# FMLGS: Fast Multilevel Language Embedded Gaussians for Part-level Interactive Agents

Xin Tan, Yuzhou Ji, He Zhu, and Yuan Xie

**Abstract**—The semantically interactive radiance field has long been a promising backbone for 3D real-world applications, such as embodied AI to achieve scene understanding and manipulation. However, multi-granularity interaction remains a challenging task due to the ambiguity of language and degraded quality when it comes to queries upon object components. In this work, we present FMLGS, an approach that supports part-level open-vocabulary query within 3D Gaussian Splatting (3DGS). We propose an efficient pipeline for building and querying consistent object- and part-level semantics based on Segment Anything Model 2 (SAM2). We designed a semantic deviation strategy to solve the problem of language ambiguity among object parts, which interpolates the semantic features of fine-grained targets for enriched information. Once trained, we can query both objects and their describable parts using natural language. Comparisons with other state-of-the-art methods prove that our method can not only better locate specified part-level targets, but also achieve first-place performance concerning both **speed** and **accuracy**, where FMLGS is 98 × faster than LERF, 4 × faster than LangSplat and 2.5 × faster than LEGaussians. Meanwhile, we further integrate FMLGS as a virtual agent that can interactively navigate through 3D scenes, locate targets, and respond to user demands through a chat interface, which demonstrates the potential of our work to be further expanded and applied in the future.

**Index Terms**—3D vision, language embedding, 3D Gaussian splatting.

## 1 INTRODUCTION

A S a specially reconstructed 3D scene, language-embedded radiance fields can be semantically interactive under continuous new visual angles. Therefore, it serves as a promising component within fields such as embodied AI and augmented reality. Given a set of posed images, a language-embedded radiance field learns an accurate and

Xin Tan, Yuzhou Ji, He Zhu, and Yuan Xie are with the School of Computer Science and Technology, East China Normal University, Shanghai 200062, China (e-mails: xtan@cs.ecnu.edu.cn; 10215102492@stu.ecnu.edu.cn; 10215102469@stu.ecnu.edu.cn; xieyuan8589@foxmail.com).
Yuan Xie is also with the Shanghai Innovation Institute, Shanghai 200062, China.



Fig. 1: Results of querying "A Button of Xbox Wireless Controller". While LangSplat [1] fails in detailed part-level localization, our method provides an accurate outcome.

efficient 3D representation of semantics, enabling scene understanding and manipulation by providing consistent semantics across any views. Although much progress has been made in object-level language-embedded radiance fields, it still remains a vital challenge to achieve high-quality part-level semantics.

The existing works mainly lie on pixel-based and mask-based methods. For pixel-based methods, they are good at detecting targets of either whole objects or small parts of interest, but cannot provide clear target outline, limiting further applications concerning demanding 3D manipulation tasks. For example, LERF [2] firstly enables pixel-aligned queries of the distilled 3D CLIP [3] embeddings, supporting long-tail open-vocabulary queries within NeRF [4]. Although LERF's heat map shows astonishing localization ability upon object parts such as "bear nose" and "engine", the results are too fuzzy to be used for downstream applications. Based on LERF's core concept, LEGaussians [5] proposes a feature lifting strategy within 3DGS [6], achieving faster queries and smoother relevancy maps, but the results are still unstable and inconsistent.

In contrast, the mask-based methods perform well in generating accurate query masks while suffering from limitedly recallable targets. For example, Gaussian Grouping [7] introduces identity encoding and 3D spatial consistency regularization, consistently lifting 2D Segment Anything Model (SAM) [8] masks into 3D scenes and provides high-quality object-level results. In FastLGS [9], we also propose cross-view grid mapping to achieve accurate results and fast queries. While these methods are based solely on *everything* mode SAM masks to ensure consistent object identity, they lose the ability to query any object parts. LangSplat [1],

on the other hand, learns hierarchical semantics by sorting and filtering SAM masks into *whole*, *part* and *subpart* groups for lifting multi-scale language embeddings. This technique enables LangSplat to further provide query ability upon certain components such as "yellow" parts of a Pikachu toy while still having high accuracy in object-level results. However, LangSplat only processes masks generated from the same $32 \times 32$ point prompts, supporting extremely limited partial queries. Meanwhile, it uses the raw CLIP embeddings which have no subordination information and also degrade after dimensionality reduction, resulting in failure concerning detailed part-level queries (see Figure 1).

The above discussion shows the dilemma of balancing query accuracy and part-level localization ability when building a language-embedded radiance field.

To get out of such a predicament, in this paper, we present the FMLGS, which provides both accurate query results and strong part-level localization ability that supports detailed description. Given a set of posed images, we obtain hierarchical semantics by processing in a subordination-compliant order. FMLGS first extracts all SAM masks in a single frame and filters redundantly overlapped ones to get object-level masks. Then for each object, we divide for its image tile and extract again to get masks of corresponding object parts. After sending them through CLIP and acquire object- and part-level embeddings separately, we designed a semantic deviation strategy to enrich object parts with subordination information. Then, we use SAM2 [10] for each object and its parts to have a consistent identity across views. While directly training a CLIP feature field in high dimensionality could be time-consuming, we map the features to lower 3D space based on identity, and train object- and part-level features in parallel. At the inference stage, open-vocabulary queries will go through both levels and generate pixel-aligned results based on relevancy scores.

During experiments, we found FMLGS not only achieves first-place performance concerning both speed and accuracy compared with other state-of-the-art 3D segmentation and semantic field methods, but also shows the best localization ability upon specified part-level targets. Meanwhile, we further integrate FMLGS as a virtual agent that can interactively navigate through 3D scenes, locate targets, and respond to user demands through a chat interface, which demonstrates the potential of our work to be further expanded and applied in the future.

In summary, the principal contributions of this work include:

- To our best knowledge, this work is the first to provide accurate part-level open-vocabulary localization ability among language-embedded radiance fields.
- We propose the multilevel feature mapping strategy with semantic deviation, which not only provides feature consistency and training efficiency, but also solves the language ambiguity issue that keeps hindering the correct target localization using natural part-level description through CLIP relevancy.
- We provide an effective 3D semantic basis along with an example of it being integrated into AI agents, proving promising future applications.

## 2 RELATED WORK

### 2.1 NeRF and 3DGS

Neural Radiance Fields (NeRF) and 3D Gaussian Splatting (3DGS) have revolutionized 3D scene modeling and rendering. After first introduction, NeRF [4] has advanced largely with further innovations from mip-nerf [11] to NeRF-MAE [12]. 3D Gaussian Splatting, detailed by [6] and expanded upon in work [13] on dynamic 3D Gaussians, optimizes the rendering of point clouds with Gaussian kernels for real-time applications. Later contributions [14]–[17] underscore the ongoing enhancements and versatility of 3DGS in handling increasingly complex rendering tasks. These works serve as a firm basis for constructing 3D language-embedded radiance fields from 2D semantics.

### 2.2 2D and 3D Segmentation

The field of 2D image segmentation has undergone remarkable progress by adopting the Transformer architecture, notably through SEgmentation TRansformer (SETR) [18], alongside studies [19]–[22]. Although works like CGRSeg [23] provides excellent accuracy and efficiency, innovations such as SAM [8] and SEEM [24] utilize various kinds of prompt for segmentation, inspired many mask-based researches. Meanwhile, image captioning has also become a promising method of retrieving image semantics [25].

Advancements in 3D segmentation have paralleled those in 2D, with all kinds of innovations including Cylinder3D [26] for LiDAR semantic segmentation in driving scenes, and 3D semantic segmentation within point clouds [27], [28]. Other developments in radiance fields have also refined the precision and applicability of 3D segmentation techniques [29]–[31]. In particular, methods based on SAM masks [32], [33] achieved part-level segmentation based on user point prompts, but open-vocabulary text-guided 3D segmentation remains a challenge due to the lack of compatible semantic scene construction.

### 2.3 Language-embedded Radiance Fields

Distilling language features into radiance fields like NeRF and 3DGS has been thoroughly explored. Zhi et al. [34] introduced semantic layers into NeRF, setting the stage for enriched scene understanding. This has been further developed in studies such as ISRF [29], DFF [35], N3F [36], and LERF [2], which refine 3D visualization and enable semantic segmentation. In special scenes, work such as MA-52 [37] and COTR [38] have also served as good examples.

Furthermore, practice through 3DGS shows more promising efficiency and accuracy. LEGaussians [5] proposes a feature-lifting strategy within 3DGS, achieving faster queries and smoother relevancy maps. The SAGA framework [39] demonstrates the integration of detailed 2D segmentation outcomes into 3D models, improving the accuracy of the query. Gaussian Grouping [7] and LangSplat [1] further explore the incorporation of SAM masks into 3DGS, allowing more consistent and accurate language processing. However, these methods fail to accurately distinguish and locate both object- and part-level targets upon open-vocabulary queries, thus further applications are still limited.

# 3 METHOD

## 3.1 Overview

As shown in Figure 2, given a set of posed images, FMLGS initializes from a single frame in the sequence. We extract SAM masks in *everything* mode and filters redundantly overlapped ones to get object-level masks. For each segmented object, another SAM extraction and filtering is conducted for corresponding part-level masks. After sending the image tiles of objects and their parts through CLIP, we will have separate language features. We use semantic deviation to enrich object parts with subordination information and generate new language features. Then, we use SAM2 for each object and its parts to have consistent identity across views, so that we can now map the high-dimensional language features to lower space. The mapped object- and part-level low-dim features will be used for separate supervision, and at the inference stage, open-vocabulary queries will go through both levels and generate pixel-aligned results based on relevancy scores.

## 3.2 Multilevel Extraction

In order to get a consistent identity of multilevel targets, we need to acquire correct lists of all objects and the corresponding parts. We achieve this purpose by extracting in subordination-compliant order.

We first use SAM mask generation on the image scale to locate object-level targets. Without manual prompts, automatically generated masks usually overlap due to point ambiguity, thus it is essential to filter out redundant masks. Specifically, all masks are sorted in descending order by area size, and assigned to empty canvas from the start. If the current mask area is completely not taken, then it results in a successful assignment. Otherwise, the current mask will be considered a part-level mask caused by point ambiguity and discarded. After this process, we have a list of different object-level targets based on used masks.

Then, we proceed to part-level extraction. For each object extracted previously, we conduct SAM mask generation solely on its image tile so that more part-level masks will be generated. When filtering masks, at this step ascending order is used to reserve part-level instead of discarding them. Meanwhile, hollow masks are further filtered, which are considered background parts besides objects and meaningful components.

After the above process, the extraction of multilevel targets and their correspondence relationship has been completed.

## 3.3 Multilevel Feature Mapping With Semantic Deviation

Generating language features of every frame for training is both time-consuming and inconsistent. Instead, we use mapped features for efficiency and accuracy. FastLGS [9] uses cross-view grid mapping to generate low-dim features, but the mapping strategy is solely made for object-level targets and cannot support part-level feature mapping. Therefore, we designed a new multilevel feature mapping strategy for both object- and part-level targets.

**Semantic Deviation For Part-level Features.** By sending the image tiles of the multilevel targets through CLIP, we will have their separate language features. Although an object-level feature can contribute to a relevancy score concerning queries containing part-level descriptions, a part-level feature has no information of its corresponding object. This is because while all parts are captured within the object image tile they belong to, no object appears in its parts' image tiles. Such a fact causes serious ambiguity issue when it comes to open-vocabulary querying upon components. For example, when querying "A Button of Xbox Wireless Controller" as shown in Figure 1, the raw feature of the Xbox controller wrongfully provides higher relevancy than that of a single button because of CLIP's "bag of words" attribute.

To tackle this problem, we propose semantic deviation. Given object-level raw feature $F_O$ and part-level raw feature $F_P$, the deviated part-level feature $F_P^{'}$ is then calculated:

$$F_P^{'} = (1 - w)F_O + wF_P \qquad (1)$$

where $w$ is the reserving weight of part-level features. While object-level features remain unchanged, they have now been incorporated into deviated part-level features.

Experiments show that this strategy allows part-level targets to be correctly located with the highest relevancy on queries containing part-level descriptions (see Sec 5.4). Furthermore, the support for targets with subordination information helps solve the issue of language ambiguity. When only queries like "buttons" are available, buttons belonging to different objects cannot be distinguished and will be provided together. This severely limits real-world applications, such as when a robot is commanded to press a specific button. With deviated semantics, queries including "A Button of Nintendo Joystick" and "A Button of Xbox Wireless Controller" will be available in the same scene, enabling more possibilities for downstream applications.

**Cross-view Feature Mapping.** FastLGS [9] uses key point and feature similarity to matching multiview targets and map for consistent low-dim features, but this process is unstable and will result in inconsistency concerning complex scenes, which also degrades to match smaller object parts. Therefore, FMLGS chooses to use an identity based cross-view feature mapping strategy. Specifically, the masks extracted from Sec.3.2 are sent as prompts (degrade to point samples) through SAM2 video segmentation module to propagate for consistent and unique identities across views. Every identity is expanded to a three-dimensional vector $\mathbf{f}$, representing the original 512-dim CLIP feature. The mapped discrete features of whole objects and parts will be separately assigned by pixel according to their masks, so that two levels of pixel-aligned semantic ground truth are generated. A mapping dictionary (less than 3MB) is also created for restoring original features at the inference stage. The detailed procedure of semantic deviation and feature mapping is shown in algorithm 1.

## 3.4 Training Features for Gaussians

Mapped low-dimensional feature $\mathbf{f}$ is trained for each gaussian $g$ as in FastLGS.
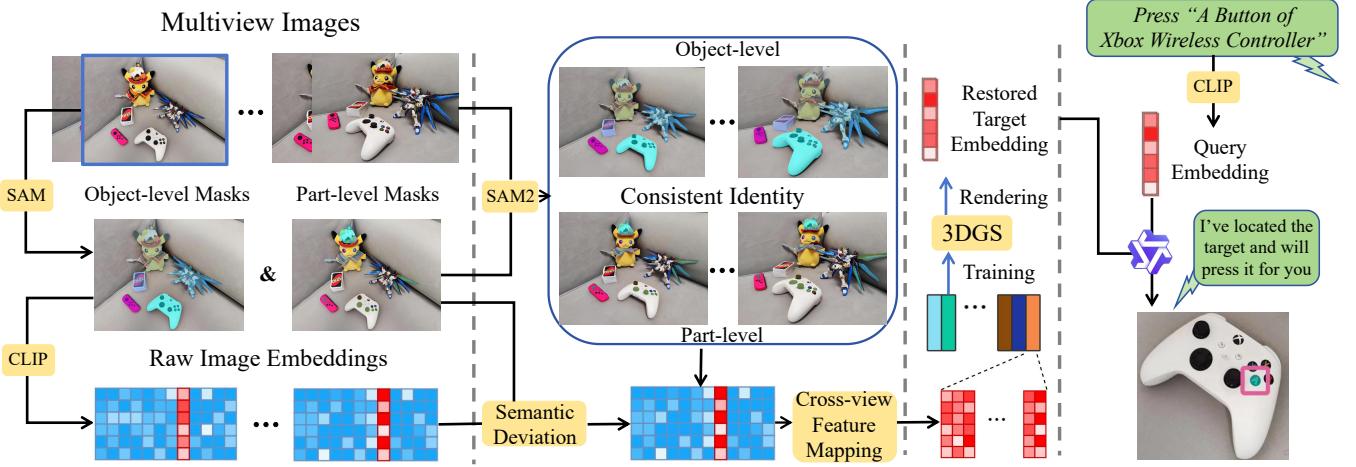
Fig. 2: FMLGS pipeline. Left: Initialization for masks and features. Mid: Feature mapping of object-level embedding and deviated part-level embedding through consistent identities for training and restoring. Right: Query using open-vocabulary prompts, supporting agent integration.

**Rendering Features** Given a camera pose $v$, we compute the feature $\mathbf{F}_{v,p}$ of a pixel by blending a set of ordered Gaussians $\mathcal{N}$ overlapping the pixel similar to the color computation of 3DGS:

$$\mathbf{F}_{v,p} = \sum_{i \in \mathcal{N}} \mathbf{f}_i a_i \prod_{j=1}^{i-1} (1 - a_j), \tag{2}$$

where $a_i$ is given by evaluating a 2D Gaussian with covariance $\sum$ multiplied with a learned per-Gaussian opacity, $\mathbf{f}_i$ is rendered low-dim feature through each Gaussian.

**Optimization** The mapped features follows 3DGS optimization pipeline and especially inherits the fast rasterization for efficient optimization and rendering. The loss function for features is $\mathcal{L}_1$ combined with a D-SSIM term:

$$\mathcal{L}_f = (1 - \lambda)\mathcal{L}_1 + \lambda\mathcal{L}_{D-SSIM}, \tag{3}$$

where $\lambda$ is also fixed to 0.2 in all cases.

### 3.5 Multilevel Localization

After training, given a viewing angle, visible features will be rendered pixel-aligned, and all targets within the field of vision can be repeatedly queried without more rendering before changing the camera position. An input user query is sent through CLIP to generate language embedding, and we calculate a similar relevancy score to LERF. We compute the cosine similarity between image embedding $\phi_{img}$ and canonical phrase embeddings $\phi_{canon}^i$, then compute the pairwise softmax between image embedding and text prompt embedding $\phi_{query}$, so that the relevancy score is:

$$S_{relev} = \min_i \frac{\exp(\phi_{img} \cdot \phi_{query})}{\exp(\phi_{img} \cdot \phi_{canon}^i) + \exp(\phi_{img} \cdot \phi_{query})}. \tag{4}$$

For object retrieval, most methods only use the simple comparison of similarity to determine the result, which is insufficient when part-level targets have become one of the candidates. The general canonical phrases may not be sensitive enough to tell an object and its parts apart, thus,

we designed a two-step localization strategy for multilevel targets.

We first set canonical phrases in equation 4 to "object", "stuff" and "texture" and compute relevancy with every image embedding in the mapping dictionary, where the highest relevancy indicates the preliminary target. Then, we determine the final results based on the following two steps:

**Step 1:** If the query is a kind of object, the preliminary target is the final result since it is the same as other object-level tasks. If the query is a kind of part and the highest relevancy is from part-level embeddings, the corresponding part is the final result. In other cases, it should move to Step 2.

**Step 2:** If the query is a kind of part and the highest relevancy is not from part-level embeddings but from object-level embeddings, the queried part should come from the corresponding object. Therefore, a detailed part-level comparison is further conducted. In this case, the object phrases will be added to the canonical phrases, and the relevancy with all parts of the corresponding object will be computed. The part with the highest relevancy provided by deviated semantics now indicates the queried object part.

After locating it, the mapping dictionary also provides the mapped feature for the corresponding target. Using the rendered pixel-aligned features, we generate a target mask by determining if each pixel's feature matches the target feature within a specified channel tolerance $t$, so that a query is complete. Meanwhile, multiple targets with the same meaning can also be queried simultaneously using top $k$ relevancy adjustably, enabling applications in more scenarios.

## 4 INTERACTIVE AGENT

The interaction with 3D semantic fields is no longer a pure text processing procedure, thus, a common agent structure such as a single chain of thought (CoT) for smart searcher cannot be directly implemented. In real-world scenarios like household robots, a core decision model is required

Fig. 3: Examples of environment-sensitive agents of four different cases.

**Algorithm 1** Semantic Deviation & Feature Mapping

**Input**: Image sequence $\{\mathbf{I}_t|t = 0, 1, ..., T\}$, object-level masks $\{OM_i|i = 0, 1, ..., n\}$, part-level masks $\{PM_{i,j}|i = 0, 1, ..., n; j = 0, 1, ..., m_i\}$, raw object-level CLIP embedding $\{\mathbf{OL}_i|i = 0, 1, ..., n\}$, raw part-level CLIP embedding $\{\mathbf{PL}_{i,j}|i = 0, 1, ..., n; j = 0, 1, ..., m_i\}$ and reserving weight $w$

**Parameter**: Deviated part-level embedding $\mathbf{PL}'$, object-level identities $\mathbf{OI}$, part-level identities $\mathbf{PI}$, prompt for mask propagation $p$, object-level mapped low-dimensional feature $f^o$, part-level mapped low-dimensional feature $f^p$

**Function**: GenPrompt($x$) generates point prompt $y$ from mask $x$, SAM2Prop($y$, $z$) uses SAM2 to propagate for consistent mask identity $id$ from prompt $y$ through image sequence $z$, MapFeature($u$,$v$) stores the $u$ to $v$ feature mapping and generates pixel-aligned mapped feature, GenGTFeature($k$,$z$) generates complete ground truth feature for image sequence $z$ from all mapped feature $k$.

**Output**: Mapping of multilevel semantics for query and ground truth multi-view low dimensional features ($f^o_{gt}$ and $f^p_{gt}$) for training in a scene

1: Let $i = 0, j = 0$.
2: **while** $i \leq n$ **do**
3:     $p = $ GenPrompt($OM_i$)
4:     $\mathbf{OI}_i = $ SAM2Prop($p$, $\mathbf{I}$)
5:     $f^o_i = $ MapFeature($\mathbf{OI}_i$, $\mathbf{OL}_i$)
6:     **while** $j \leq m_i$ **do**
7:         $p = $ GenPrompt($PM_{i,j}$)
8:         $\mathbf{PI}_{i,j} = $ SAM2Prop($p$, $\mathbf{I}$)
9:         $\mathbf{PL}'_{i,j} = (1 - w)\mathbf{OL}_i + w\mathbf{PL}_{i,j}$
10:        $f^p_{i,j} = $ MapFeature($\mathbf{PI}_{i,j}$, $\mathbf{PL}'_{i,j}$)
11:        $j = j + 1$
12:     **end while**
13:     $i = i + 1$
14: **end while**
15: $f^o_{gt} = $ GenGTFeature($f^o$, $\mathbf{I}$)
16: $f^p_{gt} = $ GenGTFeature($f^p$, $\mathbf{I}$)

to decompose a complete user requirement into one or more proxy task sequences that the intelligent agent can execute, and it should also be able to update task sequences based on perceived changes of both user commands and environments.

As Figure 3 shows, there are many challenges when executing the user commands in a 3D space. For example, given a view without seeing the queried objects/parts, it is useful to adaptively adjust the views based on the language-embedded radiance fields via an agent. In addition, user commands may not be executed for reasons such as the command being unclear, the command violating objective facts, or execution may lead to security risks. Therefore, in this section, we further designed a smart agent structure to implement FMLGS for language-guided 3D scene interactions, a suitable agent execution framework as shown in Figure 4.

Once a user command or query has been entered, an outer loop will be initiated for this complete task, where the current execution process is stored. Within each outer loop, one to multiple inner loops are executed. The inner loop corresponds to breaking down a complete user requirement into one or more proxy task sequences that the agent can execute and sequentially complete by calling the accessed functional modules in the inner loop. The inner loop will make decisions on subsequent task sequences based on the feedback of task modules and perceived environmental changes.

### 4.1 Scene Initialization

Apart from reconstruction for scene appearance and semantics, more preparation is needed for agent modules to function properly in a 3D scene.

**Data Preprocess.** The original data contains posed camera information for all training views, which are the most reliable anchors for camera positioning. However, when it comes to actual user interactions, where cameras need to navigate through scenes and constantly update viewing angles, it is impractical to stick to these limited numbers of fixed points. To generate more camera keypoints, we dilate the training camera coordinates by adding new keypoints in the surrounding six directions (up, down, left, right, front and back) for all training cameras, with each new point at $\frac{s}{2}$ away from the original ones ($s$ is average distance of adjacent training camera coordinates). In Figure 5, a comparison of dilating the original camera coordinates in the MIP-360 dataset is shown. It is clear that more varied and available paths can be created with extended keypoints.

Although we can place virtual cameras anywhere in the reconstructed scenes, in reality, we have to consider collision
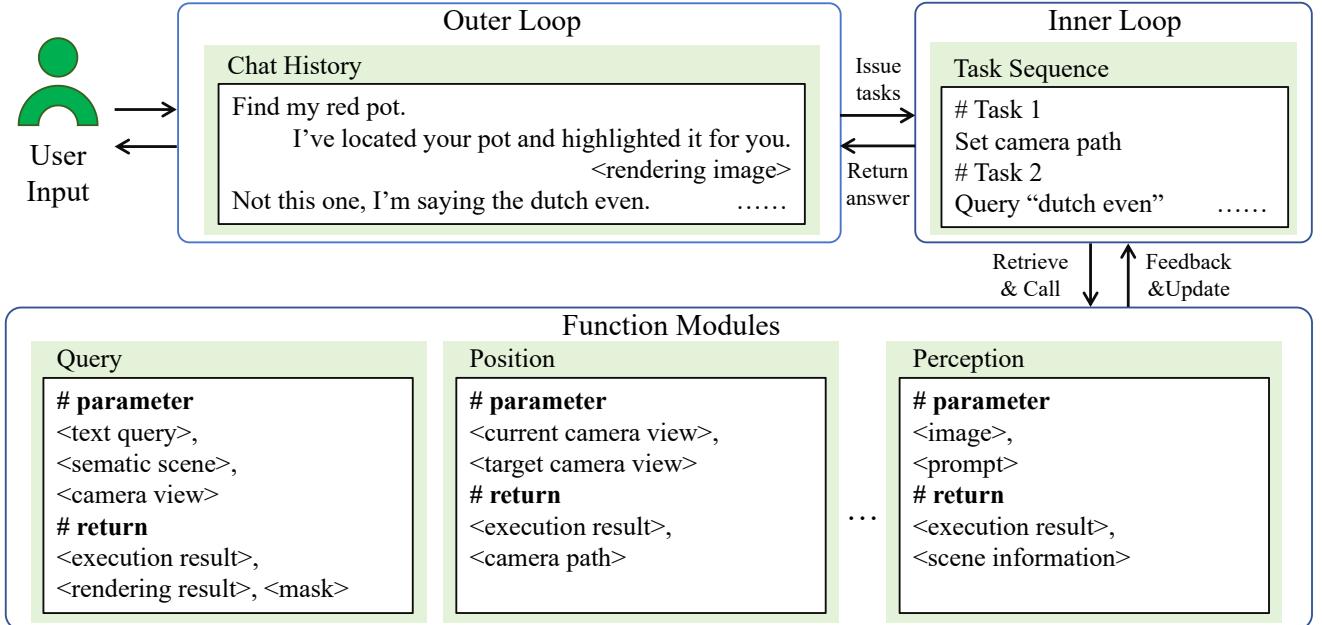
Fig. 4: Agent execution framework. Up: two-stage nested loop with the outer loop issuing main task after accepting user input, and the inner loop disassemble single task to executable task sequence. Down: function modules for different subtasks to be called by the inner loop execution.
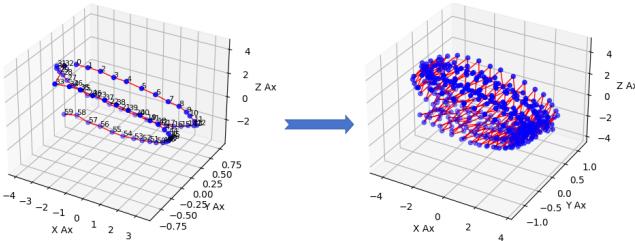


Fig. 5: Illustration of dilating keypoints to generate more view anchors in the scene.

detection and cannot place cameras inside of any objects, which is even more important in actual robot implementation. Therefore, the depth information is also necessary for realistic interactions. In this paper, we use 2DGS [40] to acquire more accurate and smoother depth predictions for the following camera path establishment.

**Path Establishment.** When performing tasks such as locating objects, a path from the current view to the view where the target exists must be shown, otherwise, the user may still be unaware of its whereabouts.

To establish an available camera path, we consider it a shortest path problem in the 3D space. While only predicted depth is provided as spatial restrictions, we regard all keypoints as discrete path points and initialize connectivity. Specifically, we sample points on the path from one keypoint to another as the camera center, and predict depth with fixed view direction. If the depth has non-positive values, it means the view frustum is likely to cut through objects or at least unreliable, then we label the points as "disconnected". For the connected ones, their distance is computed by Euclidean distance. After initializing connectivity, we can establish the shortest path using the Dijkstra algorithm.

## 4.2  Visible Interaction

The agent connects to the function modules through data exchange, but also needs to respond to users with visible feedback. In this scenario, visible feedback is constantly updated in scene rendering results through moving cameras.

**View Generation.** In Sec 4.1, a path to queried targets is established, but it only provides several keypoints, which is not enough to show users an unabridged navigation process. To generate a natural view change process, we interpolate the original path sequence. Given path $\{k_1, k_2, ..., k_n\}$, the transformation matrix $T$ for view sequence is calculated by:

$$T_i = T^{k_1} + \frac{i}{M} \sum_{v \in (2,n)} (T^{k_v} - T^{k_v - 1}), i \in [0, M], \quad (5)$$

where $M$ is the set frame count. The rotation matrix $R$ is also interpolated in this way. In this paper, we set $M$ to 150 and use the interpolated camera setting sequence to render view change frame set, so that we could generate 2 to 4 seconds smooth navigation processes of more than 30 FPS.

**Initiative Movement.** The automatic provided query results shown in Figure 3 may not always satisfy user requirements. For example, sometimes, users want to take a closer look at the target or other surrounding objects. Therefore, we further enable the camera to freely move around the scene based on the user commands. A rotation matrix $R$ is described by directional vectors, thus, a forward camera movement can be calculated by updating transformation matrix $T$:

$$T_{forward} = T + d \cdot R \cdot [0, 0, 1]^T, \quad (6)$$

where $d$ is the moving distance, and movement in other directions can be calculated similarly.

The users can demand a movement by specifying a direction and optionally with a distance, then the camera parameters can be updated and show detailed view change process through the above method.

## 5 EXPERIMENTS

In this section, we first show the speed and quality of open-vocabulary object retrieval in comparison with other state-of-the-art methods through quantitative experiments. We also provide illustration for superior part-level localization. Ablation studies are conducted to demonstrate the rationality of semantic deviation-based design. Meanwhile, we further provide examples of integrating FMLGS with the large language model to serve as a part-level interactive 3D agent for real-world applications.

### 5.1 Basic Setups

**Datasets.** For quantitative experiments, we train and evaluate the models on datasets, including SPIn-NeRF [41], LERF [2] and 3D-OVS [42]. We also tested on the MIP-360 dataset [43] for examples and downstream applications.

**Implementation Details.** We use the same OpenClip ViT-B/16 model as that in LERF [2] and the SAM ViT-H model as that in LangSplat [1]. We train the features and scenes in 3DGS for 30,000 iterations. Reserving weight $w$ is set to 0.7. While the original time calculation of LangSplat puts aside feature rendering time and feature reconstructing time, here we compute the query time of the whole query process as LERF does. The tested LERF masks are regions with relevancy higher than 20% after normalization. The normalization for each query is from 50% (less relevant than canonical phrases) to the maximum relevancy, which is identical to the visualization strategy of LERF. All results are reported running on a single TITAN RTX GPU. For agent tests, we implement Qwen-Plus as our core language model in the above-detailed framework for decision-making.

### 5.2 Comparisons

For quantitative experiments, we test different methods on the SPIn-NeRF dataset, the LERF dataset and the 3D-OVS dataset. We show the quality of FMLGS-generated masks by comparing them with other state-of-the-art 3D segmentation methods, including the multi-view segmentation of SPIn-NeRF (MVSeg) [41] and SA3D [32], and compare the language retrieval ability with LERF [2], LEGaussians [5], LangSplat [1] and FastLGS [9].

**SPIn-NeRF dataset.** We evaluate the IoU and pixel accuracy of masks with provided ground truth (1008 × 567), and also show the time consumption for each query, which is omitted in MVSeg and SA3D because they do not support single-view queries. For LERF, LangSplat, LEGaussians, FastLGS and FMLGS, we use the same text queries for the target segmentation objects to generate masks. Other methods all follow their original settings when tested on this dataset. We also provide the 2D segmentation results generated by SAM based on manual point prompts of original scene images for comparison. As shown in Table 1, FMLGS generates masks of competitive quality compared

TABLE 1: Quantitative Results on SPIn-NeRF dataset.

| Method | mIoU (%) | mPAcc (%) | mTime (s) |
|---|---|---|---|
| SAM(2D) [8] | 95.7 | 99.2 | 0.05 |
| MVSeg [41] | 89.5 | 94.6 | - |
| SA3D [32] | 90.9 | 97.4 | - |
| LERF [2] | 81.0 | 85.2 | 30.2 |
| LEGaussians [5] | 89.3 | 94.7 | 1.04 |
| LangSplat [1] | 92.2 | 98.1 | 1.43 |
| FastLGS [9] | 93.1 | 98.3 | 0.31 |
| FMLGS | **94.2** | **98.7** | **0.30** |

TABLE 2: Quantitative Results of localization accuracy on LERF dataset.

| Method | ramen | figurines | teatime | kitchen | o.a. |
|---|---|---|---|---|---|
| LERF [2] | 61.9 | 75.5 | 84.8 | 70.2 | 73.1 |
| LEGaussians [5] | 78.6 | 73.7 | 85.6 | 90.1 | 82.0 |
| LangSplat [1] | 73.2 | 80.4 | 88.1 | 95.5 | 84.3 |
| FastLGS [9] | 84.2 | 91.4 | 95.0 | 94.7 | 91.3 |
| FMLGS | **89.2** | **94.3** | **96.7** | **96.2** | **94.1** |

with other methods, and also shows the query speed of the first level.

**LERF dataset.** The LERF dataset contains several in-the-wild scenes and is much more challenging, which strongly requires zero-shot abilities. Visualized examples in both LERF and MIP-360 scenes are shown in Figure 6, where FMLGS provides the most complete and smoothest results for object retrieval. We report localization accuracy for the 3D object localization task following LERF [2] with ground truth annotations provided by LangSplat [1] (resolution around 985 × 725). While localization accuracy is reaching saturation for later methods, we further provide IoU results for comparison. Data results are shown in Table 2 and 3, which demonstrate FMLGS's advantages in natural language retrieval.

**3D-OVS dataset.** We also compare with 2D-based open-vocabulary segmentation methods including ODISE [44] and OV-Seg [45] along with 3D-based methods including 3D-OVS [42], LERF [2], LEGaussians [5], LangSplat [1] and FastLGS [9]. Results are provided in Table 4, where our method can outperform both 2D and 3D methods. Figure 6 also provides a visualization of part-level results, where

TABLE 3: Quantitative Results of mIoU scores (%) on LERF dataset.

| Method | ramen | figurines | teatime | kitchen | o.a. |
|---|---|---|---|---|---|
| LERF [2] | 28.2 | 38.6 | 45.0 | 37.9 | 37.4 |
| LEGaussians [5] | 34.2 | 47.2 | 58.6 | 50.1 | 47.5 |
| LangSplat [1] | 51.2 | 44.7 | 65.1 | 44.5 | 51.4 |
| FastLGS [9] | 56.2 | 61.4 | 59.3 | 48.5 | 56.4 |
| FMLGS | **73.2** | **72.4** | **81.8** | **64.3** | **72.9** |

TABLE 4: Quantitative Results of mIoU scores (%) on 3D-OVS dataset.

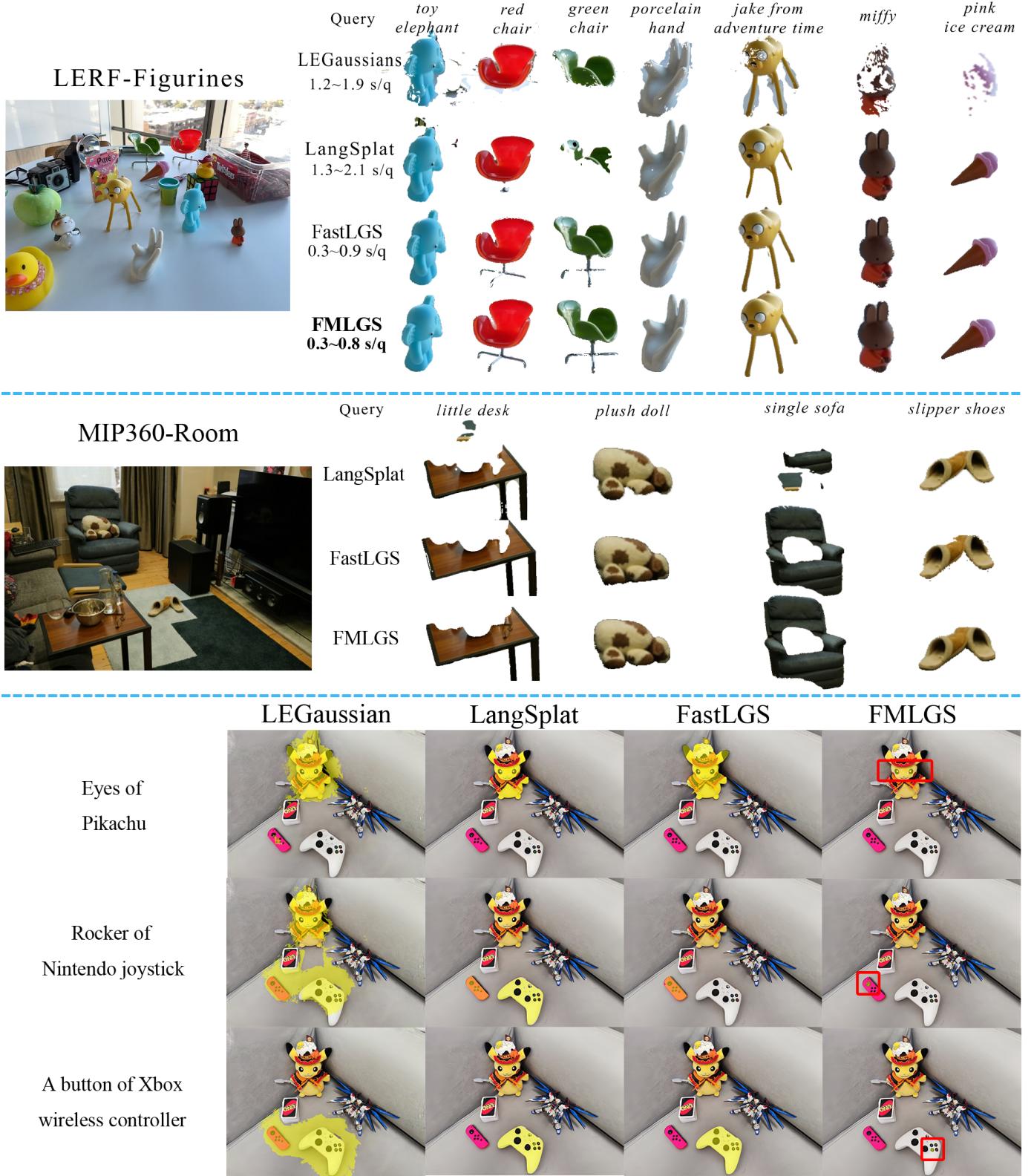| Method | bed | bench | room | sofa | o.a. |
|---|---|---|---|---|---|
| ODISE [44] | 55.6 | 30.1 | 53.5 | 49.3 | 47.1 |
| OV-Seg [45] | 79.8 | 88.9 | 71.4 | 66.1 | 76.6 |
| 3D-OVS [42] | 89.5 | 89.3 | 92.8 | 74.1 | 86.4 |
| LERF [2] | 76.2 | 59.1 | 56.4 | 37.6 | 57.3 |
| LEGaussians [5] | 45.7 | 47.4 | 44.7 | 48.2 | 46.5 |
| LangSplat [1] | 92.6 | 93.2 | 94.1 | 89.3 | 92.3 |
| FastLGS [9] | 94.7 | 95.1 | 95.3 | 90.6 | 93.9 |
| FMLGS | **95.7** | **96.3** | **96.8** | **95.2** | **96.0** |

Fig. 6: Visualized results of object retrieval and part-level localization on different scenes.

LangSplat fails even with its multilevel masks and FMLGS proves to be the only method that supports part-level localization with detailed description.

## 5.3 Downstream Applications

FMLGS's strong target retrieval ability is also capable of supporting various downstream 3D applications. In this section, we apply FMLGS to both language driven 3D segmentation and object inpainting.

Fig. 7: Results of language driven 3D segmentation and object inpainting.

TABLE 5: Ablation on feature mapping. FM is for FastLGS's feature matching strategy, and IM is for our identity-based matching.

| Component | | | Performance | | |
|---|---|---|---|---|---|
| 3DGS | FM | IM | Pre(min) | mIoU(%) | mTime(s) |
| | | | 10 | 83.3 | 20.1 |
| ✓ | | | OOM | OOM | OOM |
| ✓ | ✓ | | 30 | 95.1 | 0.98 |
| ✓ | ✓ | ✓ | 4 | 97.2 | 0.73 |

**Language Driven 3D Segmentation.** Most 3D segmentation methods are based on manual positional prompting such as clicks or frame selection, which cannot be directly applied to real-world applications such as embodied AI, because it is impossible for an automatic robot to rely on constant human input when it is required to retrieve an object. Ideally, specified targets in a 3D scene should be immediately processed based on user description. Therefore, we integrate FMLGS with Segment Any 3D Gaussians (SAGA) [39] and use FMLGS's strong open-vocabulary localization ability to prompt segmentation. Results are shown in Figure 7 (up), where we can obtain segmented 3D objects solely using a text prompt.

**Language Driven Object Inpainting.** When virtually interacting with scene objects, it is also necessary to update the scenes. If an object should be removed, inpainting is one way of updating the 3D representations. While 3D object inpainting requires accurate multi-view segmentation masks, the consistently built FMLGS provides perfect masks for guidance. We integrate FMLGS with the SPIn-NeRF [41] inpainting pipeline to achieve direct language-driven object inpainting. As shown in Figure 7 (down), we can accurately inpaint the described scene targets.

In conclusion, all the above experiments show that FMLGS has a very strong open-vocabulary localization ability along with interactive efficiency. Compared with other state-of-the-art methods, it further supports grounding and querying part-level semantics. Experiments on integration also prove its capability of being applied to downstream applications.

## 5.4 Ablation Study

In this section, we conduct ablation to validate the necessity of our feature mapping strategy and design for part-level retrieval.

**Feature Mapping.** We compare the performance of training with raw CLIP features in NeRF and 3DGS, and

TABLE 6: Ablation on multilevel extraction (ME), semantic deviation (SD) and multilevel localization (ML).

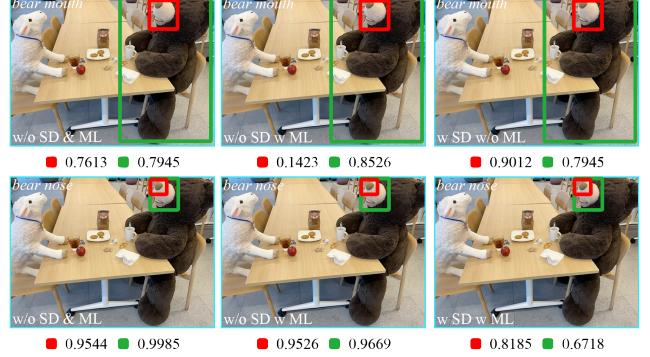| Component | | | Performance | |
|---|---|---|---|---|
| ME | SD | ML | Recall(%) | mTime(s) |
| | | | 1.1 | 0.67 |
| ✓ | | | 5.3 | 0.69 |
| ✓ | ✓ | | 41.6 | 0.69 |
| ✓ | ✓ | ✓ | 66.7 | 0.75 |



Fig. 8: Ablation on semantic deviation and multilevel localization.

mapping features with/without identity based matching. The results are shown in Table 5. Directly rendering CLIP features from NeRF is time-consuming, and inconsistent interpolated features will also lead to low accuracy. While training with raw features in 3DGS will lead to out-of-memory issues, mapping features using FastLGS's matching strategy requires too much preprocess time and also has certain inconsistencies. Our mapping strategy through identity based matching not only provides the least preprocess time and highest accuracy, but also leads to faster queries compared with feature matching, which adds redundant matching information caused by inconsistency.

**Part-level Design.** We compare the performance of multilevel extraction, semantic deviation and multilevel localization for part-level retrieval. We report the recall rate and average query time for part-level queries in Table 6. Using only automatically generated SAM masks for filter makes it an object oriented method, which will fail in almost all part-level queries and provide wrong object masks as FastLGS results shown in Figure 6. After extracting multilevel masks, simple parts in controlled scenes can be localized, but because of incorrect semantics, the results are often wrong part-level targets as LangSplat results shown in Figure 6. Using semantic deviation helps locate parts with clearer view and more information, solving more queries with subordination descriptions. However, the "bag of words" behavior of CLIP may still lead to possible wrong localization on objects instead of their parts. By using the two-step multilevel localization strategy in Section 3.5, this problem is successfully resolved and supports more queries upon part-level targets. Comparison examples are shown in Figure 8, where queries (top left on images) and similarity for framed targets are provided. For the first query of "bear mouth", the whole bear will provide higher relevancy than the mouth part. But solely adding a multilevel localization module cannot solve this problem because this module only
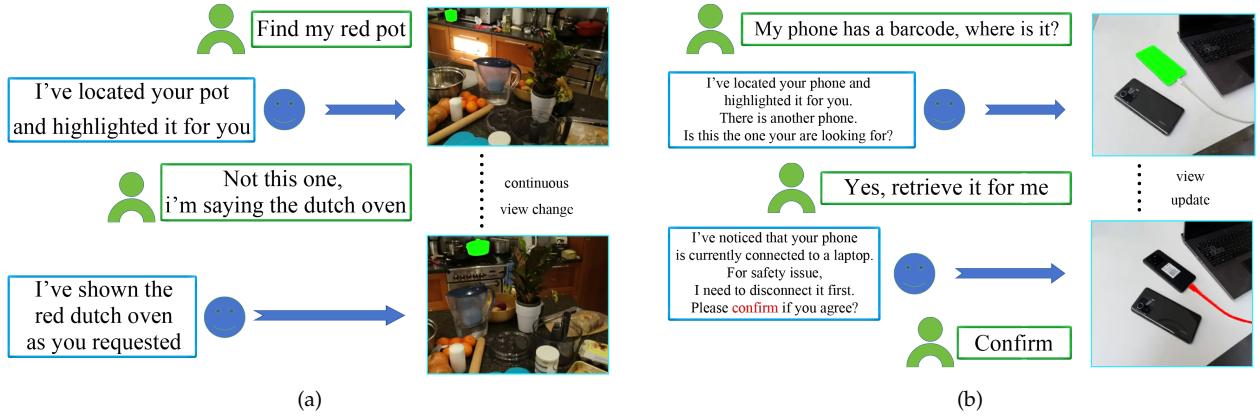
Fig. 9: Illustration of continuous agent execution based on user requirements. (a) Case I. (b) Case II.

helps with part-level comparison. With semantic deviation, the part-level target will correctly gain higher relevancy. However, when it comes to part-level comparison as "bear nose" is queried, the "bear mouth" part will wrongfully provide higher similarity than the nose because this similarity is from comparing with regular canonical phrases. Solely adding multilevel localization also fails in this situation, because their original language features are not sufficient enough for distinguishing. After completely implementing semantic deviation and multilevel localization, the resulting similarity shows clear localization.

### 5.5 Interactive Agents

We conduct experiments on our design of a 3D interactive agent. We examine different prompts and search for the most suitable one for this task, and conduct tests within different scenes for pure text-guided 3D environment interactions. Experiments show that our method can continuously locate various targets and naturally moving around viewing camera and show results to the user.

The illustration is provided in Figure 9. An example in the MIP-360 dataset [43] is shown in Figure 9a, where we ask the agent to find a target and gradually add more description. The agent is able to revise its result and give the correct target. Another example in Figure 9b shows that our method can discover risks such as unauthorized or dangerous operations and automatically update the following task sequences.

## 6 DISCUSSION

The development of this research has revealed many valuable insights. In this section, we discuss a couple of questions that may contribute to future research concerning this field.

**Do we really need feature compression?** Many previous methods, such as LERF and LangSplat, tend to compress the original CLIP features. They either use averaged CLIP features to reduce feature quantity or use MLPs to reduce and restore feature dimensionality. The result is that these methods all have degraded language features and certain inconsistency, leading to failure in demanding queries and poor capability in part-level localization. Instead, we believe

it is necessary to retain the original features as long as a consistent mapping is created. In FastLGS and FMLGS, no compression is implemented and they have the first-tier performance.

**Limitation.** This method relies on accurate and consistent masks. Although the current segmentation models can easily provide object-level masks, generating masks for irregular common object parts remains a challenge. For some parts, even if they are successfully segmented in most views, they could still be ignored in some views where the parts are too far away from the camera, leading to inconsistency and degraded results (eg. some buttons). Therefore, in real-world applications, a single reconstruction may not be enough to obtain all consistent and accurate part-level information (depending on how detailed the scene is captured). For such scenarios, a close-up recapture and update should be necessary for demanding part-level interaction tasks.

Meanwhile, although this language embedded gaussians ground detailed semantics for each target, no natural connection has been established within the fields. Future tasks will require the model to understand beyond simple target description, such as "give me the largest bottle placed on the right of the book shelf". Later works in this field may need to combine the vision language model to fill this gap.

## 7 CONCLUSION

In this paper, we present FMLGS, an approach that supports part-level open-vocabulary query within 3D Gaussian Splatting (3DGS). We propose an efficient pipeline for building and querying consistent object- and part-level semantics based on the Segment Anything Model 2 (SAM2). We also designed a semantic deviation strategy to solve the problem of language ambiguity among object parts. Once trained, FMLGS can query both objects and their describable parts using natural language. Comparisons with other state-of-the-art methods prove that our method can not only better locate specified part-level targets, but also achieve first-place performance concerning both **speed** and **accuracy**. Moreover, we further designed an agent framework that integrates FMLGS as a virtual assistant that can interactively locate targets and respond to user demands within 3D scenes, showcasing the potential of our work to be further expanded and applied in the future.

# REFERENCES

[1] M. Qin, W. Li, J. Zhou, H. Wang, and H. Pfister, "Langsplat: 3d language gaussian splatting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024.

[2] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik, "Lerf: Language embedded radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 19 729–19 739.

[3] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning (ICML)*. PMLR, 2021, pp. 8748–8763.

[4] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[5] J.-C. Shi, M. Wang, H.-B. Duan, and S.-H. Guan, "Language embedded 3d gaussians for open-vocabulary scene understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024.

[6] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, July 2023.

[7] M. Ye, M. Danelljan, F. Yu, and L. Ke, "Gaussian grouping: Segment and edit anything in 3d scenes," in *European Conference on Computer Vision (ECCV)*, 2024.

[8] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollar, and R. Girshick, "Segment anything," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 4015–4026.

[9] Y. Ji, H. Zhu, J. Tang, W. Liu, Z. Zhang, X. Tan, and Y. Xie, "Fastlgs: Speeding up language embedded gaussians with feature grid mapping," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2025.

[10] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer, "Sam 2: Segment anything in images and videos," *arXiv preprint arXiv:2408.00714*, 2024. [Online]. Available: https://arxiv.org/abs/2408.00714

[11] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, "Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 5855–5864.

[12] M. Z. Irshad, S. Zakharov, V. Guizilini, A. Gaidon, Z. Kira, and R. Ambrus, "Nerf-mae: Masked autoencoders for self-supervised 3d representation learning for neural radiance fields," in *European Conference on Computer Vision (ECCV)*, 2024.

[13] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan, "Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis," in *International Conference on 3D Vision (3DV)*, 2024.

[14] Y. Liu, H. Guan, C. Luo, L. Fan, J. Peng, and Z. Zhang, "City-gaussian: Real-time high-quality large-scale scene rendering with gaussians," in *European Conference on Computer Vision (ECCV)*, 2024.

[15] Z. Zhu, Z. Fan, Y. Jiang, and Z. Wang, "Fsgs: Real-time few-shot view synthesis using gaussian splatting," in *European Conference on Computer Vision (ECCV)*, 2024.

[16] Y. Chen, H. Xu, C. Zheng, B. Zhuang, M. Pollefeys, A. Geiger, T.-J. Cham, and J. Cai, "Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images," in *European Conference on Computer Vision (ECCV)*, 2024.

[17] Q. Tian, X. Tan, Y. Xie, and L. Ma, "Drivingforward: Feed-forward 3d gaussian splatting for driving scene reconstruction from flexible surround-view input," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.

[18] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, and L. Zhang, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 6881–6890.

[19] B. Cheng, A. Schwing, and A. Kirillov, "Per-pixel classification is not all you need for semantic segmentation," in *Advances in Neural Information Processing Systems (NeurIPS)*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34. Curran Associates, Inc., 2021, pp. 17 864–17 875.

[20] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," in *Advances in Neural Information Processing Systems (NeurIPS)*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34. Curran Associates, Inc., 2021, pp. 12 077–12 090.

[21] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Masked-attention mask transformer for universal image segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 1290–1299.

[22] T. Sun, Z. Zhang, X. Tan, Y. Peng, Y. Qu, and Y. Xie, "Uni-to-multi modal knowledge distillation for bidirectional lidar-camera semantic segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2024.

[23] Z. Ni, X. Chen, Y. Zhai, Y. Tang, and Y. Wang, "Context-guided spatial feature reconstruction for efficient semantic segmentation," in *The 18th European Conference on Computer Vision (ECCV)*, 2024.

[24] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Wang, L. Wang, J. Gao, and Y. J. Lee, "Segment everything everywhere all at once," in *Advances in Neural Information Processing Systems (NeurIPS)*, A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 19 769–19 782.

[25] X. Yang, Y. Wu, M. Yang, H. Chen, and X. Geng, "Exploring diverse in-context configurations for image captioning," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36. Curran Associates, Inc., 2023, pp. 40 924–40 943. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2023/file/804b5e300c9ed4e3ea3b073f186f4adc-Paper-Conference.pdf

[26] H. Zhou, X. Zhu, X. Song, Y. Ma, Z. Wang, H. Li, and D. Lin, "Cylinder3d: An effective 3d framework for driving-scene lidar semantic segmentation," *arXiv preprint arXiv:2008.01550*, 2020.

[27] X. Tan, Q. Ma, J. Gong, J. Xu, Z. Zhang, H. Song, Y. Qu, Y. Xie, and L. Ma, "Positive-negative receptive field reasoning for omni-supervised 3d segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2023.

[28] T. Sun, Z. Zhang, X. Tan, Y. Qu, and Y. Xie, "Image understands point cloud: Weakly supervised 3d semantic segmentation via association learning," *IEEE Transactions on Image Processing (TIP)*, 2024.

[29] R. Goel, D. Sirikonda, S. Saini, and P. J. Narayanan, "Interactive segmentation of radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 4201–4211.

[30] S. Tang, W. Pei, X. Tao, T. Jia, G. Lu, and Y.-W. Tai, "Scene-generalizable interactive segmentation of radiance fields," in *Proceedings of the 31st ACM International Conference on Multimedia (ACM MM)*, 2023, pp. 6744–6755.

[31] Q. Gu, Z. Lv, D. Frost, S. Green, J. Straub, and C. Sweeney, "Ego-lifter: Open-world 3d segmentation for egocentric perception," in *European Conference on Computer Vision (ECCV)*, 2024.

[32] J. Cen, Z. Zhou, J. Fang, W. Shen, L. Xie, D. Jiang, X. Zhang, Q. Tian *et al.*, "Segment anything in 3d with nerfs," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36, pp. 25 971–25 990, 2023.

[33] C. M. Kim, M. Wu, J. Kerr, M. Tancik, K. Goldberg, and A. Kanazawa, "Garfield: Group anything with radiance fields," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[34] S. Zhi, T. Laidlow, S. Leutenegger, and A. J. Davison, "In-place scene labelling and understanding with implicit scene representation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 15 838–15 847.

[35] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola, "Nerf-supervision: Learning dense object descriptors from neural radiance fields," in *International Conference on Robotics and Automation (ICRA)*, 2022, pp. 6496–6503.

[36] V. Tschernezki, I. Laina, D. Larlus, and A. Vedaldi, "Neural feature fusion fields: 3d distillation of self-supervised 2d image representations," in *International Conference on 3D Vision (3DV)*, 2022, pp. 443–453.

[37] D. Guo, K. Li, B. Hu, Y. Zhang, and M. Wang, "Benchmarking micro-action recognition: Dataset, method, and application," *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. 34, no. 7, pp. 6238–6252, 2024.

[38] Q. Ma, X. Tan, Y. Qu, L. Ma, Z. Zhang, and Y. Xie, "Cotr: Compact occupancy transformer for vision-based 3d occupancy prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024, pp. 19 936–19 945.

[39] J. Cen, J. Fang, C. Yang, L. Xie, X. Zhang, W. Shen, and Q. Tian, "Segment any 3d gaussians," *arXiv preprint arXiv:2312.00860*, 2023.

[40] B. Huang, Z. Yu, A. Chen, A. Geiger, and S. Gao, "2d gaussian splatting for geometrically accurate radiance fields," in *ACM SIGGRAPH*, 2024.

[41] A. Mirzaei, T. Aumentado-Armstrong, K. G. Derpanis, J. Kelly, M. A. Brubaker, I. Gilitschenski, and A. Levinshtein, "Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 20 669–20 679.

[42] K. Liu, F. Zhan, J. Zhang, M. Xu, Y. Yu, A. E. Saddik, C. Theobalt, E. Xing, and S. Lu, "Weakly supervised 3d open-vocabulary segmentation," *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[43] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 5470–5479.

[44] J. Xu, S. Liu, A. Vahdat, W. Byeon, X. Wang, and S. De Mello, "Open-vocabulary panoptic segmentation with text-to-image diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 2955–2966.

[45] F. Liang, B. Wu, X. Dai, K. Li, Y. Zhao, H. Zhang, P. Zhang, P. Vajda, and D. Marculescu, "Open-vocabulary semantic segmentation with mask-adapted clip," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 7061–7070.

**Xin Tan** is currently the Research Professor (Zijiang Young Scholar) with School of Computer Science and Technology, East China Normal University, China. Before that, he was the Associate Research Professor at ECNU. He received dual Ph.D. degrees in Computer Science from Shanghai Jiao Tong University and City University of Hong Kong in 2022. He received his B.Eng. degree in Automation from Chongqing University, China in 2017. His research interests lie in computer vision and deep learning. He serves as a program committee member/reviewer for CVPR, ICCV, ECCV, AAAI, IJCAI, IEEE TPAMI, TIP and IJCV. He was selected for the Young Elite Scientists Sponsorship Program by CAST. He also serves as the associate editor for Pattern Recognition and Visual Computer.



**Yuzhou Ji** is currently a forth-year undergraduate student at the School of Computer Science and Technology, East China Normal University, China. He is going to pursue his master degree at the Department of Computer Science and Engineering, Shanghai Jiao Tong University in 2025. His research interests cover 3D reconstruction and scene understanding.



**He Zhu** is now a forth-year undergraduate student in the School of Computer Science and Technology, East China Normal University. He is going to pursue a master's degree in Electronic Engineering at Shanghai Jiao Tong University in 2025. His research interests cover scene reconstruction and neural rendering.



**Yuan Xie** received the PhD degree in Pattern Recognition and Intelligent Systems from the Institute of Automation, Chinese Academy of Sciences (CAS), in 2013. He is currently a full professor with the School of Computer Science and Technology, East China Normal University, Shanghai, China. His research interests include image processing, computer vision, machine learning and pattern recognition. He has published around 85 papers in major international journals and conferences including the IJCV, IEEE TPAMI, TIP, TNNLS, TCYB, and NIPS, ICML, CVPR, ECCV, ICCV, etc. He also has served as a reviewer for more than 15 journals and conferences. Dr. Xie received the National Science Fund for Excellent Young Scholars 2022.