# FMGS: Foundation Model Embedded 3D Gaussian Splatting for Holistic 3D Scene Understanding

Xingxing Zuo[1], Pouya Samangouei[1], Yunwen Zhou[1], Yan Di[1], Mingyang Li[1]

[1]Google.

Contributing authors: xingxingzuo@google.com; samangouei@google.com; verse@google.com; yanditum@google.com; mingyangli@google.com;

**Abstract**

Precisely perceiving the geometric and semantic properties of real-world 3D objects is crucial for the continued evolution of augmented reality and robotic applications. To this end, we present Foundation Model Embedded Gaussian Splatting (FMGS), which incorporates vision-language embeddings of foundation models into 3D Gaussian Splatting (GS). The key contribution of this work is an efficient method to reconstruct and represent 3D vision-language models. This is achieved by distilling feature maps generated from image-based foundation models into those rendered from our 3D model. To ensure high-quality rendering and fast training, we introduce a novel scene representation by integrating strengths from both GS and multi-resolution hash encodings (MHE). Our effective training procedure also introduces a pixel alignment loss that makes the rendered feature distance of same semantic entities close, following the pixel-level semantic boundaries. Our results demonstrate remarkable multi-view semantic consistency, facilitating diverse downstream tasks, beating state-of-the-art methods by **10.2** percent on open-vocabulary language-based object detection, despite that we are **851×** faster for inference. This research explores the intersection of vision, language, and 3D scene representation, paving the way for enhanced scene understanding in uncontrolled real-world environments. We plan to release the code on the [project page].

**Keywords:** Gaussian Splatting, Vision-Language Embeddings, Foundation Models, Open-Vocabulary Semantics

## 1 Introduction

3D scene understanding is a critical task in various computer vision and robotics applications. Yet, most existing methods primarily concentrate on either 3D geometry and appearance estimation [45, 38, 26] or 3D object detection and scene segmentation trained on datasets with closed sets of classes [15, 19, 41]. However, for an intelligent agent to interact smoothly with the physical world, merely understanding a subset of the space characterized by pre-identified labels is insufficient. Inspired by the latest advancements in foundation models (FMs) with impressive language and vision semantics [44, 1], this paper endeavors to develop a more natural 3D scene representation supporting open-world visual recognition and understanding. It integrates both geometric and open-vocabulary semantic information, facilitating easy querying for downstream tasks such as object detection and semantic segmentation in open-world scenarios.

In this paper, we utilize Gaussian Splatting [26] as the backbone for reconstructing 3D geometry and appearance, which has demonstrated superior performance in terms of rendering quality for novel-view image synthesis and training efficiency. To assist open-vocabulary 3D scene understanding, we rely on pre-train 2D vision-language CLIP [44] and lift the corresponding information into 3D by a novel multi-view training procedure. We note that, in research communities, the system that is most similar to us is LEFR [27], which integrates implicit NERF [38] based scene representation and CLIP embeddings. Compared to LERF, our system develops a different architecture, provides a variety of technical contributions ranging from high efficiency to 3D consistent query, and obtains significantly better results (approximately 10.2 percent in representative key metrics).

A straightforward approach to enhance 3D Gaussian Splatting with vision-language FM embeddings is to attach each Gaussian with a learnable feature vector, which can be trained through image rasterization to formulate loss functions. However, maintaining high-quality rendering with GS typically requires millions of Gaussians in a nominal room-scale environment. Employing per-Gaussian feature vectors inevitably results in excessive memory consumption and significantly slows down training, limiting the practical applications of this system. Motivated by iNGP [40], we model our system by using 3D Gaussians together with multi-resolution hash encoding (MHE) to distill the foundation model embeddings. Specifically, to obtain the language embedding from the Gaussians, we utilize their mean values to query the MHE field at corresponding positions. Subsequently, this queried MHE is processed through a Multi-Layer Perceptron (MLP) to generate the output language embedding.

In the training phase, we employ a supervision mechanism on the MHE-based language FM CLIP feature field using a hybrid feature map. This map is derived from the average of multi-scale image crops obtained from various viewpoints. This approach enables the embedding to effectively capture language features corresponding to each scale ensuring a comprehensive representation. For instance, the embedding might represent a 'red book' when viewed up-close, while depicting a 'library' from a more distant perspective. It is noteworthy that CLIP embeddings are designed to encapsulate the overall concept presented in a 2D image, exhibiting minimal variation across individual pixels. Additionally, CLIP embeddings are not perfectly multi-view consistent, i.e., when a 3D object observed by a moving camera via different views, the difference between computed CLIP embeddings across frames are not explicitly minimized. To solve the above-mentioned problems, we rely on multi-view consistency training process to ensure that 3D models, when rendered from different image views, exhibit minimal variations. Additionally, to allow pixel-aligned query experience, DINO [7] embeddings are used together with CLIP embeddings similar to LERF [27]. By carefully analyzing the properties in both CLIP and DINO embeddings, we design an additional pixel alignment loss to further improve the object localization and scene understanding capabilities. This loss is grounded in the dot product similarity of CLIP/DINO features between the central pixel and its surroundings, guiding the rendered CLIP feature map to replicate the same similarity pattern observed in the DINO feature map.

This research paves the way for enhanced real-world applications, such as augmented reality experiences where users can interact with objects using natural language and robotic systems that can navigate and manipulate environments based on linguistic commands. By bridging the gap between language and 3D representation, FMGS opens up new possibilities for understanding and interacting with our surroundings.

Our contributions can be summarized as follows:

- Novel semantic scene representation: We introduce a novel approach combining 3D Gaussians (parameterized by mean, covariance, opacity, and spherical harmonics) for geometry and appearance representation, with MHE for efficient semantic embedding. This approach addresses memory constraints in room-scale scenes including millions of 3D Gaussians.
- Multi-view consistent language embeddings: Our training process utilizes Gaussian-splatting based rendering from multiple views, ensuring consistency across 3D space in static scenarios.

Language embeddings remain invariant to viewpoints, enforcing local proximity consistency within Gaussian volumes.

- Addressing pixel misalignment: We address pixel alignment challenges of CLIP features by extracting and aggregating them at multiple resolutions for a hybrid CLIP feature, which is used for supervising the training. Regularization with pixel-aligned DINO features and a novel dot-product similarity loss enhances spatial precision and object differentiation.

- State-of-the-art performance: Our methods demonstrate superior performance in open-vocabulary semantic object localization, outperforming existing state-of-the-art approaches with quantitative and qualitative results by a wide margin, despite being hundreds of times faster.

## 2 Related Works

We review three main areas of related articles: 3D scene representation, open-vocabulary object recognition and scene understanding, and combined 3D scene representation and semantic understanding.

### 3D Scene Representation

Scene representation in 3D can be roughly categorized by mesh based, voxel based, point based, and implicit ones. Voxel based methods typically discretize 3D space into regular grid cell elements where each grid cell corresponds to a voxel. To estimate the dense 3d voxel cells, probabilistic fusion methods were firstly [22] used and researchers also developed end-to-end learn-able methods [50], by using either depth sensors [22] or monocular camera systems [59]. To visualize estimated voxel fields, they are typically converted into a mesh based representation. This enables efficient rendering on modern computer graphics systems. While alternative methods, such as those using 3D meshes [47, 32], have achieved notable success in various fields, their discrete scene representation, whether voxel-based or mesh-based, imposes limitations on the ability to achieve photo-realistic reconstruction and rendering performance.

Neural implicit representation, e.g., NeRF series [38, 4, 5, 6], represent 3D scenes by fully-connected neural networks, in which volume density and radiance can be queried by input position and view direction vectors. To improve the training and rendering efficiency of NeRFs, 3D space can be discretized by using MHE similar to the concept used in voxel based methods [40]. TensoRF [10] models radiance fields as 4D tensors, factorizing them into compact low-rank tensor components using CP decomposition and introducing novel vector-matrix (VM) decomposition for improved rendering quality, reduced memory footprint, and faster reconstruction.

Finally, point-based methods are originally widely used for directly processing data from depth sensors, for performing geometrical and semantic computer vision tasks [42, 24]. Point-NeRF [56] efficiently combines point cloud and NeRF to achieve impressive fast view synthesis results. Recently, 3D Gaussian Splatting (GS) has been proposed to model points as 3D Gaussians for scene representation [26], and achieved state-of-the-art novel view synthesis rendering quality. However, in [26], the number of Gaussians used for scene representation can easily surpass one million, which introduces strict memory and computational requirements for downstream use cases.

### Open-Vocabulary Object Detection and Scene Understanding

Advancements in open-vocabulary object detection in 2D images have been made by leveraging natural language prompts. LSeg [30] employs a text encoder for semantic label embeddings and a transformer-based image encoder for dense pixel embeddings, using contrastive alignment to achieve zero-shot image segmentation and generalization to unseen categories. CRIS [55] leverages CLIP for image segmentation, employing a vision-language decoder to align text and pixel-level features, and text-to-pixel contrastive learning to enforce similarity between text and relevant pixel features. CLIP-Seg [37] leverages CLIP as a backbone, employs a transformer-based decoder for dense prediction, and generates image segmentation based on arbitrary text or image prompts. OV-Seg [31] improves open-vocabulary semantic segmentation by finetuning CLIP on masked image regions and text descriptions from noisy

captions, achieving promising performance without dataset adaptations.

Current approaches often employ region proposal or mask prediction methods to guide open-vocabulary classification models. OpenSeg [17] employs mask representations to facilitate visual grouping and align captions with predicted segmentation masks for open-vocabulary image segmentation. ViLD [21] advances open-vocabulary object detection by distilling knowledge from a pretrained image classification model (teacher) into a two-stage detector (student), aligning region embeddings of detected boxes with text and image embeddings inferred by the teacher. Detic [65] expands object detectors' vocabulary by training their classifiers on image classification data, outperforming prior methods on open-vocabulary and long-tail detection benchmarks, achieving generalization to new datasets without finetuning and enabling detectors trained on all ImageNet classes. OVIR-3D [35] enables open-vocabulary 3D object instance retrieval by fusing text-aligned 2D region proposals into 3D space, leveraging 2D datasets.

Open-vocabulary scene understanding has also been explored by using point cloud as sensor inputs. PointCLIP [63] aligns CLIP-encoded point cloud with 3D category texts, transferring knowledge from 2D to 3D recognition by projecting point cloud into multi-view depth maps, using an inter-view adapter for global feature extraction and few-shot knowledge fusion. ULIP series [57, 58] learn a unified representation for images, texts, and 3D point cloud by leveraging pre-trained vision-language models and automatically synthesized triplets, improving the performance of various 3D backbones. Lu et al. [36] leverage pre-trained image and vision-language models and cross-modal contrastive learning for open-vocabulary 3D point cloud detection without 3D annotations.

### Combined 3D Scene Representation and Semantic Understanding

Language has been incorporated into 3D scene understanding in various ways. For the task of visual question answering, systems like iQA [18], ScanQA [3], and SimVQA [8] leverage 3D information to answer queries about the environment. For object recognition enhancement, language and shape information can be combined to improve object recognition, as seen in [14] and [52].

Inspired by the success of implicit neural reconstruction [38, 4, 5], researchers also start to explore incorporating language guidance into 3d neural scene representation. LERF [27] enables open-ended language queries in 3D by incorporating language embeddings from models, e.g. CLIP, into NeRF. 3D-OVS [34] leverages pre-trained CLIP and DINO models in a weakly supervised manner, distilling multi-modal knowledge and object reasoning into a neural radiance field (NeRF) for segmentation task.

Tschernezki et al. [54] leverage a pre-trained 2D image feature extractor to train a 3D student network, boosting performance in analyzing multiple images forming a 3D scene. FFD [29] tackles scene editing by distilling knowledge from pre-trained 2D image feature extractors into a 3D feature field that guides local editing based on user queries. VL-Fields [53], a neural implicit spatial representation fusing scene geometry and vision-language features, enables open-vocabulary semantic queries without requiring prior object class knowledge. FeatureNeRF [60] distills pre-trained vision models (DINO, Latent Diffusion) to learn generalizable NeRFs, leveraging neural rendering for 2D-to-3D mapping and extracting deep features from NeRF MLPs.

Additionally, ConceptFusion [23] enables open-set and multimodal reasoning in 3D scene representations by fusing foundation model features with SLAM and multi-view fusion. ConceptGraphs [20] leverages 2D foundation models and multi-view association to capture semantic and spatial relationships for efficient task-driven planning. OpenMask3D [51] aggregates per-mask features using the multi-view fusion of CLIP-based image embeddings guided by predicted class-agnostic 3D instance masks. SA3D [9] enables 3D segmentation of target objects in neural radiance fields (NeRF) through one-shot manual prompting, leveraging density-guided inverse rendering, cross-view self-prompting, and an iterative process to project 2D segmentation masks onto 3D mask grids. PVLFF [12] generates a scene's feature field, combining vision-language and hierarchical instance features through contrastive loss from 2D instance segment proposals.

CLIP-Fields [48] learns a spatial mapping to semantic embeddings via weak supervision from

web-trained language and vision models, enabling tasks like object identification and robot navigation without direct human labeling. GNFactor [62], a multi-task robotic manipulation agent, leverages a shared 3D voxel representation and language-augmented neural fields for generalizable visual behavior cloning.

Our work is close and directly comparable to LERF [27] in terms of assumptions about information available at training phase and query time. For example, it does not assume a priori knowledge of query categories at training time which is assumed 3D-OVS [33].

Recently, several concurrent works have emerged, investigating the distillation of semantic features into GS scene representations [11, 16]. LEGaussians [49] introduces a method for quantizing high-dimensional concatenated CLIP and DINO features into compact ones to conserve memory, and attach to each individual Gaussian. Langsplat [43] segments images by SAM [28] and then inputs hierarchical semantic segments into CLIP to extract semantic features for the segments. To address memory constraints, Langsplat incorporates a scene-specific language autoencoder to encode CLIP features into lower dimensions, which are also attached to each individual Gaussian. Feature 3DGS [64] distills pixel-aligned features from 2D foundation models, such as SAM [28] and LSeg [30], into GS by associating each Gaussian with a learnable vector. However, this suffers from a loss of semantic understanding capability, particularly for long-tail semantics [27]. Different from all the aforementioned methods, which attach semantic features to each Gaussian, we propose to seamlessly integrate 3D Gaussian scene representation together with MHE for efficient semantic encodings. Notably, our method does not necessitate additional scene-specific quantization or autodecoder steps, thereby preserving the semantic features' representation capability from foundation models.

# 3 Background Methods

## 3.1 3D Gaussian Splatting

GS [26] represents an environment using a set of 3D Gaussians, each defined by a mean $\mu \in \mathbb{R}^3$, an anisotropic covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{3\times3}$, an alpha value $\alpha \in [0,1]$ representing opacity, and spherical harmonics coefficients (SH). Given a 3D position $\mathbf{x} \in \mathbb{R}^3$, the probability density function of 3D Gaussian is defined as:

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \mathbf{\Sigma}^{-1}(\mathbf{x}-\mu)} \quad (1)$$

where $(\cdot)^T$ represents a transpose operation and $(\cdot)^{-1}$ denotes matrix inversion. To render 3D Gaussians in 2D, we project their mean positions by point projection, and project their covariance using the following equation:

$$\mathbf{\Sigma}' = \mathbf{J}\mathbf{W}\,\mathbf{\Sigma}\,\mathbf{W}^T\mathbf{J}^T \quad (2)$$

where $\mathbf{W} \in \mathbb{R}^{3\times3}$ is the viewing transformation and $\mathbf{J} \in \mathbb{R}^{3\times3}$ is the Jacobian of the affine approximation of the projective transformation [66]. To optimize covariance matrices, we use an equivalent representation:

$$\mathbf{\Sigma} = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T \quad (3)$$

where $\mathbf{R} \in \mathbb{R}^{3\times3}$ and $\mathbf{S} \in \mathbb{R}^{3\times3}$ are rotation and scaling matrices, respectively. GS also includes spherical harmonics coefficients to model the appearance of the scene. Gradients for all parameters are derived explicitly to avoid overhead during training.

Each Gaussian encodes the color $c$ using spherical harmonics, which gives a value depending on the viewing directions. The $\alpha-$blending point-based rendering for a pixel color $\mathbf{c}$ is done by blending $\mathcal{N}$ points in the depth order from front to back:

$$\mathbf{c} = \sum_{i \in \mathcal{N}} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1}(1-\alpha_j), \quad (4)$$

where $\alpha_i$ is given by a 2D Gaussian multiplied by a learned per Gaussian opacity [61].

Note that although the image rendering model is similar across NeRFs and GS, the rendering algorithm is much more efficient in GS. NeRFs need to march along the ray to integrate volume, however, GS rendering uses a point-based $\alpha-$blending approach. This allows GS to include a real-time rendering solution that leverages GPU sorting algorithms and draws inspiration from tile-based rasterization. By using a 3D Gaussian representation, anisotropic splatting can be performed while respecting visibility order. This is achieved through sorting and alpha-blending.

Additionally, a fast and accurate backward pass is enabled by tracking the traversal of sorted splats.

## 3.2 Multi-resolution Hash Encoding

Representing a 3D feature field can have many forms. A naive method is to attach a feature vector (or multiple) to each Gaussian, which can be optimized along with other Gaussian parameters (position, covariance, and so on). However, this is extremely costly in terms of computational cost and memory consumption especially when a large number of Gaussians are generated for scene representation. In fact, adding a $512 \times 1$ feature vector per Gaussian will increase the number of optimized parameters to be $9.83 \times$ under authentic GS parameterization [26] (10 geometric parameters and 48 spherical harmonic appearance parameters per Gaussian) and $65.0 \times$ under simplified GS parameterization [25] (4 geometric parameters and 4 appearance parameters per Gaussian).

To mitigate this problem, we are motivated by multi-resolution hash embedding (MHE) [40], which provides efficient scene representation that consists of two trainable components. The first component first hashes a given position $\mathbf{x} \in \mathbb{R}^3$, and then looks up into a trainable hash table for the corresponding embedding. The second component is an MLP that takes the corresponding embeddings and makes predictions such as color and density. The representation contains multiple hash tables, one per each scale. Specifically, MHE first encodes a given position $\mathbf{q} = MHE_\theta(\mathbf{x})$. To do so, it contains a hash table with $L$ levels. Each level contains up to $E$ feature vectors with dimensionality $D$. Resolution of each level is determined by $N_l = \lfloor N_{\min} \cdot b^l \rfloor$ where $N_{\min}$ is the coarsest resolution, $N_{\max}$ is the finest resolution, and $b$ is a growth factor.

To get $\mathbf{q}$ for a given position $\mathbf{x}$, we query MHE at all scales and concatenate the resulting features. For each scale, we find the enclosing voxel for $\mathbf{x}$. Then, each corner entry of the voxel is mapped into a feature vector with dimensionality $D$ according to the trainable hash table. MHE trilinearly interpolates the queried corner entries according to their relative position of $\mathbf{x}$ in its hypercube for each level. This ensures the continuity of the encoded input and its composition with the neural network, avoiding grid-aligned discontinuities and blocky appearance. After this mapping is done, the features from all scales are concatenated to each other, and the auxiliary inputs $\psi \in \mathbb{R}^K$ which results in a feature vector $\mathbf{q}$ of size $L \times D + K$. The resulting encoding then goes to the second component which is an MLP network, $MLP_\Phi(\mathbf{q})$, produces the final output. This architecture significantly reduces the number of weights that are trained for each view while having an $O(1)$ GPU look up for hashing. Overall this results in significant improvements in quality and speed of training.
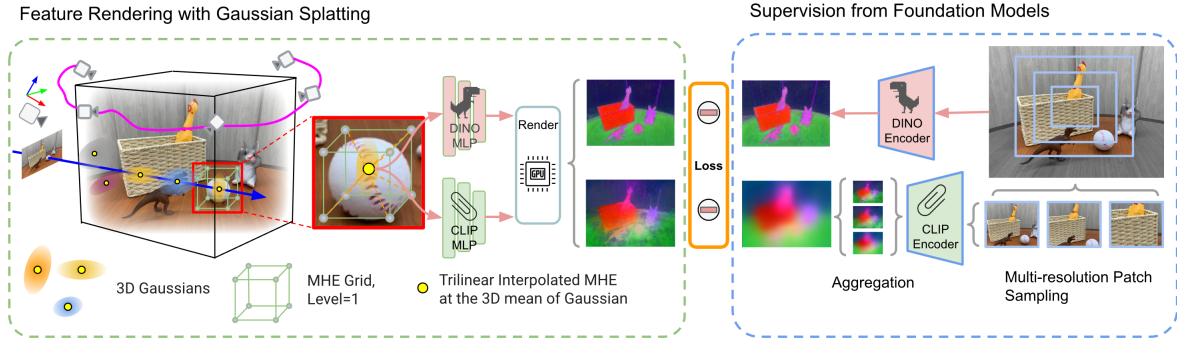
## 4 Method

Our method, i.e. Foundation Model Embedded Gaussian Splatting (FMGS), leverages strengths of both GS and MHE. We rely on GS for efficient and accurate scene geometry representation and on MHE for representing the scene's language content in a light-weighted manner. Given a set of input images, we compute the corresponding camera poses and 3D sparse visual points using an off-the-shelf structure from motion system, e.g., COLMAP [46]. After that we train GS and acquire 3D Gaussians.

Subsequently, we train the feature embedding field (MHE) in 3D by grounding 2D CLIP embeddings. This requires us to generate pixel-aligned features on a set of calibrated input images. However, CLIP embeddings are global in nature and not suitable for pixel-aligned feature extraction. To overcome this challenge, we introduce a framework to learn a volumetric language embedding field that embeds over the 3D Gaussians. The field effectively generate features that is the average CLIP features across all views that include that 3D Gaussian. To supervise our dense feature field, we create a hybrid feature map based on CLIP embeddings across multi-scale crops of training views. Figure 1 provides an overview of our training pipeline.

### 4.1 Feature Field Architecture

3D Gaussian Splatting produces millions of Gaussians to enable high quality rendering of a room-scale scene. This makes it very inefficient to have one CLIP feature per Gaussian since these features are high dimensional and keeping all of these features in GPU memory is not feasible.

**Fig. 1 FMGS Training pipeline:** *Left*: Shows how FMGS' feature field renders CLIP and DINO feature maps for loss calculation. The feature field is a multi-resolution hash encoder (MHE) [40] that embeds semantic information into 3D Gaussians acquired from 3D Gaussian Splatting [26]. *Right*: Shows the target DINO feature map and hybrid CLIP feature map from the foundation models. Note, for visualization simplicity, we only show a single-level MHE here but in implementation we have used multiple levels and concatenate their encodings.

To this end, we parameterize our feature field efficiently using MHE. For a given 3D Gaussian $G(\mathbf{x})$ with mean position $\mathbf{x}$, we first encode $\mathbf{x}$ to a feature vector $\mathbf{q} = MHE_\theta(\mathbf{x})$ where $\theta$ is our multi-resolution hash table parameters. We subsequently feed this output into an MLP, which generates our language embedding $\hat{\mathbf{f}} = MLP_\phi^{CLIP}(\mathbf{q})$, with $\hat{\mathbf{f}}$ belonging to $\mathbb{R}^D$. We also normalize $\hat{\mathbf{f}}$ to make it a unit vector.

## 4.2 Embed the Foundation Models

We embed the semantic embeddings from foundation models to our scene representation. Training the semantic embedding has three aspects. First, we use our scene representation to render a predicted feature map $\hat{\mathbf{F}} \in \mathbb{R}^{W \times H \times D}$ where $W$ is the width, $H$ is the height, and $D$ is the dimension of the feature map. Second, we generate a target feature map $\mathbf{F}$ by feeding the view to a FM. Finally we need to ensure that the predicted feature map is aligned with the corresponding target pixels and follows the same object boundaries in terms of feature similarity.
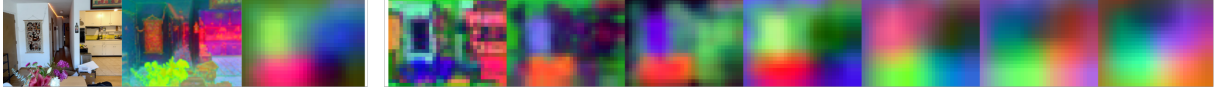
### *Hybrid CLIP Feature for Supervision*

To supervise our feature field outputs, given a calibrated input image, we first rasterize the features into a 2D feature map $\hat{\mathbf{F}}$ where the $(i, j)$th feature is acquired by point-based $\alpha-$blending:

$$\hat{\mathbf{f}}_{i,j} = \sum_{k \in \mathcal{N}} \hat{\mathbf{f}}_k \alpha_k \prod_{l=1}^{i-1}(1 - \alpha_l) \qquad (5)$$
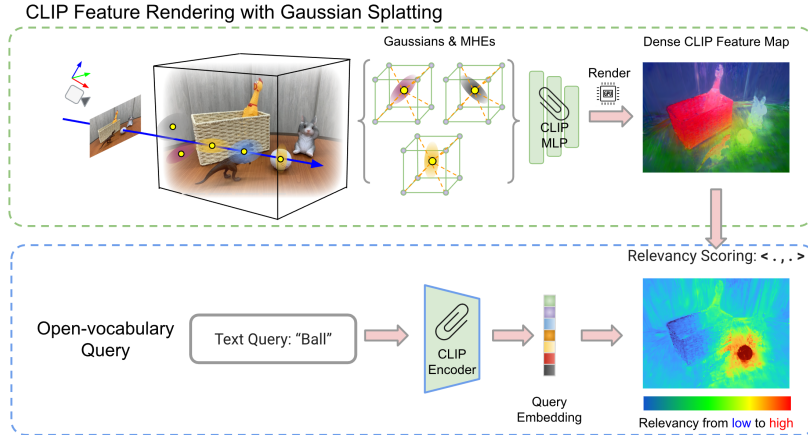
To generate our target CLIP feature map, denoted as $\mathbf{F}$, we initially pre-compute a multi-scale feature pyramid of CLIP embeddings, similar to the approach used in LERF [27]. This involves feeding image patches at various sizes into the CLIP foundation model. However, in contrast to LERF, which trains its scene representation by interpolating embeddings from the pre-computed CLIP feature pyramid at random scales, we rely on a single hybrid CLIP feature map for training our scene representation. We scale up the embeddings of the smaller scales in the pre-computed CLIP feature pyramid bilinearly to the largest scale feature map, and generate the hybrid feature map by averaging them. We define our CLIP loss by the following Huber loss:

$$\mathcal{L}_{CLIP} = \begin{cases} 0.5|\hat{\mathbf{F}} - \mathbf{F}|^2, & \text{if } |\hat{\mathbf{F}} - \mathbf{F}| < \delta \\ \delta \cdot (|\hat{\mathbf{F}} - \mathbf{F}| - 0.5 \cdot \delta), & \text{otherwise} \end{cases}$$
$$(6)$$

where $\delta$ is a hyperparameter, which is set to be 1.25 empirically. As seen in Figure 2 where we use PCA to visualize feature maps following FFD [29], we notice that the target CLIP feature map is not fine-grained enough when embedding similarities of neighboring pixels are considered. This results in poor pixel-alignment gradient signals on Gaussians that are not relevant semantically. On the other hand, DINO [7] features give sharp boundaries between objects [2] in terms of embedding similarity, which can be used for additional regularization.

**Fig. 2 The features extracted from foundation models.** The left three subfigures include the RGB image, extracted DINO features from the foundation model, and the hybrid CLIP feature, which is an average of multi-scale CLIP feature maps shown on the right. On the right, the shown seven CLIP feature maps are the extracted from an image pyramid at multiple scales using the foundation model. The resolution of CLIP features decreases from left to right.



**Fig. 3 FMGS Query pipeline:** *Top*: Given a query view to localize a query, FMGS first renders the dense CLIP feature map. *Bottom*: given an open-vocabulary query, FMGS generates a relevancy map highlighting the relevant part of the rendered CLIP feature map to the query embedding. The highest relevant is colored as red while the lowest relevant part is colored as blue. Note, for visualization simplicity, we show a single-level MHE in this figure while used multiple in implementations.

### Regularization with DINO Feature

To transfer the characteristic of DINO features while maintaining the CLIP embedding semantics, we (a) add a DINO feature field loss and (b) define a pixel-alignment loss between the DINO and CLIP feature fields. The DINO feature field shares the same hash grid parameters as CLIP and gives the same encoding $\mathbf{q}$ for a given $\mathbf{x}$. Then the DINO feature field outputs $\hat{\mathbf{d}} = MLP_{\psi}^{DINO}(\mathbf{q})$ where $\psi$ denotes the parameters of the MLP that are not shared with $MLP_{\phi}^{CLIP}$. This feature field is supervised by passing the *sampled image* once to the pre-trained DINO model without scaling, yielding $\mathbf{D} \in \mathbb{R}^{W \times H \times L}$ where $L$ is the DINO feature dimension. We then render $\hat{\mathbf{D}}$ using the same approach as rendering $\hat{\mathbf{F}}$. The DINO regularization loss is as follows:

$$\mathcal{L}_{DINO} = |\hat{\mathbf{D}} - \mathbf{D}|^2 \qquad (7)$$

### Pixel-alignment with Dot Product Similarity

We define a pixel-alignment loss by defining a kernel around every pixel and enforce the dot product similarity in normalized embedding spaces (between DINO and CLIP) are consistent across the center pixel and surrounding ones. We normalize both rendered features to unit norm, and then compute the loss:

$$\mathcal{L}_{pixel} = \frac{1}{K^2 - 1} \sum_{i \in \mathcal{P}} \sum_{\substack{j \in \mathcal{N}(i), \\ j \neq i}} |\hat{\mathbf{d}}_i^T \hat{\mathbf{d}}_j - \hat{\mathbf{f}}_i^T \hat{\mathbf{f}}_j| \quad (8)$$

where $\mathcal{P}$ denotes the set of all the pixels in the image, and $\mathcal{N}(i)$ is the $K \times K$ patch kernel around the rendered feature at pixel $i$. This makes the rendered CLIP feature follow the same similarity pattern as the DINO feature. Note that we stop the gradient back-propagation through the rendered DINO features in this training loss, which means $MLP_{\psi}^{DINO}$ would not be affected by this loss. $\mathcal{L}_{pixel}$ is also termed as "dotsim" loss for the

rest of the paper, since it is formulated by dot product similarity.

### Training Loss

Overall our total loss is

$$\mathcal{L}_{total} = \lambda \mathcal{L}_{CLIP} + (1 - \lambda)\mathcal{L}_{DINO} + \gamma \mathcal{L}_{pixel} \quad (9)$$

We take the mean reduction over all the pixels in the image plane when computing different loss terms. We also empirically find out adding $\mathcal{L}_{pixel}$ in later iterations during training produces the best results. In Figure 4, we provide examples of features extracted from foundation models for training and the rendered features generated by our trained hybrid semantic scene representation. It is evident that the rendered feature maps exhibit higher quality when compared to the raw feature maps obtained directly from the foundation models, owing to our training process enforces multiple-view consistency.

## 4.3 Relevancy Score

At query time, when provided with a query prompt and a viewing direction, FMGS generates a relevancy map that assigns high scores to semantically relevant locations (see Figure 3). To obtain this relevancy map, we first render the feature map $\hat{\mathbf{F}}$ using our learned semantic feature field via GS rasterization. Then, we calculate the CLIP embedding $\mathbf{f}_{query}$ corresponding to the query prompt.

To obtain the dense relevancy map, we define a set of canonical phrases with CLIP embeddings $\mathcal{F}_{can}$ following the methodology similar to [27]. Then, we compute pairwise softmax scores based on the cosine similarity between the prompt embedding and $\hat{\mathbf{f}}_{i,j}$, representing the $\hat{\mathbf{F}}$ at location $(i, j)$, as well as the canonical embeddings for canonical phrases. We take the minimum value of the softmax over all canonical prompts and deem it the relevancy score $r$:

$$r_{i,j} = \min_{n} \frac{\exp(\hat{\mathbf{f}}_{i,j}^T \mathbf{f}_{query})}{\exp(\hat{\mathbf{f}}_{i,j}^T \mathbf{f}_{query}) + \exp(\hat{\mathbf{f}}_{i,j}^T \mathbf{f}_{can}^n)}, \mathbf{f}_{can}^n \in \mathcal{F}_{can} \quad (10)$$

With the above definition, the relevancy score is higher when a query embedding is closer to the rendered feature than the canonical features. We follow [27] and choose the following canonical prompts: "object", "stuff", "things", and "texture". We also find that these work well for a wide range of queries removing the need for tuning these canonical terms. In Figure 4, we present representative relevancy maps generated by matching the query embedding with our rendered CLIP feature map and the target CLIP feature map from the foundation model used in our training. It's evident that the relevancy map derived from our rendered CLIP feature map overall exhibits finer granularity and higher quality.

## 4.4 Implementation Details

Our approach employs a hash grid for representing language features, which is notably larger than a typical RGB hash grid. This hash grid comprises 24 layers, spanning resolutions from 16 to 512, and possesses a hash table size of $2^{20}$ with an associated feature dimension of 8. The architecture of the CLIP and DINO MLP models used for $MLP_{\phi}^{CLIP}$ and $MLP_{\psi}^{DINO}$ aligns with that of LERF [26]. Furthermore, we leverage the Open-CLIP [13] ViT-B/16 model, which has undergone training on the LAION-2B dataset. Notably, this model operates with an image pyramid that varies in scale from 0.05 to 0.5 of image size, encompassing a total of seven scales for pre-computing CLIP feature pyramid. The pre-computed feature pyramid is subsequently processed by average pooling to generate the final hybrid CLIP feature for training our semantics embedded field.
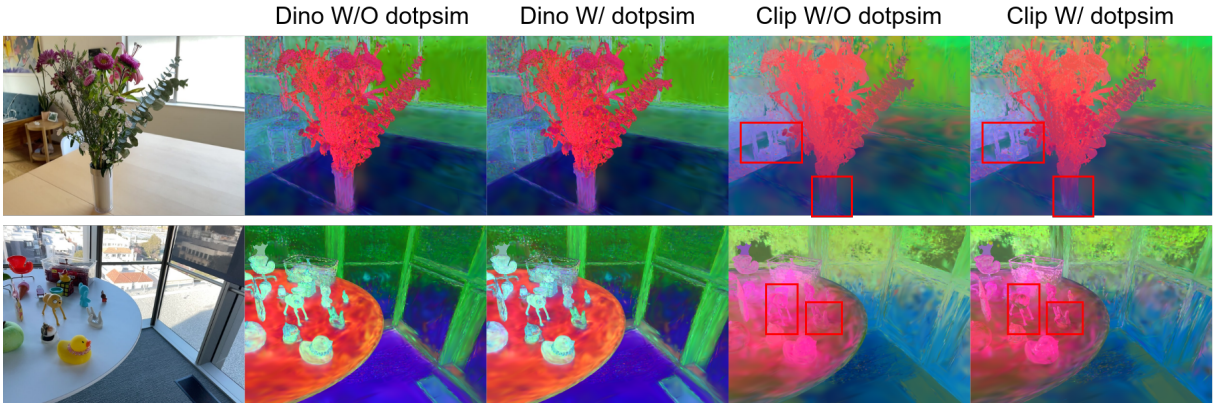
Initially, we train the Vanilla Gaussian Splatting scene representation [26] through a total number of 30K iterations, with approximately 10 minutes total time for a room-scale scene. It's worth noting that representing such a scene requires the utilization of millions of Gaussians. Subsequently, we maintain the frozen states of the geometric attributes and spherical harmonics associated with these Gaussians throughout the subsequent training process for semantic embedding fields.

To mitigate GPU memory constraints, we strategically select approximately 40% of the Gaussians based on criteria such as high opacity values and a 2D radius of projected Gaussian exceeding 2 pixels in at least one training view. Only these selected Gaussians are involved in the rendering process when we train the semantic embeddings. For optimization, we employ the

**Fig. 4 Features for Training and Rendered Views. Left**: From left to right, the figures show the RGB image, the rendered DINO feature map, the raw DINO feature map extracted for training, the rendered CLIP feature map, and the raw CLIP feature map used for training. **Right**: We display the relevancy scores for the rendered and raw CLIP feature maps with the text query 'flower', where the color bar indicates relevancy scores normalized within the 0-255 range. Notably, querying the raw CLIP feature map is much inferior to querying the rendered CLIP feature map.



**Fig. 5 Effect of dot product similarity (dotpsim) loss.** From left to right: RGB image, rendered DINO feature without dotpsim, rendered DINO feature with dotpsim, rendered CLIP without dotpsim, and rendered CLIP feature map with dotpsim. The DINO feature maps do not have significant differences with or without dotpsim. From the CLIP feature maps, we can see that objects can be further distinguished from each other and the background. Differences are highlighted in the red boxes.

RAdam optimizer with a weight decay of $10^{-9}$. We incorporate an exponential learning rate scheduler, which spans from an initial value of $5 \times 10^{-3}$ and gradually decreases to $4 \times 10^{-3}$ over the course of 4.2K training steps (after the initial 30K original GS training steps). In our training regimen, all models initially undergo 2.5K steps without the pixel alignment loss being enabled. These training and testing procedures are executed on an NVIDIA RTX A5000 GPU with 24GB of GPU RAM. The semantic feature field training time with a total of 4.2K steps takes about 1.4 hours. During training, we use weighting factors to balance the CLIP loss ($\lambda = 0.2$) and the pixel-alignment loss ($\gamma = 0.01$).
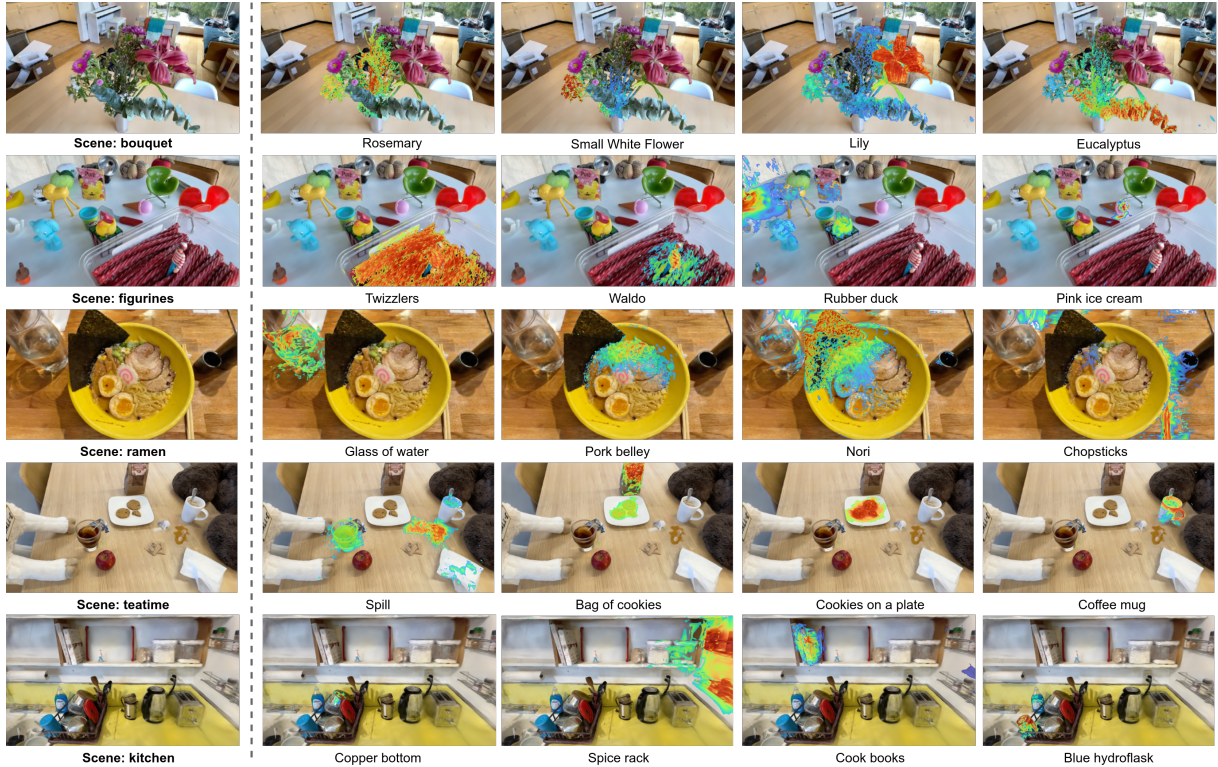
## 5 Experiments

Our hybrid semantic scene representation, FMGS, seamlessly integrates the 3D Gaussians and multi-resolution hashing encoding and supports both photo-realistic rendering and open-vocabulary object detection. In this section, we carefully evaluate the performance of open-vocabulary object detection (or localization) of our proposed method in uncontrolled real-world scenarios. To showcase the embedding quality of our method, we also evaluate it out-of-the-box on the open-vocabulary semantic segmentation task. We compare our method to other SOTA approaches for each experiment and show significant improvement over their results.

### 5.1 Object Detection in the Wild

By distilling the language embeddings extracted from off-the-shelf vision-language model, CLIP, our FMGS is applicable for associating a wide range of textual prompts with the relevant vision clues. We test the open-vocabulary object understanding capability of our method by object detection experiments.

**Dataset:** We use the same dataset as used in the LERF [27] for object detection evaluation,

**Fig. 6 Relevancy score for object detection.** *Left*: The rendered RGB image at novel view from 5 scenes on LERF dataset [27]. *Right*: Visualization of relevancy scores with the given text queries shown below the figures. We overlay them on the RGB images.

| Scene | FFD-LSeg [29] | OWL-ViT [39] | LERF [27] | Ours |
|---|---|---|---|---|
| bouquet | 50.0% | 66.7% | 83.3% | 100.0 % |
| figurines | 8.9% | 38.5% | 87.2% | 89.7% |
| ramen | 15.0% | 92.5% | 62.5% | 90.0 % |
| teatime | 28.1% | 75.0% | 96.9% | 93.8% |
| kitchen | 13.0% | 42.6% | 85.2% | 92.6 % |
| Average Acc. | 18.0% | 54.8% | 83.0% | 93.2% |
| Inference FPS | - | - | 0.1214 | 103.4 |

**Table 1 Accuracy and runtime efficiency (Frames Per Second, FPS) of object detection with open-vocabulary queries.** comparison between Feature Fields Distillation [29] using LSeg [30] features (FFD-Lseg), OWL-ViT [39], LERF [27] and Ours FMGS. We highlight the best , second-best accuracy scores. Please find more details on scenes and text queries for LERF dataset in [27].
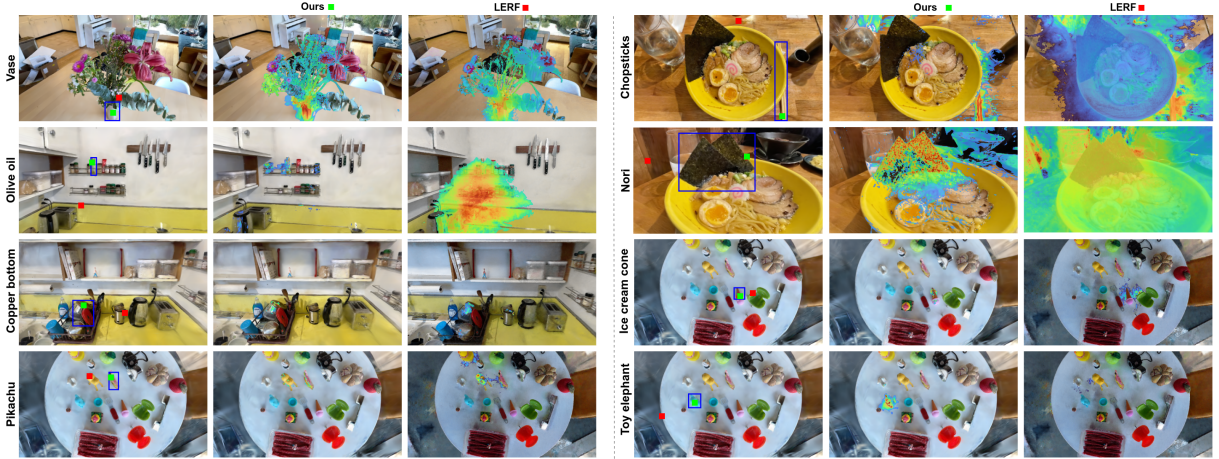
for the purpose of fair comparison. It consists of five labelled scenes with 2D bounding boxes of objects associated with text prompts. There are objects including both common and long-tail ones with different sizes, and the queries for objects are quite diverse, like 'vase', 'eucalyptus', 'big white crinkly flower', 'pikachu', 'twizzlers', 'spoon handle', 'power outlet', 'waldo', 'stuffed bear', 'cookies on a plate', etc. The location of queried images are labelled by bounding boxes in the test images, which are rendered at novel views from trained

NeRF models of individual scenes. The scenes in LERF dataset are collected by an iPhone, and each scene comprise $\sim 200$ images. The provided poses of images from Ploycam app are with significant noises in some scenes. Thus we regenerate the poses of images by running COLMAP [46], which also yields sparse 3D visual points serving as input to initialize 3D Gaussians in our method. The poses of the officially-provided test images are also properly transferred to our COLMAP trajectory by Sim(3) alignment between officially-provided image poses and our COLMAP poses.

**Evaluation Protocol:** Following LERF [26], the evaluation metric for object detection is the accuracy rate. We redeem the query is a success if the highest relevancy pixel locates inside the target box. The relevancy score at each pixel is obtained by matching the rendered CLIP feature map with the language embedding of given text query as described in Sec. 4.3.

**Baselines:** We compare against FFD-LSeg that embeds pixel-aligned LSeg feature [30] into NeRF (NeuralStudio 'neurfacto' implementation

11

**Fig. 7 Object detection results.** The left and right groups of figures illustrate the detection results for four open-vocabulary queries, respectively. Each group comprises three subfigures: *Left* displays the ground-truth bounding boxes (blue), our detected highest-relevancy pixel (green) and the one detected by LERF (red) [27]. *Middle* showcases our relevancy score corresponding to the given text query. The text query is shown at the far left of each row. *Right* showcases LERF's relevancy score corresponding to the given text query. Our computed relevancy score is more focused on the target objects linked to the query.

| Methods | bed | | sofa | | lawn | | room | | bench | | table | |
| | mIoU | mAP | mIoU | mAP | mIoU | mAP | mIoU | mAP | mIoU | mAP | mIoU | mAP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| OV-Seg [31] | 79.8 | 40.4 | 66.1 | 69.6 | 81.2 | 92.1 | 71.4 | 49.1 | 88.9 | 89.2 | 80.6 | 65.3 |
| 3D-OVS [33] | 89.5 | 96.7 | 74.0 | 91.6 | 88.2 | 97.3 | 92.8 | 98.9 | 89.3 | 96.3 | 88.8 | 96.5 |
| **LERF** [27] | 33.5 | 25.6 | 28.1 | 45.6 | 49.8 | 82.0 | 26.3 | 49.1 | 55.2 | 79.5 | 31.1 | 33.3 |
| **Ours** | 38.0 | 50.1 | 56.6 | 82.0 | 64.9 | 90.5 | 57.0 | 85.3 | 62.1 | 84.1 | 63.6 | 85.3 |
| **Ours-Refined** | 80.6 | 85.5 | 90.8 | 97.4 | 92.6 | 98.5 | 87.9 | 97.8 | 84.5 | 94.8 | 89.4 | 97.2 |

**Table 2 Segmentation Evaluation.** We report the mIoU(↑) scores and the mAP(↑) scores of the following methods in 6 scenes of 3D-OVS dataset [33]. Note that 3D-OVS is a weakly supervised method, which knows the segmentation annotations in training and specially designed for segmentation task. Our method and LERF are 3D method training without any segmentation annotations, relying only on the relevancy between class query and the rendered CLIP features. OV-Seg [31] is a supervised method for segmentation task. Our method and LERF are unsupervised methods, under apple-to-apple comparison. We can further post-process and refine our 2D segmentation results by SAM [28] and get the results shown as 'Ours-Refined'.

by feature fields distillation method [29], OWL-ViT [39] that is a 2D method based on Vision Transformer encoder and fine-tuned for object detection, LERF [27] that embeds CLIP and DINO features into NeRF. The 3D methods, FFD-LSeg and LERF, share the same evaluation protocol as our FMGS. For the 2D method, OWL-ViT, we regard it as a success if the center of the predicted bounding box locates in the target box.

**Evaluation Results:** The quantitative evaluation results on all sequences of LERF dataset are presented in Table 1, and representative relevancy score maps of the proposed method are shown in

| Methods | bouquet | figurines | ramen | teatime | kitchen | **Average** |
|---|---|---|---|---|---|---|
| **Ours** | 100.0 | 89.7 | 90.0 | 93.8 | 92.6 | **93.2** |
| **W/O dotpsim** | 100.0 | 91.0 | 85.0 | 90.6 | 85.2 | 90.4 |
| **W/O hybrid CLIP** | 54.2 | 32.1 | 52.5 | 6.3 | 9.3 | 30.8 |
| **W/ LERF CLIP** | 91.7 | 70.5 | 72.5 | 72.5 | 87.0 | 78.8 |
| **W/O MHE** | 91.7 | 71.8 | 90.0 | 90.6 | 77.8 | 84.4 |

**Table 3 Ablation study.** Object detection comparison between our full method, ours without dot product similarity (dotpsim) loss, and ours without hybrid CLIP features by averaging at multiple scales for supervision, using single scale CLIP feature at the finest-resolution instead, as well as ours with LERF CLIP at multiple individual scales [27].

Figure 6. The detailed results demonstrate significant advantages of FMGS's integration of language embeddings in detecting objects associated with long-tail prompts. While LSeg [30], trained

on a small dataset to learn pixel-aligned CLIP features, exhibits diminished open-vocabulary language understanding capabilities, the approach of FFD-LSeg, which distills LSeg features into radiance fields, struggles with comprehending long-tail queries and consequently exhibits poorer performance. In terms of open-vocabulary 2D detection, Owl-ViT, which utilizes full-HD NeRF views and selects bounding boxes based on the highest confidence scores for text queries, outperforms FFD-Lseg. However, when faced with long-tail queries, Owl-ViT's performance falls short in comparison to the robust and versatile FMGS.

We also conducted a comparison with the closest method, LERF, which distills DINO and CLIP features into neural radiance fields represented solely by MHEs. As depicted in Table 1, our FMGS outperforms LERF significantly, achieving an accuracy improvement of 10.2 percentage points. Note that our tested LERF results, obtained using the officially released code, slightly surpasses those reported in the original paper [27].

In Figure 7, we present side-by-side comparisons with LERF [27]. The object detection results are visualized, highlighting the superior quality of the relevance map produced by our FMGS. It notably focuses more on the queried target objects, as opposed to LERF. This outcome stems from our hybrid representation, which combines 3D Gaussians and MHEs for semantic scene representation. The 3D Gaussians represent both the geometry and appearance of the scene, naturally dividing 3D structures of objects and the scene into distinct Gaussian volumes. This partitioning feature aids in distinguishing objects from each other and from the background. In FMGS, we assign an identical MHE embedding to a Gaussian volume, further promoting semantic consistency in local proximity. This, in turn, contributes to the focusing of relevance on the target object. Taking the query 'Pikachu' in Figure 7 as an example, where 'Pikachu' is depicted on the side of a paper bag. Even when observing from a challenging viewpoint with almost no visibility of 'Pikachu', FMGS successfully maintains high relevance at the target location, due to its 3D consistency and fine-grained scene understanding. In contrast, LERF fails to detect 'Pikachu' and mistakenly identifies a visually similar object.
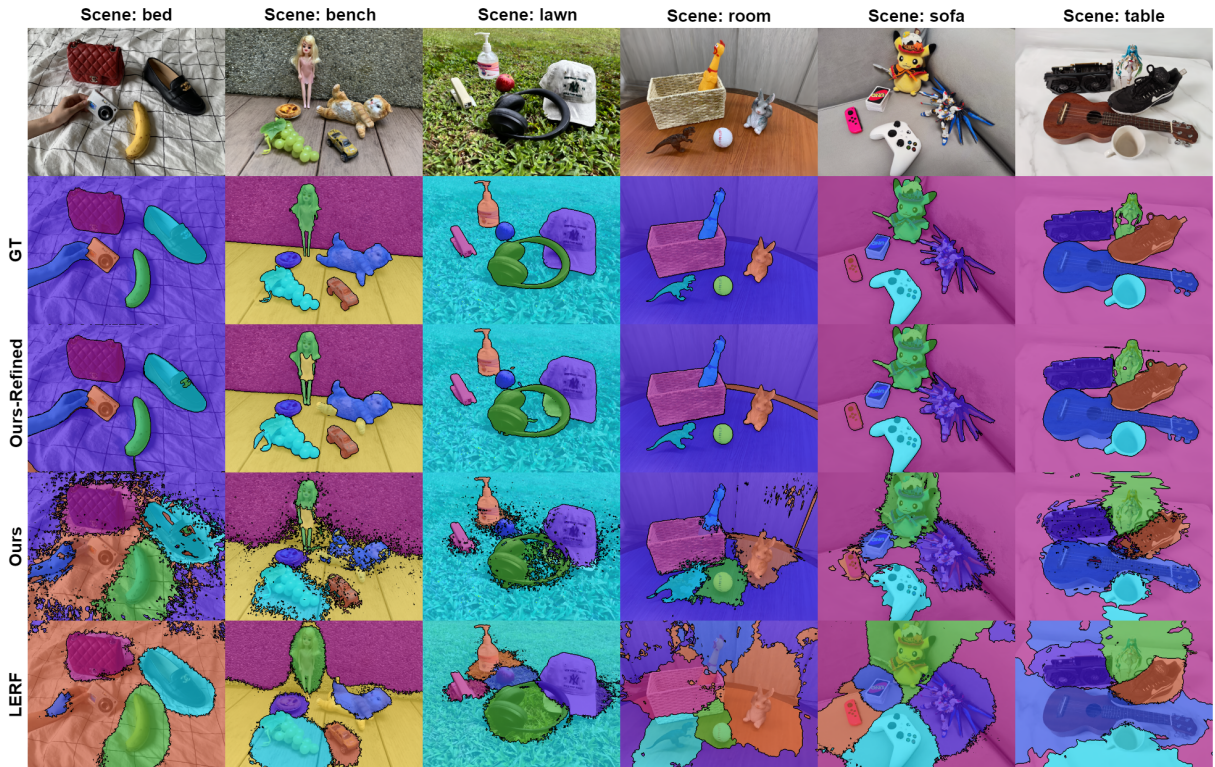
**Inference Runtime:** Our FMGS, relying on 3D Gaussian Splatting rendering [26], excels in efficiently rendering RGB images. We've implemented our rendering method for CLIP and DINO feature maps based on a CUDA implementation of Gaussian Splatting rendering. Even when rendering deep features with high dimensions, which can significantly increase computation time, our FMGS remains remarkably fast. It can render the $480 \times 270$ CLIP feature map, DINO feature map, and RGB image jointly at an impressively high rate of 103.4 FPS during inference, even with our unoptimized implementation. In contrast, LERF operates at a significantly slower pace, achieving a mere 0.1214 FPS during inference (see Table 1). This slowness stems from LERF's need to perform a brute-force search for the best scales when rendering CLIP features, spanning a range from 0 to 2 meters with 30 increments. Consequently, we are **851.73 times faster** than LERF in rendering CLIP features, enabling efficient real-time open-vocabulary queries after our scene representation is trained.

## 5.2 Unsupervised Segmentation

In following experiments we use FMGS to segment queries and evaluate their segmentation masks. **Note** that our method is not delicately designed for the segmentation task. We lack a dedicated segmentation header for predicting segmentation masks, nor do we explicitly partition the scene at the object level. We have examined the open-vocabulary language understanding capability of FMGS in the object detection experiments discussed in the above section. Our primary objective for doing this segmentation evaluation is to assess the pixel-level accuracy of the rendered CLIP features obtained from the trained scene representation. Segmentation relies on matching these rendered CLIP features to the embeddings of the provided semantic labels.

**Dataset:** We conducted our segmentation evaluation using the 3D-OVS dataset [33], which consists of six scenes with labeled ground-truth semantic segmentation masks for test image views. These scenes are characterized by their cleanliness, with clear backgrounds and well-defined foreground objects. Each scene comprises approximately 30 images with predefined poses and sparse points computed using COLMAP [46]. The

13

**Fig. 8 Semantic segmentation results.** In the rows from top to bottom, we display RGB images, ground-truth (GT) segmentation masks, our refined segmentation results, our segmentation results, and the segmentation results obtained by LERF [27] scene representation. It's essential to note that neither our method nor LERF was initially intended for the segmentation task. Our primary aim is to evaluate the pixel accuracy of the relevance map computed from the rendered CLIP features. We can further post-process and refine our 2D segmentation results by SAM [28] and get the results shown as 'Ours-Refined'.

dataset includes a variety of objects, including many long-tail objects like 'Gundam,' 'Pikachu,' 'stapler', and more. For further details about the scenes and semantic labels, please refer to [33].

**Evaluation Protocol:** In terms of our evaluation protocol, we rely on the annotated ground-truth masks for the test views. These masks serve as a reliable benchmark for both qualitative and quantitative assessments of segmentation performance. We calculate the mean Intersection over Union (mIOU) scores and mean Average Precision (AP) metrics by comparing the segmentation results with these ground-truth masks.

**Baselines:** We conduct a direct comparison of our method with LERF [27]. To perform semantic segmentation, we initially obtain relevancy scores by computing the cosine similarity between the rendered CLIP feature and the embeddings of all class labels (this is different from the relevancy score calculation with auxiliary canonical phrases involved in Sec. 4.3.). These relevancy scores serve

as segmentation logits, and we subsequently apply the softmax function to convert them into probabilities. Each pixel is then assigned a semantic class label corresponding to the maximum probability. Note that LERF [27] requires a scale factor when rendering CLIP features, and we report the best segmentation results that can be achieved by LERF by selecting the best scales for each ray. It's also important to note that both LERF and our method encounter challenges in discerning the semantic labels of backgrounds when presented with visibility-limited close views and lack of context. Therefore, we have replaced the original background labels, including 'white sheet', 'wood wall', 'grey sofa', and 'lime wall', with a more general label 'background' when testing LERF and our method.

Additionally, for comprehensive reference, we present results obtained using the dedicated 3D-OVS method [33] for the segmentation task. However, it is worth emphasizing that comparing

14

object detection methods like ours and LERF [27] to 3D-OVS is not entirely equitable, as acknowledged in the paper of 3D-OVS [33]. 3D-OVS [33] has prior access to segmentation class labels and distill class-related information into the radiance field during training. In contrast, neither LERF nor our methods have access to class labels during scene representation training. Consequently, the trained 3D-OVS scene representation can only be effectively employed for querying the classes known before training, and does not support arbitrary semantic queries beyond the trained classes. Furthermore, we compare to a 2D ceiling approach [31], OV-Seg, which is directly trained for open-vocabulary semantic segmentation by fine-tuning CLIP on masked image regions and text descriptions. OV-Seg is supervised with mask-category pairs, while ours and LERF are completely unsupervised.

**Evaluation Results** The segmentation experiment results are presented in Table 2 and Figure 8. Notably, our approach outperforms LERF [27] by a significant margin across all cases. This superior performance can be attributed to the higher quality of our rendered CLIP feature compared to the one produced by LERF. Our method exhibits more concentrated high relevancy around the queried objects, showcasing the advantage of our semantic scene representation, which maintains high semantic consistency in local proximity. We refine our 2D segmentation results using SAM [28] by assigning class labels to SAM segments through a majority voting based on pixel labels obtained from our method. The refined result is displayed in the 'Ours-Refined' row of Table 2 and Figure 8.

## 5.3 Ablations

We conducted an ablation study on the object detection task, as it serves as a key indicator of our method's open-vocabulary semantic understanding capabilities. The results are presented in Table 3.

### 5.3.1 Hybrid CLIP feature

In this ablation study, we investigated using a single scale of CLIP features at the finest scale level, rather than our hybrid CLIP features, which are obtained by averaging multiple-scale

CLIP features extracted from patches at different resolutions. As demonstrated in Table 3, the hybrid CLIP feature for supervision is greatly important. The scene understanding capability is severely compromised when employing only a single-scale CLIP feature for supervision (denoted as 'W/O Hybrid CLIP'). The inferior performance observed with the use of a single-scale CLIP feature stems from its potential inadequacy in capturing sufficient contextual information within the image. By contrast, our hybrid CLIP feature encompasses a significantly larger receptive field, enhancing its contextual awareness.

To conduct a comprehensive analysis, we also compare our method of using a single hybrid CLIP supervision ('Ours') to the approach of utilizing CLIP supervision at multiple individual scales, as proposed by LERF [27], which we denote as 'W/ LERF CLIP'. During its training, randomly sampled scale factors are concatenated with the intermediate MHE feature vector $\mathbf{q}$ (refer to Sec. 4.1) to decode CLIP features at each scale. At inference time, the optimal scale factors at each pixel location are determined by an exhaustive search among 30 uniformly distributed scales. The best scale factor is selected based on its relevance score to the query. Consequently, this approach is at least 30 times slower than our method with a single hybrid CLIP supervision during inference. Despite its inefficiency, the accuracy achieved by ' W/ LERF CLIP' (78.8%) significantly lags behind that of 'Ours' (93.2%), as shown in Table 3, underscoring the efficacy and superior efficiency of our proposed method. It is worth noting that while it is plausible that 'W/ LERF CLIP' could achieve better performance with much more training iterations, this aspect falls beyond the scope of our current investigation.

### 5.3.2 Pixel-alignment loss

To assess the effectiveness of our proposed pixel alignment loss, we conducted an ablation study by training our semantic scene representation without this loss. The impact of omitting the pixel alignment loss on the accuracy of the object detection task is shown in Table 3. Furthermore, we provide qualitative results in Figure 5, which indicates that CLIP features from a scene representation trained with pixel-alignment loss are

better at distinguishing between different objects and separating objects from the background.

### 5.3.3 Scene Representation

To evaluate the effectiveness of our hybrid scene presentation, which integrates 3D Gaussians for geometry and appearance representation alongside MHE for efficient semantic embedding, we compare it with a vanilla scene representation method that solely employs Gaussians without MHE by attaching semantic embeddings to each Gaussian ('*W/O MHE*'). Notably, for fair comparisons, we consider only the selected Gaussians with sufficient opacity and 2D radius in the vanilla method (as detailed in Sec. 4.4) identical to our method ('*Ours*'). These attached semantic embeddings share the same dimensionality as the intermediate MHE feature vector $\mathbf{q}$ in our method (see Sec. 4.1). They are also decoded by two MLPs, $MLP_\phi^{CLIP}$ and $MLP_\psi^{DINO}$, to obtain CLIP and DINO features while rendering the feature maps.

As indicated in Table 3, '*W/O MHE*' achieves an object detection accuracy of 84.4%, which is inferior to '*Ours*' with its hybrid scene representation. The superior performance of '*Ours*' can be attributed to its ability to render CLIP feature maps at a higher quality over '*W/O MHE*'. Additionally, in terms of scene representation complexity, the MHE component of our method maintains a constant number of parameters across the five scenes of the LERF dataset, whereas the parameter count of the semantic embeddings in '*W/O MHE*' varies with the number of Gaussians and increases by 84.1%, 113.5%, 2.14%, 174.2%, and 136.9% compared to '*Ours*'. This substantial memory demand of '*W/O MHE*' can pose even greater challenges in large-scale scenarios.

## 6 Discussion and Limitations

When comparing FMGS to LERF[27], both methods distill Clip and Dino features from foundation models into 3D scene representations. However, their rendering algorithms and scene representations differ significantly. These distinctions lead to rapid and high-quality language feature acquisition using common hyperparameters, such as the feature field architecture. An additional key advantage of FMGS is that it employs the same feature embedding for each Gaussian, regardless of the viewing direction. This feature enables direct 3D localization of vision-language queries. It's important to note that FMGS not only facilitates the localization of language queries in 3D but also allows for finding a given image of the scene using the 3D Gaussian embeddings. LERF, on the other hand, does not offer such 3D localization capabilities out of the box.

In terms of limitations, FMGS currently relies heavily on the presence of high-quality and calibrated input images, a limitation shared with NeRF-based approaches. Additionally, the performance of FMGS is entirely contingent on the quality of the base foundation models used for training the feature fields. It is conceivable that a model better suited for localizing language within images could yield improved feature field quality. Furthermore, to enhance performance in semantic segmentation tasks, it is advisable to embed a specialized segmentation foundation model, such as SAM [28, 9], into our scene representation. Unlike using SAM solely for post-refinement of our 2D segmentation results, directly distilling SAM results into the 3D space offers the advantage of enforcing multi-view consistency, potentially leading to improved performance.

## 7 Conclusions

Foundation Model Embedded Gaussian Splatting (FMGS) contributes to scene understanding by seamlessly merging vision-language embeddings and 3D representation. This novel 3D scene representation achieves multi-view semantic consistency through self-supervised distillation and pixel alignment of CLIP features. The resulting feature embedded 3D Gaussians achieve state-of-the-art performance in comparison to previous methods. By bridging vision, language, and 3D, FMGS paves the way for unprecedented object comprehension in real-world environments, opening exciting possibilities for augmented reality, robotics, and beyond.

## 8 Acknowledgement

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

[2] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *arXiv preprint arXiv:2112.05814*, 2(3):4, 2021.

[3] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19129–19139, 2022.

[4] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.

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

[6] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023.

[7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.

[8] Paola Cascante-Bonilla, Hui Wu, Letao Wang, Rogerio S Feris, and Vicente Ordonez. Simvqa: Exploring simulated environments for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5056–5066, 2022.

[9] Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Chen Yang, Wei Shen, Lingxi Xie, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs. In *NeurIPS*, 2023.

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

[11] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024.

[12] Haoran Chen, Kenneth Blomqvist, Francesco Milano, and Roland Siegwart. Panoptic vision-language feature fields. *arXiv preprint arXiv:2309.05448*, 2023.

[13] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. *arXiv preprint arXiv:2212.07143*, 2022.

[14] Rodolfo Corona, Shizhan Zhu, Dan Klein, and Trevor Darrell. Voxel-informed language grounding. *arXiv preprint arXiv:2205.09710*, 2022.

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

[16] Ben Fei, Jingyi Xu, Rui Zhang, Qingyuan Zhou, Weidong Yang, and Ying He. 3d gaussian as a new vision era: A survey. *arXiv preprint arXiv:2402.07181*, 2024.

[17] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Open-vocabulary image segmentation. *arXiv preprint arXiv:2112.12143*, 2021.

[18] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4089–4098, 2018.

[19] Margarita Grinvald, Fadri Furrer, Tonci Novkovic, Jen Jen Chung, Cesar Cadena, Roland Siegwart, and Juan Nieto. Volumetric instance-aware semantic mapping and 3d object discovery. *IEEE Robotics and Automation Letters*, 4(3):3037–3044, 2019.

[20] Qiao Gu, Alihusein Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, Chuang Gan, Celso Miguel de Melo, Joshua B. Tenenbaum, Antonio Torralba, Florian Shkurti, and Liam Paull. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. *arXiv*, 2023.

[21] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021.

[22] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568, 2011.

[23] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, Joshua B. Tenenbaum, Celso Miguel de Melo, Madhava Krishna, Liam Paull, Florian Shkurti, and Antonio Torralba. Conceptfusion: Open-set multimodal 3d mapping, 2023.

[24] Peter Karkus, Shaojun Cai, and David Hsu. Differentiable slam-net: Learning particle slam for visual navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2815–2825, 2021.

[25] Nikhil Keetha, Jay Karhade, Krishna Murthy Jatavallabhula, Gengshan Yang, Sebastian Scherer, Deva Ramanan, and Jonathon Luiten. Splatam: Splat, track & map 3d gaussians for dense rgb-d slam. *arXiv preprint arXiv:2312.02126*, 2023.

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

[27] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023.

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

[29] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. In *Advances in Neural Information Processing Systems*, volume 35, 2022.

[30] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022.

[31] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. *arXiv preprint arXiv:2210.04150*, 2022.

[32] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1954–1963, 2021.

[33] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. Weakly supervised 3d open-vocabulary segmentation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[34] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. 3d open-vocabulary segmentation with foundation models. *arXiv preprint*

*arXiv:2305.14093*, 2023.

[35] Shiyang Lu, Haonan Chang, Eric Pu Jing, Abdeslam Boularias, and Kostas Bekris. OVIR-3d: Open-vocabulary 3d instance retrieval without training on 3d data. In *7th Annual Conference on Robot Learning*, 2023.

[36] Yuheng Lu, Chenfeng Xu, Xiaobao Wei, Xiaodong Xie, Masayoshi Tomizuka, Kurt Keutzer, and Shanghang Zhang. Open-vocabulary point-cloud object detection without 3d annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.

[37] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *CVPR*, pages 7086–7096, 2022.

[38] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.

[39] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022.

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

[41] Gaku Narita, Takashi Seno, Tomoya Ishikawa, and Yohsuke Kaji. Panopticfusion: Online volumetric semantic mapping at the level of stuff and things. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4205–4212, 2019.

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

[43] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2024.

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

[45] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.

[46] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.

[47] Thomas Schöps, Torsten Sattler, and Marc Pollefeys. Surfelmeshing: Online surfel-based mesh reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2494–2507, 2019.

[48] Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022.

[49] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. LEGaussians: Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2024.

[50] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15598–15607, 2021.

[51] Ayça Takmaz, Elisabetta Fedele, Robert W. Sumner, Marc Pollefeys, Federico Tombari, and Francis Engelmann. OpenMask3D: Open-Vocabulary 3D Instance Segmentation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[52] Jesse Thomason, Mohit Shridhar, Yonatan Bisk, Chris Paxton, and Luke Zettlemoyer. Language grounding with 3d objects. In *Conference on Robot Learning*, pages 1691–1701,

2022.

[53] Nikolaos Tsagkas, Oisin Mac Aodha, and Chris Xiaoxuan Lu. Vl-fields: Towards language-grounded neural implicit spatial representations. *arXiv preprint arXiv:2305.12427*, 2023.

[54] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In *3DV*, 2022.

[55] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11686–11695, 2022.

[56] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. *CoRR*, abs/2201.08845, 2022.

[57] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. Ulip: Learning unified representation of language, image and point cloud for 3d understanding. *arXiv preprint arXiv:2212.05171*, 2022.

[58] Le Xue, Ning Yu, Shu Zhang, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. Ulip-2: Towards scalable multimodal pre-training for 3d understanding, 2023.

[59] Xingbin Yang, Liyang Zhou, Hanqing Jiang, Zhongliang Tang, Yuanbo Wang, Hujun Bao, and Guofeng Zhang. Mobile3drecon: real-time monocular 3d reconstruction on a mobile phone. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3446–3456, 2020.

[60] Jianglong Ye, Naiyan Wang, and Xiaolong Wang. Featurenerf: Learning generalizable nerfs by distilling pre-trained vision foundation models. *arXiv preprint arXiv:2303.12786*, 2023.

[61] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019.

[62] Yanjie Ze, Ge Yan, Yueh-Hua Wu, Annabella Macaluso, Yuying Ge, Jianglong Ye, Nicklas Hansen, Li Erran Li, and Xiaolong Wang. Multi-task real robot learning with generalizable neural feature fields. *CoRL*, 2023.

[63] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. *arXiv preprint arXiv:2112.02413*, 2021.

[64] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya You, Zhangyang Wang, and Achuta Kadambi. Feature 3DGS: Supercharging 3d gaussian splatting to enable distilled feature fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2024.

[65] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, pages 350–368. Springer, 2022.

[66] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization*, 2001.