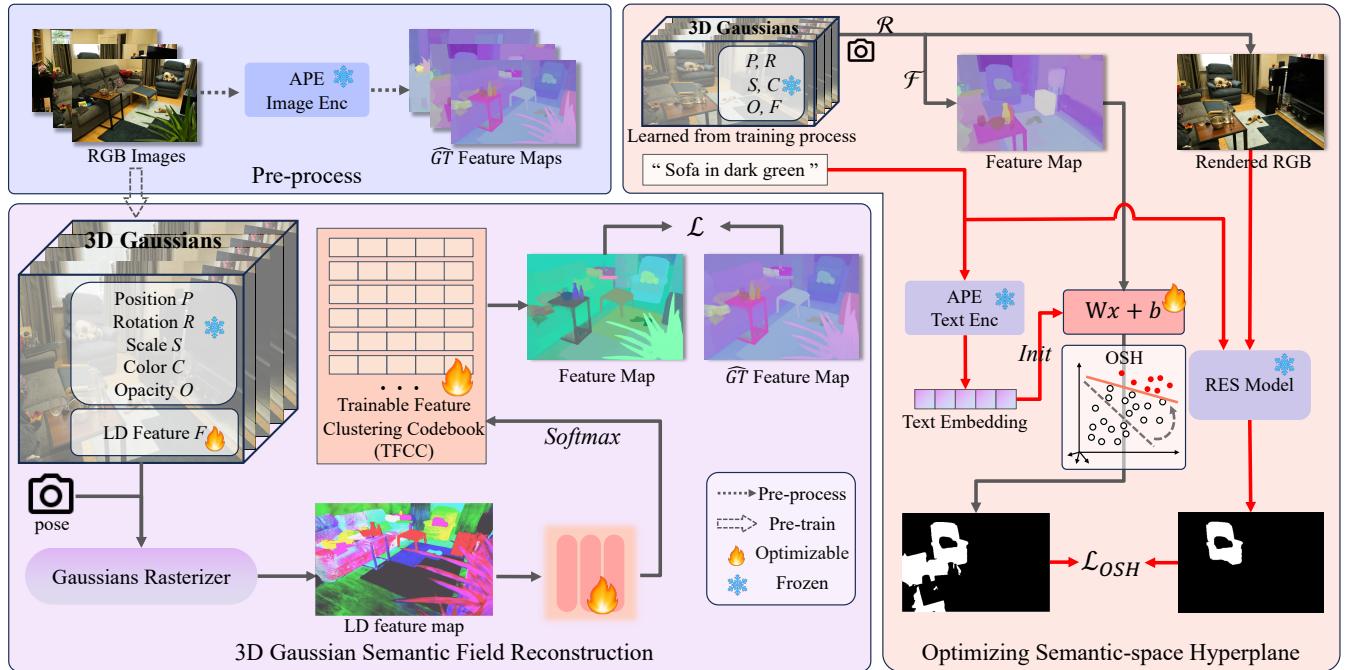
**Figure 2: The framework of our GOI. Top left: Reconstruction of a 3D Gaussian scene [20], encoding multi-view images. Bottom left: The optimization process. For each training view, a low-dimensional (LD) feature map is rendered through Gaussian Rasterizer and transformed into a predicted feature map via the Trainable Feature Clustering Codebook (TFCC). Right: The pipeline illustrates open-vocabulary querying. The processes denoted by $\mathcal{R}$ and $\mathcal{F}$ correspond to rendering and feature map prediction, respectively. The red line indicates operations exclusive to the initial query with a new text prompt. During these operations, the Optimizable Semantic-space Hyperplane (OSH) is fine-tuned to more precisely delineate the target region.**

However, using multiple models for feature extraction substantially increases the complexity for training and image prepocessing.

We leverage the Aligning and Prompting Everything All at Once model (APE) [48], which has the ability to efficiently align the features of vision and language. In APE, a fixed language model formulates language features, and a visual encoder is trained from scratch. The core of the visual encoder, derived from the DeformableDETR [61], provides APE with formidable detection and localization capacities. Additionally, APE possesses specially designed modules for vision-language fusion and vision-language alignment. The modules diminish cross-modal interaction and subsequently reducing computational costs. Therefore, APE presents a robust solution for feature lifting. For this purpose, we make minor modifications to the APE model to extract pixel-aligned features with fine boundaries efficiently (~2s per image). We treat the encoded pixel-aligned feature maps as the pseudo ground truth features, denoted as $\widehat{GT}$.

We extract APE feature maps from all training viewpoints and embed them into each 3D Gaussian to reconstruct a 3D semantic field. During the open-vocabulary querying process, we use the language model from pretrained APE to encode the language prompts.

### 3.4 Trainable Feature Clustering Codebook

Due to APE being trained on mass data and the need to align text and image features, it results in a higher feature dimensionality

(256). As the previous works [40, 49] have mentioned, directly lifting high-dimensional semantic features into each 3D Gaussian results in excessive storage and computational demands. The semantics of a single scene cover only a small portion of the original CLIP feature space. Therefore, leveraging scene priors for compression can effectively reduce storage and computational costs. On the other hand, due to the inherent multi-view inconsistency of 2D semantic feature map encoded by visual encoders, Gaussians tend to overfit each training viewpoint, inheriting this inconsistency and causing discrepancies between 3D and 2D within an object. Therefore, we introduce the Trainable Feature Clustering Codebook (TFCC), which leverages scene priors to compress the semantic space of a scene and encode it into a $N$ length codebook. Features similar in the feature space are explicitly constrained to the same entry in the table. Each entry in the codebook has a feature dimension equivalent to the dimension of the semantic features. This approach effectively reduces redundant and noisy semantic features while preserving sufficient scene information and clear semantic boundaries.

### 3.5 3D Gaussian Semantic Fields

We introduce a low-dimensional semantic feature, symbolized as $f$, into each 3D Gaussian, capitalizing on the redundancy of high-dimensional semantics across the scene and dimensions to facilitate efficient rendering. To create a 2D semantic representation, we employ a volumetric rendering process similar to color rendering

(Sec. 3.2) onto the low-dimensional semantic feature.

$$\hat{f} = \sum_{i \in G} f_i \alpha_i T_i. \tag{4}$$

$\hat{f}$ is the pixel-wise low-dimensional feature. We utilize an MLP as a feature decoder to obtain logits $e$, which are subsequently activated by the Softmax function to find the corresponding TFCC entry's index. This process acquires the feature $v$ in the high-dimensional semantic space for each $\hat{f}$. Given that volumetric rendering is essentially a process of weighted averages, the 3D Gaussian feature $f$ and the rendered 2D pixel-wise feature $\hat{f}$ are fundamentally equivalent. The low-dimensional feature $\hat{f}$ and $f$ can both be recovered to semantic feature $v$ through the MLP decoder $\mathcal{D}$ and the TFCC $\mathcal{T}$ with $N$ entries,

$$v = \mathcal{T}\left[\arg\max_i(e_i)\right], \tag{5}$$

where $e = \mathcal{D}(\hat{f})$ and $e \in \mathbb{R}^N$. Thus, both 2D and 3D features can be restrained to a compact and finite semantic space.

Initially in the semantic field optimization, we focus on learning the TFCC from $\widehat{GT}$ features. To enhance reconstruction efficiency, we adopt $k$-means clustering through $\widehat{GT}$ feature maps $V$ for the codebook initialization. Also, we find some resemblance between the learning of TFCC and the contrastive pre-training from CLIP: Features in the codebook are to align with the $\widehat{GT}$ features, and each $\widehat{GT}$ feature, denoted as $v_{gt}$, is assigned to one TFCC entry with the highest similarity. However, the assignment of a pixel feature to a particular entry is not predetermined, rather it pivots on similarity. Therefore, we devise a self-supervised loss function aimed at reducing the self-entropy of the clustering process.

$$\mathcal{L}_{ent} = -\sum_{i=1}^{N} p_i \log(p_i), \tag{6}$$

where $p_i = \text{Softmax}\left(\cos\langle v_{gt}, \mathcal{T}[i] \rangle \cdot \tau\right)$ and $\tau$ is the annealing temperature. To accelerate the process, we additionally optimize the entry with the highest similarity, introducing a loss function similar to [49],

$$d = \arg\max_i \left(\cos\langle v_{gt}, \mathcal{T}[i] \rangle\right), \tag{7}$$

$$\mathcal{L}_{max} = 1 - \cos\langle v_{gt}, \mathcal{T}[d] \rangle. \tag{8}$$

Thus, the loss in optimizing the TFCC is

$$\mathcal{L}_T = \lambda_{ent} \mathcal{L}_{ent} + \lambda_{max} \mathcal{L}_{max}. \tag{9}$$

Subsequently, we undertake a joint optimization of the low-dimensional features $\hat{f}$ and the MLP decoder $\mathcal{D}$. Ideally, the feature recovered from low-dimensional feature should closely correlate with the $\widehat{GT}$ feature $v_{gt}$. As a result, we impose a stronger constraint geared towards aligning the entries' logits of the low-dimensional features with the assigned $\widehat{GT}$ entry $d$,

$$\mathcal{L}_{joint} = \|e - \text{onehot}(d)\|_2^2. \tag{10}$$

Finally, to bolster the robustness of this procedure, we introduce an end-to-end regularization, directly optimizing the cosine similarity of 2D semantic feature and corresponding ground truth,

$$\mathcal{L}_{e2e} = 1 - \cos\langle v_{gt}, v \rangle. \tag{11}$$

The comprehensive loss function designated for our semantic field reconstruction process is represented as $\mathcal{L}$,

$$\mathcal{L} = \mathcal{L}_T + \lambda_{joint} \mathcal{L}_{joint} + \lambda_{e2e} \mathcal{L}_{e2e}. \tag{12}$$

## 3.6 Optimizable Semantic-space Hyperplane

Thanks to the vision-language models like CLIP and APE, which align features well in image and text spaces. Our 3D Gaussian semantic field, once trained, supports open-vocabulary 3D queries with any text prompt. Most existing methods enable open-vocabulary queries by computing the cosine similarity between semantic and text features, defined as follows: $\cos(\theta) = \frac{\phi_{img} \cdot \phi_{text}}{\|\phi_{img}\|\|\phi_{text}\|}$, where $\phi_{img}$ and $\phi_{text}$ represent the image and text features, respectively. After normalizing the features, the score can be simplified as $Score = \phi_{img} \cdot \phi_{text}$. The higher the score, the greater the similarity between the two features. By manually setting an empirical threshold $\tau$, regions with score exceeding $\tau$ are retained, thus enabling open-vocabulary queries. The aforementioned process can be conceptualized as a hyperplane separating semantic features into two categories: features of interest and features not of interest, based on the queried text feature and $\tau$. The hyperplane is represented as follows:

$$\widetilde{W}x + \bar{b} = 0. \tag{13}$$

Here $\widetilde{W}$ denotes the queried text feature, $x$ represents semantic features and $\bar{b}$ is the bias derived from $\tau$. However, the empirical parameter $\tau$ is not universally applicable to all queries, often resulting in an inability to precisely locate target areas. Consequently, we propose the Optimizable Semantic-space Hyperplane (OSH). Utilizing a RES model, such as Grounded-SAM [45], we obtain a 2D binary mask of the target area and optimize the hyperplane via one-shot logistic regression. This optimization ensures that the classification results of the hyperplane more closely align with the target area of the query.

As shown on the right side of Figure 2, From a specific camera pose, an RGB image and a feature map are obtained through the rgb and semantic feature rendering processes described in Sec. 3.5, respectively. For a text query $t$, the text encoder of APE generates a text embedding $\phi_{text}$, which is used as the initial weight of the hyperplane $Wx + b = 0$. The Feature Map is classified by the hyperplane, resulting in the prediction of a binary mask $m$. The text query $t$ and the RGB image are processed by the RES Model to generate a binary mask $\hat{m}$ of the target area as the pseudo-label. This mask is subsequently used with $m$ in logistic regression to optimize $W$ and $b$. We fine-tune the OSH with the objective:

$$\mathcal{L}_{OSH} = -\frac{1}{P} \sum_{i=1}^{P} [w \cdot \hat{m}_i \log(\sigma(m_i)) + (1 - \hat{m}_i) \log(1 - \sigma(m_i))], \tag{14}$$

where $P$ denotes all samples, $\sigma(\cdot)$ denotes Sigmoid function. Following the one-shot logistic regression, the optimized Semantic-space Hyperplane can be represented by

$$\widetilde{W}'x + b' = 0. \tag{15}$$

Note that the parameters of the 3D Gaussians remain frozen during this process. The red lines in Figure 2 indicate operations that only occur upon the initial query with a new text prompt. Subsequently, the OSH can be used to delineate regions of interest in

both 2D feature maps rendered from novel views and in 3D Gaussians. Specifically, for a semantic feature $F$, derived either from a 2D semantic feature map at pixel $p$ or from a 3D Gaussian $g$, if $\widetilde{W}'F + b' > 0$, it indicates that $F$ is sufficiently close to the queried text, warranting retention of $p$ or $g$ in the query results set.

## 4 IMPLEMENTATION DETAILS

Our method is implemented based on 3D Gaussian Splatting[20]. We modified the CUDA kernel to render semantic features on the 3D Gaussians, ensuring that the extended semantic feature attributes of each 3D Gaussian support gradient backpropagation. Our model, based on a 3D Gaussian Scene reconstructed via vanilla 3D Gaussian Splatting [20], can be trained on a single 40G-A100 GPU in approximately 10 minutes.

## 5 EXPERIMENTS

### 5.1 Evaluation Setup

**Datasets.** To assess the effectiveness of our approach, we conduct experiments on two datasets: The Mip-NeRF360 dataset [1] and the Replica dataset [50]. Mip-NeRF360 is a high-quality real-world dataset that contains a number of objects with rich details. It is extensively used in 3D reconstruction. We selected four scenes (Room, Bonsai, Garden, and Kitchen), both indoors and outdoors, for our evaluations. Additionally, we designed an open-vocabulary semantic segmentation test set under these scenes. We manually annotated a few relatively prominent objects in each scene, providing their 2D masks and descriptive phrases, such as "sofa in dark green". Replica is a 3D synthetic dataset that features high-fidelity indoor scenes. Each scene comprises RGB images along with corresponding semantic segmentation masks. We conducted reconstruction and evaluation in four commonly used scenes from the Replica dataset [50]: office0, office1, room0, and room1. For a given viewpoint image, our evaluation concentrates on assessing the effectiveness of single-query results within an open-vocabulary context rather than obtaining a similarity map for all vocabularies in a closed set and deciding mask regions based on similarity scores [27, 30, 60]. Therefore, in designing our experiments, we drew inspiration from the methodologies of refCOCO and refCOCOg[58]. For each semantic ground truth in the Replica test set, we cataloged the class names present and sequentially used these class names as text queries to quantitatively measure the performance metrics.

**Baseline Methods and Evaluation Metrics.** To assess the accuracy of open-vocabulary querying results, we employ mean Intersection over Union (mIoU), mean Pixel Accuracy (mPA), and mean Precision (mP) as evaluation metrics. Additionally, to evaluate model performance metrics, we measure the training duration and the rendering time.

### 5.2 Comparisons

We conduct a comparative evaluation of our approach in contrast with LangSplat[40], Gaussian Grouping[57], Feature 3DGS[60], and LERF[21].

**Qualitative Results.** We present the qualitative results produced by our method alongside comparisons with other approaches. Figure 3 offers a detailed showcase of the open-vocabulary query performance on the Mip-NeRF360 test data. It especially highlights

**Table 1: Evaluation metrics for comparing our method with others on Mip-NeRF360 [1] evaluation dataset.**

| Method | mIoU | mPA | mP |
|---|---|---|---|
| LERF [21] | 0.2698 | 0.8183 | 0.6553 |
| Feature 3DGS [60] | 0.3889 | 0.8279 | 0.7085 |
| GS Grouping [57] | 0.4410 | 0.7586 | 0.7611 |
| LangSplat [40] | 0.5545 | 0.8071 | 0.8600 |
| Ours | **0.8646** | **0.9569** | **0.9362** |

the utilization of phrases that describe the appearance, texture, and relative positioning of different objects.

LeRF [21] generates imprecise and vague 3D features, which hinder the clear discernment of boundaries between the target region and others. Feature 3DGS [60] employs a 2D semantic segmentation model LSeg [24] as its feature extractor. However, like LSeg, it lacks proficiency in handling open-vocabulary queries. It frequently queries all objects related to the prompt and struggles with complex distinctions, like distinguishing between a sofa and a toy resting on it. Gaussian Grouping [57] leverages the instance mask via SAM [22] to group 3D Gaussians into 3D instances devoid of semantic information. It uses Grounding DINO [32] to pinpoint regions of interest for enabling 3D open-vocabulary queries. However, this approach leads to granularity issues, often identifying only a fraction of the queried object, such as the major part of "green grass" or the flower stem from the "flowerpot on the table". LangSplat [40] uses SAM to generate object segmentation masks and subsequently employs CLIP to encode these regions. However, this strategy results in CLIP encoding only object-level features, leading to an inadequate understanding of the correlations among objects within a scene. For instance, when querying "the tablemat next to the red gloves", it erroneously highlights the "red gloves" rather than the intended "tablemat". Similar to Gaussian Grouping, LangSplat also encounters granularity issues, such as failing to segment all "green grass" and improperly dividing the "sofa" into multiple parts.

Our methodology is notably effective as it harnesses the power of semantic redundancy to cluster features into a TFCC, enabling the efficient encoding of diverse object features. Consequently, this approach precisely pinpoints objects such as the sofa, grass, and road while maintaining accurate boundaries. Our strategy further excels at discerning the intricate interrelationships among various objects within a scene. Unlike LangSplat, we encode entire images with the image encoder to integrate scene-level information into the semantic features. Additionally, we deploy dynamically optimize a semantic-space hyperplane, effectively filtering out unnecessary objects from the 3D Gaussians of Interest. For instance, in the cases of "flowerpot on the table" and "the tablemat next to the red gloves", we successfully segment the primary subjects of the phrase rather than the secondary objects.

**Quantitative Results.** Table 1 and Table 2 provide a comparative analysis of the efficacy of our work relative to other projects across multiple datasets. As displayed, our segmentation precision significantly exceeds that of LERF and open-vocabulary 3DGS-based methods. We observed a substantial mean Intersection over

Figure 3: Visualization comparisons of open-vocabulary querying results are presented. From top to bottom: Ground truth, querying results from LERF [21], Feature 3DGS [60], Gaussian Grouping [57], LangSplat [40], and our method. From left to right, the images display the querying results corresponding to text descriptions, which are noted at the bottom line.

Table 2: Evaluation metrics for comparing our method with others on Replica [50] evaluation dataset.

| Method | mIoU | mPA | mP |
|---|---|---|---|
| LERF [21] | 0.2815 | 0.7071 | 0.6602 |
| Feature 3DGS [60] | 0.4480 | 0.7901 | 0.7310 |
| GS Grouping [57] | 0.4170 | 0.73699 | 0.7276 |
| LangSplat [40] | 0.4703 | 0.7694 | 0.7604 |
| Ours | **0.6169** | **0.8367** | **0.8088** |

Union (mIoU) improvement of 30% on the Mip-NeRF360 dataset and 12% on the Replica dataset, respectively.

Moreover, Table 3 underscores the effectiveness of our approach. We detail the pre-processing encoding time for extracting 2D semantic feature maps, scene reconstruction duration, total training time, and rendering frame rates for each approach under consideration. By deriving a highly efficient visual encoder from APE, we reduced the image encoding time to ~2 seconds. Furthermore, unlike LERF, Feature 3DGS, and LangSplat, which start training from scratch, both our method and Gaussian Grouping build on 3D semantic fields from scenes that are pre-trained using 3D Gaussian Splatting [20]. To ensure fairness, the time required for pre-training scenes using 3D Gaussian Splatting (25 minutes) is included in our overall training time calculation. Through meticulous TFCC design

**Table 3: Time evaluation for training and rendering on Mip-NeRF360 [1] dataset.**

| Method | Pre-process | Training | Total | FPS |
|---|---|---|---|---|
| LERF [21] | **3min** | 40min | **43min** | 0.17 |
| Feature 3DGS [60] | 25min | 10h 23min | 10h 48min | ~10 |
| GS Grouping [57] | 27min | 25+113min | 165min | **~100** |
| LangSplat [40] | 50min | 99min | 149min | ~30 |
| Ours | 8min | **25+12min** | 45min | ~30 |

**Table 4: Evaluation metrics for ablation studies on Mip-NeRF360 [1] dataset.**

| Setting | mIoU | mPA | mP |
|---|---|---|---|
| Baseline | 0.4753 | 0.8638 | 0.7577 |
| w/o OSH | 0.6282 | 0.9464 | 0.8157 |
| w/o TFCC | 0.7537 | 0.9011 | 0.9115 |
| Full model | **0.8646** | **0.9569** | **0.9362** |

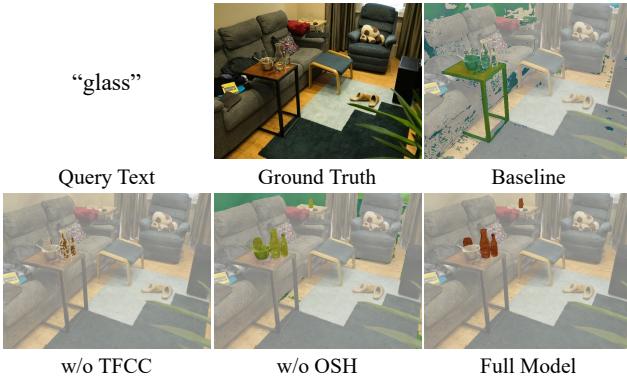and training regularization, we successfully reconstruct a semantic field in under 12 minutes.



Figure 4: Visualization comparison of ablation experiments using the query text "glass".

## 5.3 Ablation Studies

To discover each component's contribution to 3D open-vocabulary scene understanding, a series of ablation experiments are conducted for the Mip-NeRF360 dataset [1] using the same 2D semantic features extracted from APE[48] image encoder. We employ the approach of lifting reduced-dimensionality semantic features into 3D Gaussians as our baseline. This is contrasted with results from models not utilizing the TFCC module, those not employing the OSH module, and the results from the complete model.

As illstrated in Table 4, OSH and TFCC are critical to the effectiveness of our approach; without them, there would be a significant deterioration in performance(-27% ~ -12% mIoU). As shown in Figure 5, the baseline model (middle-left) struggles due to its scattered features, making it difficult for the model with the OST module (middle-right) to identify a suitable hyperplane. In contrast, the

model with TFCC (bottom-left) demonstrates more clustered features and distinct semantic boundaries.

To investigate the impact of 2D foundation models on 3D open-vocabulary understanding, Figure 4 compares the effects of using the CLIP model to extract 2D semantic features against our baseline, which utilizes the APE model for feature extraction. Additionally, the figure illustrates the performance of each setting when integrated with TFCC module proposed by us. The pure CLIP setting struggles with imprecise and vague 3D features, which are alleviated after integrating the TFCC module. Although the baseline setting has more distinct contours, it exhibits disorganized semantic features; however, significant improvement is observed when it is combined with the TFCC module.
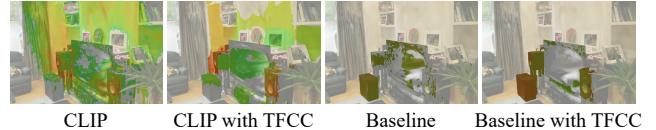


Figure 5: Comparison of different 2D Foundation Models: CLIP and APE, using the query text "speakers".

## 5.4 Application

Our method can be applied to a variety of downstream tasks, with the most direct application being the editing of 3D scenes. As shown in the figure 6, we use the text query "Flowerpot on the table" to locate the 3D Gaussians of interest. Our method enables the highlighting of target areas, localized deletion, and movement. Furthermore, by integrating with Stable-Diffusion[46], We can employ the Score Distillation Sampling (SDS) [38] loss function to achieve high-quality 3D generation tasks in specific areas.
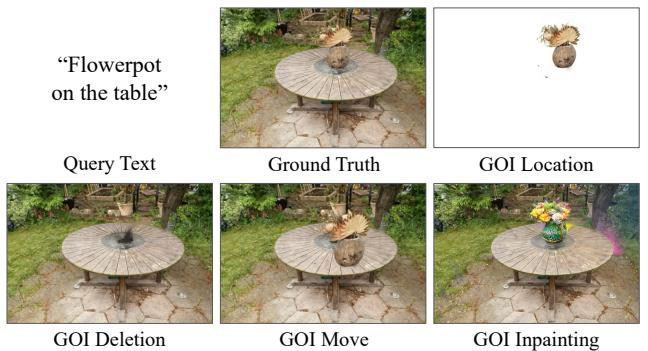


Figure 6: Visualization of scene manipulation results using our method. The query text is used to locate the 3D Gaussians of interest (GOI). "A beautiful vase" is used as the prompt for the 3D inpainting process after locating the GOI.

## 6 CONCLUSION

In this paper, we introduce GOI, a method for reconstructing 3D semantic fields, capable of delivering precise results in 3D open-vocabulary querying. By leveraging the Trainable Feature Clustering Codebook, GOI effectively compresses high-dimensional

semantic features and integrates these lower-dimensional features into 3DGS, significantly reducing memory and rendering costs while preserving distinct semantic feature boundaries. Moreover, moving away from traditional methods reliant on fixed empirical thresholds, our approach employs an Optimizable Semantic-space Hyperplane for feature selection, thereby enhancing querying accuracy. Through extensive experiments, GOI has demonstrated improved performance over existing methods, underscoring its potential for downstream tasks, such as localized scene editing.

# REFERENCES

[1] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 5460–5469. https://doi.org/10.1109/CVPR52688.2022.00539

[2] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *ArXiv preprint* abs/2108.07258 (2021). https://arxiv.org/abs/2108.07258

[3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging Properties in Self-Supervised Vision Transformers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 9630–9640. https://doi.org/10.1109/ICCV48922.2021.00951

[4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging Properties in Self-Supervised Vision Transformers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 9630–9640. https://doi.org/10.1109/ICCV48922.2021.00951

[5] Paola Cascante-Bonilla, Hui Wu, Letao Wang, Rogério Feris, and Vicente Ordonez. 2022. Sim VQA: Exploring Simulated Environments for Visual Question Answering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 5046–5056. https://doi.org/10.1109/CVPR52688.2022.00500

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

[7] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S Ryoo, Austin Stone, and Daniel Kappler. 2023. Open-vocabulary queryable scene representations for real world planning. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 11509–11522.

[8] Forrester Cole, Kyle Genova, Avneesh Sud, Daniel Vlasic, and Zhoutong Zhang. 2021. Differentiable Surface Rendering via Non-Differentiable Sampling. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 6068–6077. https://doi.org/10.1109/ICCV48922.2021.00603

[9] Peng Dai, Yinda Zhang, Xin Yu, Xiaoyang Lyu, and Xiaojuan Qi. 2023. Hybrid neural rendering for large-scale scenes with motion blur. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 154–164.

[10] Qiancheng Fu, Qingshan Xu, Yew Soon Ong, and Wenbing Tao. 2022. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems* 35 (2022), 3403–3416.

[11] Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. 2022. Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In *2022 International Conference on 3D Vision (3DV)*. IEEE, 1–11.

[12] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. 2021. FastNeRF: High-Fidelity Neural Rendering at 200FPS. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 14326–14335. https://doi.org/10.1109/ICCV48922.2021.01408

[13] Jianfei Guo, Nianchen Deng, Xinyang Li, Yeqi Bai, Botian Shi, Chiyu Wang, Chenjing Ding, Dongliang Wang, and Yikang Li. 2023. Streetsurf: Extending multi-view implicit surface reconstruction to street views. *arXiv preprint arXiv:2306.04988* (2023).

[14] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. 2023. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 19740–19750.

[15] Jie Hu, Liujuan Cao, Yao Lu, ShengChuan Zhang, Yan Wang, Ke Li, Feiyue Huang, Ling Shao, and Rongrong Ji. 2021. Istr: End-to-end instance segmentation with transformers. *arXiv preprint arXiv:2105.00637* (2021).

[16] Jie Hu, Linyan Huang, Tianhe Ren, Shengchuan Zhang, Rongrong Ji, and Liujuan Cao. 2023. You only segment once: Towards real-time panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17819–17829.

[17] Jie Hu, Yao Lu, Shengchuan Zhang, and Liujuan Cao. 2024. ISTR: Mask-Embedding-Based Instance Segmentation Transformer. *IEEE Transactions on Image Processing* (2024).

[18] 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 arXiv:2404.13921* (2024).

[19] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. 2023. Visual language maps for robot navigation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 10608–10615.

[20] 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.

[21] Justin* Kerr, Chung Min* Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. 2023. LERF: Language Embedded Radiance Fields. In *International Conference on Computer Vision (ICCV)*.

[22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4015–4026.

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

[24] 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

[25] Xinyang Li, Zhangyu Lai, Linning Xu, Jianfei Guo, Liujuan Cao, Shengchuan Zhang, Bo Dai, and Rongrong Ji. 2024. Dual3D: Efficient and Consistent Text-to-3D Generation with Dual-mode Multi-view Latent Diffusion. *arXiv preprint arXiv:2405.09874* (2024).

[26] 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. *arXiv preprint arXiv:2406.17601* (2024).

[27] 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

[28] Jianghang Lin, Yunhang Shen, Bingquan Wang, Shaohui Lin, Ke Li, and Liujuan Cao. 2024. Weakly supervised open-vocabulary object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 3404–3412.

[29] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. 2023. Weakly supervised 3d open-vocabulary segmentation. *Advances in Neural Information Processing Systems* 36 (2023), 53433–53456.

[30] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. 2023. Weakly supervised 3d open-vocabulary segmentation. *Advances in Neural Information Processing Systems* 36 (2023), 53433–53456.

[31] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499* (2023).

[32] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *ArXiv preprint* abs/2303.05499 (2023). https://arxiv.org/abs/2303.05499

[33] Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. 2022. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *European Conference on Computer Vision*. Springer, 210–227.

[34] Shiyang Li, Haonan Chang, Eric Pu Jing, Abdeslam Boularias, and Kostas Bekris. 2023. Ovir-3d: Open-vocabulary 3d instance retrieval without training on 3d data. In *Conference on Robot Learning*. PMLR, 1610–1620.

[35] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

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

[37] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. *ArXiv preprint* abs/2304.07193 (2023). https://arxiv.org/abs/2304.07193

[38] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *ArXiv preprint* abs/2209.14988 (2022). https://arxiv.org/abs/2209.14988

[39] Sergey Prokudin, Qianli Ma, Maxime Raafat, Julien Valentin, and Siyu Tang. 2023. Dynamic point fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7964–7976.

[40] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. 2023. LangSplat: 3D Language Gaussian Splatting. *ArXiv preprint* abs/2312.16084 (2023). https://arxiv.org/abs/2312.16084

[41] 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.

[42] 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

[43] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 14315–14325. https://doi.org/10.1109/ICCV48922.2021.01407

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

[45] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. 2024. Grounded sam: Assembling open-world models for diverse visual tasks. *ArXiv preprint* abs/2401.14159 (2024). https://arxiv.org/abs/2401.14159

[46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. arXiv:2112.10752 [cs.CV]

[47] William Shen, Ge Yang, Alan Yu, Jansen Wong, Leslie Pack Kaelbling, and Phillip Isola. 2023. Distilled feature fields enable few-shot language-guided manipulation. *ArXiv preprint* abs/2308.07931 (2023). https://arxiv.org/abs/2308.07931

[48] Yunhang Shen, Chaoyou Fu, Peixian Chen, Mengdan Zhang, Ke Li, Xing Sun, Yunsheng Wu, Shaohui Lin, and Rongrong Ji. 2023. Aligning and Prompting Everything All at Once for Universal Visual Perception. *ArXiv preprint* abs/2312.02153 (2023). https://arxiv.org/abs/2312.02153

[49] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. 2023. Language Embedded 3D Gaussians for Open-Vocabulary Scene Understanding. *ArXiv preprint* abs/2311.18482 (2023). https://arxiv.org/abs/2311.18482

[50] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. 2019. The Replica dataset: A digital replica of indoor spaces. *ArXiv preprint* abs/1906.05797 (2019). https://arxiv.org/abs/1906.05797

[51] 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.

[52] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 27171–27183. https://proceedings.neurips.cc/paper/2021/hash/e41e164f7485ec4a28741a2d0ea41c74-Abstract.html

[53] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. 2022. Hf-neus: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems* 35 (2022), 1966–1978.

[54] Yuze Wang, Junyi Wang, and Yue Qi. 2024. WE-GS: An In-the-wild Efficient 3D Gaussian Representation for Unconstrained Photo Collections. *arXiv preprint arXiv:2406.02407* (2024).

[55] Yuze Wang, Junyi Wang, Yansong Qu, and Yue Qi. 2023. Rip-nerf: learning rotation-invariant point-based neural radiance field for fine-grained editing and compositing. In *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*. 125–134.

[56] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5438–5448.

[57] 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

[58] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*. Springer, 69–85.

[59] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J. Davison. 2021. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 15818–15827. https://doi.org/10.1109/ICCV48922.2021.01554

[60] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya You, Zhangyang Wang, and Achuta Kadambi. 2023. Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields. *ArXiv preprint* abs/2312.03203 (2023). https://arxiv.org/abs/2312.03203

[61] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2021. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=gZ9hCDWe6ke

# A ADDITIONAL IMPLEMENTATIONAL DETAILS

Our work is based on pretrained vanilla Gaussian scenes. Subsequent to this fundamental step, we embark on a procedure of semantic field optimization, comprising 1500 iterations. Throughout this period, our principal focus is on the optimization of the semantic field, while maintaining the stasis of other parameters. In this stage, we resort to the default values of the unrelated hyperparameters in 3D Gaussian Splatting [20] for anything outside of semantic field optimization.

## A.1 Trainable Feature Clustering Codebook

We incorporate a low-dimensional semantic feature with 10 dimensions $f$ within each 3D Gaussian. By default, the Trainable Feature Clustering Codebook (TFCC) is configured with $N = 300$ entries. As a result, the input dimension of MLP decoder $\mathcal{D}$ is set to 10, while the output logits $e$ from $\mathcal{D}$ are a 300-dimensional vector. Importantly, the decoder $\mathcal{D}$ is simplified to contain solely a lone fully-connected layer, deemed sufficient for efficacious feature decoding.

In order to augment the efficiency of reconstruction, $k$-means clustering is employed for initializing the TFCC. Between 30 to 50 feature maps are sampled from densely observed viewpoints. Subsequently, for each pixel-wise feature, we adopt the $k$-means clustering based on the cosine similarity amid features.

The resultant loss in the course of the TFCC and low-dimensional feature $f$ optimization is

$$\mathcal{L} = \mathcal{L}_T + \lambda_{joint}\mathcal{L}_{joint} + \lambda_{e2e}\mathcal{L}_{e2e}$$
$$= \lambda_{ent}\mathcal{L}_{ent} + \lambda_{max}\mathcal{L}_{max} + \lambda_{joint}\mathcal{L}_{joint} + \lambda_{e2e}\mathcal{L}_{e2e}, \quad (16)$$

We allocate a weightage of $\lambda_{ent} = 0.3$ for $\mathcal{L}_{ent}$, whilst the remainder are set as 1. The annealing temperature $\tau$ derived from $\mathcal{L}_{ent}$ begins at 1, escalating to 2 post 1000 iterations.

## A.2 Optimizable Semantic-space Hyperplane

We use the Grounded-SAM [45] model as our Referring Expression Segmentation (RES) model. The text query $t$ and the RGB image are processed by the RES model to generate a binary mask $\hat{m}$ of the target area as the pseudo-label. This mask is subsequently used with $m$ in logistic regression to optimize $W$ and $b$. We fine-tune the OSH with the objective:

$$\mathcal{L}_{OSH} = -\frac{1}{P}\sum_{i=1}^{P}[w \cdot \hat{m}_i \log(\sigma(m_i)) + (1 - \hat{m}_i)\log(1 - \sigma(m_i))], \quad (17)$$

where $P$ denotes all samples, $\sigma(\cdot)$ denotes Sigmoid function, $w$ is a hyperparameter. Considering that regions of interest tend to be significantly smaller than non-interest regions, we set $w = \frac{1}{10}$ to increase the penalty weight for misclassifying target areas, thereby accelerating convergence.

# B EXPERIMENTAL DETAILS

## B.1 Expanding the Mip-NeRF360 Dataset

Within each of the four selected scenes (Room, Bonsai, Garden, and Kitchen) from the Mip-NeRF360 dataset [1], we've identified four notably distinctive objects. For every individual object, we've

established ten distinct viewpoints in the scenario, and employed the SAM [22] ViT-H model to generate object masks for these preselected perspectives. Moreover, we present textual descriptions founded on either the appearance of the chosen objects (e.g., "sofa in dark green"), or their spatial relationship with other objects (e.g., "table under the bowl"). Consequently, our expanded evaluation set for Mip-NeRF360 includes tuples encapsulating the viewpoint image, ground truth mask, and a concise text description.

We have listed the textual descriptions of each individual object selected within the scenes in Table 5. Additionally, in Figure 8, we exhibit the ground truth segmentation masks pertinent to select objects in our expanded Mip-NeRF360 evaluation dataset.

| Scene | Text Description |
|---|---|
| Room | bowl on the table, brown slipper, sofa in dark green, table under the bowl |
| Bonsai | black chair, flowerpot on the table, orange bottle, purple table |
| Garden | brown table, flowerpot on the table, green football, green grass |
| Kitchen | chair, red gloves, table mat, wooden table |

Table 5: Text description for select objects of each scene in our extended version of the Mip-NeRF360 evaluation dataset.

## B.2 More Results

*B.2.1 Qualitative Results.* Figure 7 serves as a visual representation of our comprehensive query results derived from the Mip-NeRF360 dataset. The effect of executing queries on an identical object, but from varying viewpoints, is lucidly demonstrated. The takeaway is that our outcomes have effectively demarcated the object boundaries and simultaneously exhibited consistency when observed from multiple viewpoints.

*B.2.2 Quantitative Results.* We base our evaluation on metrics such as mean Intersection over Union (mIoU), mean Pixel Accuracy (mPA), and mean Precision (mP), akin to the LEGaussian [49] method. The efficiency and efficacy of our approach have previously been demonstrated. Furthermore, Tables 6 and 7 provide a detailed exposition of our scene-level metrics derived from the Mip-NeRF360 [1] and Replica [50] datasets. Notably, our proposed methodology consistently outperforms, irrespective of the scene encompassing the datasets. Additionally, we provide a video that juxtapose our methodology with others, facilitating a more effective elucidation of our superior performance.

## B.3 3D Manipulations

As addressed in Sec. 3.5, the low-dimensional feature $f$ in 3D Gaussians and the rendered 2D pixel-wise feature $\hat{f}$ are fundamentally equivalent. We can also retrieve the high-dimensional semantic feature $v$ for the feature $f$, as depicted in the following equation.

$$v = \mathcal{T}\left[\underset{j=1,2,\dots,N}{\operatorname{argmax}}(e_j)\right], \text{ where } e = \mathcal{D}(f) \quad (18)$$
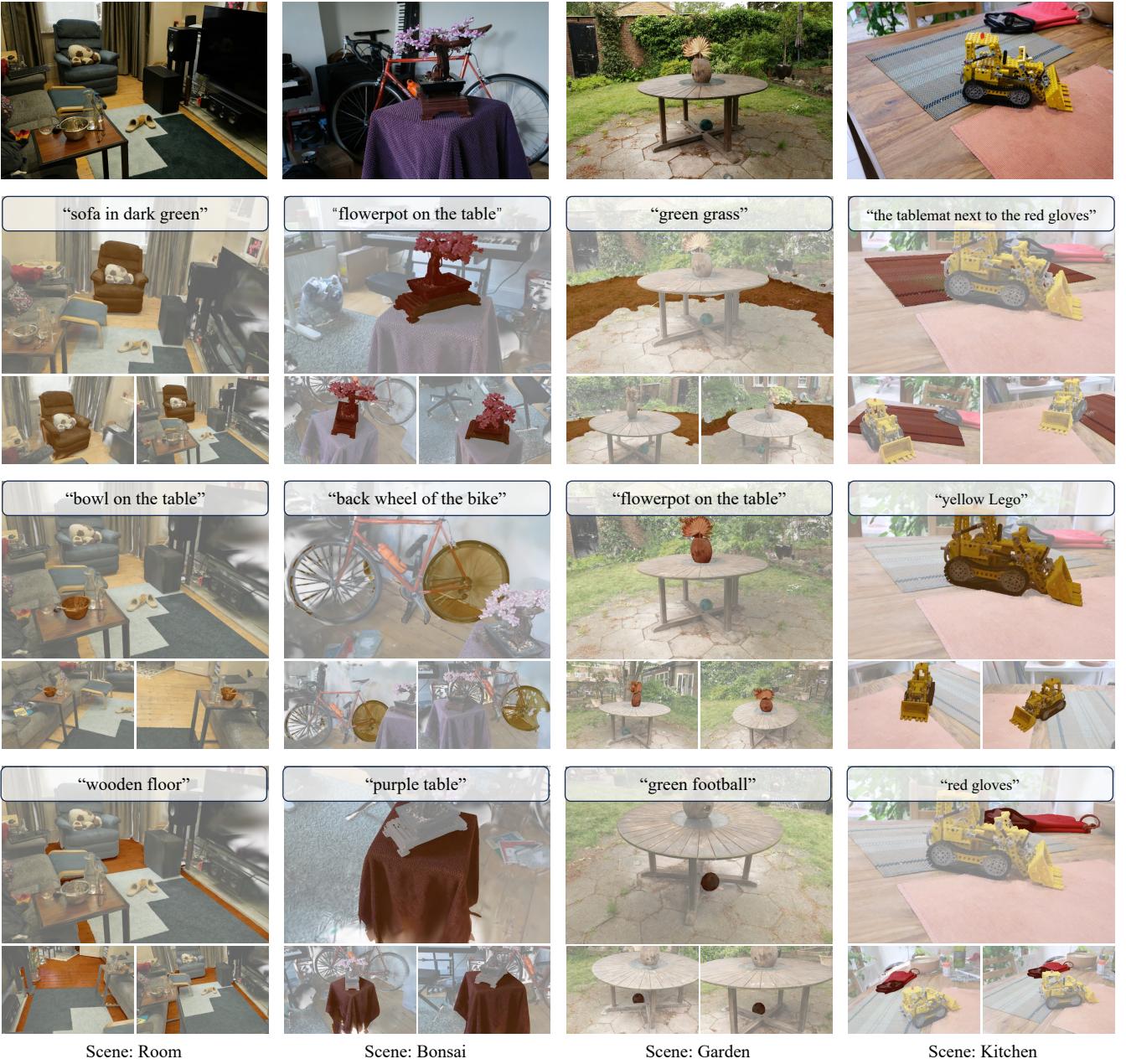
**Figure 7: Extensive query visualization on the Mip-NeRF360 dataset. In each column, the images delineated on the top row and the descriptions in the bottom line typify the scene under examination. Within each depicted scene, we have identified three distinct objects to constitute our query. Three distinctive viewpoints from the same scene are exhibited for every given prompt.**

wherein $\mathcal{T}$ and $\mathcal{D}$ are the TFCC and the MLP decoder, and the subscript $j$ iterates over the elements of the logits $e$, ascending from 1 up to its length $N$.

Through this process, we are able to comprehend the 3D Gaussian-level semantic feature. Subsequently, via the Optimizable Semantic-space Hyperplane, we can effectively extract the Gaussians of interest. Consequently, our GOI approach can be harnessed for downstream tasks, enabling efficient 3D manipulations such as deletion, localization, and inpainting.
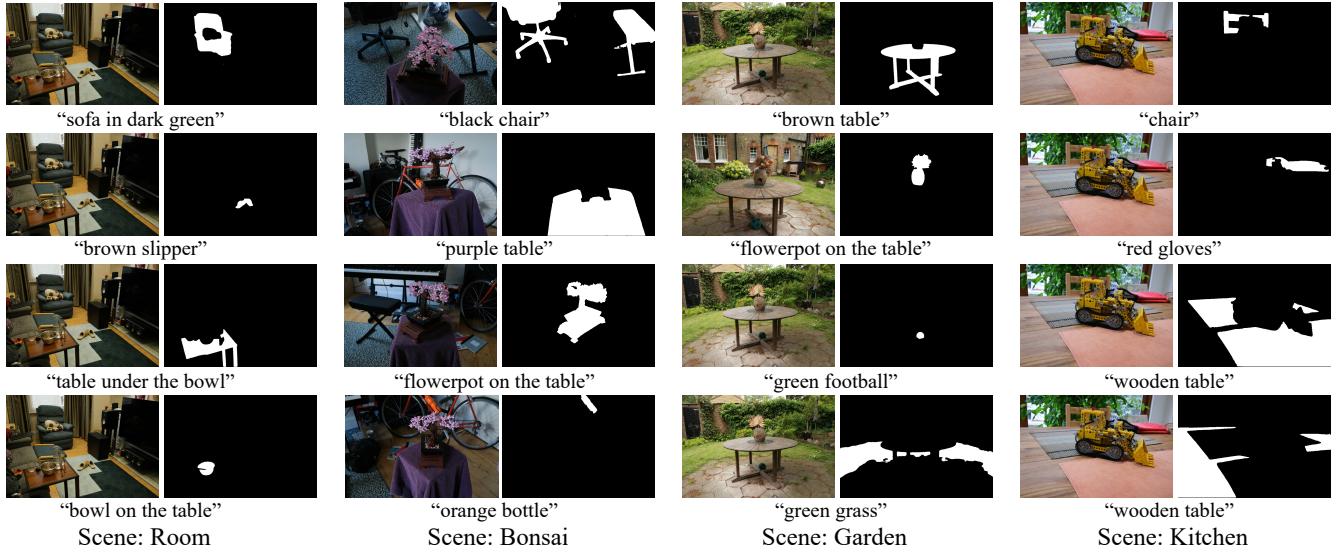
"sofa in dark green"　　　"black chair"　　　"brown table"　　　"chair"

"brown slipper"　　　"purple table"　　　"flowerpot on the table"　　　"red gloves"

"table under the bowl"　　　"flowerpot on the table"　　　"green football"　　　"wooden table"

"bowl on the table"　　　"orange bottle"　　　"green grass"　　　"wooden table"

Scene: Room　　　Scene: Bonsai　　　Scene: Garden　　　Scene: Kitchen

**Figure 8: Ground truth segmentation masks for select objects in our extended version of the Mip-NeRF360 evaluation dataset.**

| Scene | Metric | Works | | | | |
|-------|--------|-------|------|------|------|------|
| | | LERF [21] | Feat. 3DGS [60] | GS Grouping [57] | LangSplat [40] | Ours |
| Room | mIoU | 0.0806 | 0.1748 | 0.4909 | 0.6263 | **0.8504** |
| | mPA | 0.8458 | 0.8246 | 0.8190 | 0.9104 | **0.9718** |
| | mP | 0.5400 | 0.5919 | 0.7663 | 0.8442 | **0.9485** |
| Bonsai | mIoU | 0.3214 | 0.4623 | 0.4305 | 0.5914 | **0.9147** |
| | mPA | 0.8852 | 0.8027 | 0.8244 | 0.8083 | **0.9630** |
| | mP | 0.6603 | 0.7793 | 0.7926 | **0.9338** | 0.9129 |
| Garden | mIoU | 0.2986 | 0.4507 | 0.4203 | 0.5006 | **0.8499** |
| | mPA | 0.8586 | 0.8863 | 0.6825 | 0.7579 | **0.9577** |
| | mP | 0.6504 | 0.7774 | 0.7302 | 0.8227 | **0.9312** |
| Kitchen | mIoU | 0.3788 | 0.4678 | 0.4222 | 0.4995 | **0.8434** |
| | mPA | 0.6837 | 0.7981 | 0.7085 | 0.7517 | **0.9351** |
| | mP | 0.7708 | 0.6853 | 0.7152 | 0.8392 | **0.9520** |
| Average | mIoU | 0.2698 | 0.3889 | 0.4410 | 0.5545 | **0.8646** |
| | mPA | 0.8183 | 0.8279 | 0.7586 | 0.8071 | **0.9569** |
| | mP | 0.6553 | 0.7085 | 0.7511 | 0.8600 | **0.9362** |

**Table 6: Per-scene and average performance on the Mip-NeRF360 dataset**

| Scene | Metric | Works | | | | |
|-------|--------|-------|--|--|--|--|
| | | LERF [21] | Feat. 3DGS [60] | GS Grouping [57] | LangSplat [40] | Ours |
| Room 0 | mIoU | 0.3095 | 0.4980 | 0.5937 | 0.4843 | **0.6589** |
| | mPA | 0.7761 | 0.8499 | 0.8872 | 0.8134 | **0.9039** |
| | mP | 0.6622 | 0.7484 | 0.8241 | 0.7734 | **0.8301** |
| Room 1 | mIoU | 0.3573 | 0.4244 | 0.4525 | 0.5819 | **0.8020** |
| | mPA | 0.7974 | 0.7826 | 0.7480 | 0.8205 | **0.9383** |
| | mP | 0.6810 | 0.7260 | 0.7667 | 0.8694 | **0.9314** |
| Office 0 | mIoU | 0.2962 | 0.5513 | 0.3388 | 0.4471 | **0.5042** |
| | mPA | 0.6736 | 0.8415 | 0.6664 | **0.7700** | 0.7597 |
| | mP | 0.7004 | 0.7786 | 0.7135 | **0.7395** | 0.7384 |
| Office 1 | mIoU | 0.1630 | 0.3181 | 0.2829 | 0.3682 | **0.5024** |
| | mPA | 0.5812 | 0.6865 | 0.6460 | 0.6736 | **0.7443** |
| | mP | 0.5971 | 0.6710 | 0.6060 | 0.6592 | **0.7353** |
| Average | mIoU | 0.2815 | 0.4480 | 0.4170 | 0.4704 | **0.6169** |
| | mPA | 0.7071 | 0.7901 | 0.7369 | 0.7694 | **0.8365** |
| | mP | 0.6602 | 0.7310 | 0.7276 | 0.7604 | **0.8088** |

Table 7: Per-scene and average performance on the Replica dataset