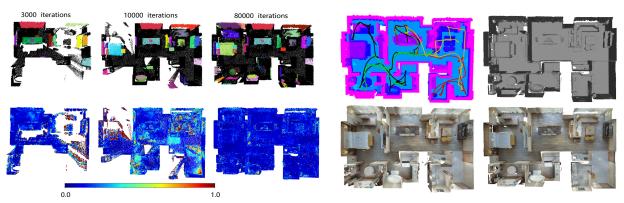
# Multi-robot autonomous 3D reconstruction using Gaussian splatting with Semantic guidance

Jing Zeng<sup>1</sup>, Qi Ye<sup>1\*</sup>, Tianle Liu<sup>1</sup>, Yang Xu<sup>1</sup>, Jin Li<sup>2</sup>, Jinming Xu<sup>1</sup>, Liang Li<sup>1</sup>, Jiming Chen<sup>1</sup>



(a) Segmented results and surface uncertainty

(b) Planned trajectories and reconstructed scene

Fig. 1: (a) Segmented results (Top) and Surface uncertainty (Bottom) from reconstruction iterations; (b) Planned trajectories (Top left, four robots); Reconstructed scene: 3DGS (Bottom left), mesh with texture (Bottom right) and without texture (Top right) from ASH.

Abstract-Implicit neural representations and 3D Gaussian splatting (3DGS) have shown great potential for scene reconstruction. Recent studies have expanded their applications in autonomous reconstruction through task assignment methods. However, these methods are mainly limited to single robot, and rapid reconstruction of large-scale scenes remains challenging. Additionally, task-driven planning based on surface uncertainty is prone to being trapped in local optima. To this end, we propose the first 3DGS-based centralized multi-robot autonomous 3D reconstruction framework. To further reduce time cost of task generation and improve reconstruction quality, we integrate online open-vocabulary semantic segmentation with surface uncertainty of 3DGS, focusing view sampling on regions with high instance uncertainty. Finally, we develop a multi-robot collaboration strategy with mode and task assignments improving reconstruction quality while ensuring planning efficiency. Our method demonstrates the highest reconstruction quality among all planning methods and superior planning efficiency compared to existing multi-robot methods. We deploy our method on multiple robots, and results show that it can effectively plan view paths and reconstruct scenes with high quality.

Index Terms—3D Gaussian splatting; Open-vocabulary semantic segmentation; Multi-robot collaboration.

# I. INTRODUCTION

RECONSTRUCTION of indoor scenes is essential for numerous applications, such as gaming, robotics, augmented and virtual reality [3], [14], [7]. Recently, implicit

neural representations have demonstrated significant potential in autonomous systems due to their precise 3D reconstruction quality [17], [24]. To achieve higher rendering speed and quality, GS-Planner [8] proposes the first active 3D reconstruction system using 3DGS with online evaluation. Despite the impressive results of these works, their primary limitation is the inherent design to operate with only one single robot, significantly reducing scanning efficiency in large indoor environments.

Multi-robot autonomous 3D reconstruction provides advantages over single-robot systems, including broader coverage and enhanced efficiency [4], [6]. To address the task assignment problem in multi-robot collaboration, a Multi-depot multi-traveling salesman problem (MDMTSP) needs to be solved. Considering that MDMTSP is NP-hard and lacks an efficient exact solution, some improved methods, such as cluster-and-assign [4] and greedy cluster assignment [6], have been employed to accelerate the process of robot mode and task assignment. Although these methods achieve relatively favorable results in mode and task assignment, they still require repeated iterations to approximate optimal assignment. When the number of tasks increases, repeated computations of TSP entail significant computational costs, reducing planning efficiency.

However, during the generation of multiple tasks in active reconstruction methods [11], [1], [6], either the coverage of frontiers [6] or the coverage of specified surface uncertainty regions [11], [1] is prioritized, leading to objects not being reconstructed with high quality due to rapid scanning without specific attention. Specifically, when defined or learned uncertainties contain noise, excessive focus on high uncertainty areas can cause the robot to move back and forth in large

<sup>&</sup>lt;sup>1</sup> College of Control Science and Engineering, Zhejiang University, Hangzhou, 310027, China.

<sup>&</sup>lt;sup>2</sup> College of Information Engineering, Zhejiang University of Technology, Hangzhou, 310023, China.

<sup>\*</sup> Qi Ye (Corresponding author, qi.ye@zju.edu.cn) is with the College of Control Science and Engineering, the State Key Laboratory of Industrial Control Technology, Zhejiang University, and the Key Lab of CS&AUS of Zhejiang Province.

scenes, potentially getting trapped in local optima. To address this problem, semantic methods [12], [4] are adopted to make viewpoint sampling more focused. Liu et al. [12] proposed a model-driven objectness metric to evaluate the similarity and completeness of segmented components extracted from objects in a 3D model database. Moreover, Guo et al. [4] utilized incompleteness scores of segmented objects derived from point cloud completion to inform and optimize task generation. Nevertheless, building a 3D model database [12] or performing point cloud completion [4] requires significant human and time resources.

To handle the first challenge, we propose a novel centralized framework for multi-robot autonomous 3D reconstruction using 3DGS. To address the second challenge, we propose an efficient multi-robot collaboration strategy incorporating global-local tasks for robot mode assignment and an improved K-means algorithm for task assignment. By separating mode assignment from task assignment and focusing task clustering solely on Euclidean distance instead of TSP distance, this method eliminates repetitive TSP calculations typically required for iterative assignment in MDMTSP. To tackle the third challenge, we integrate online open-vocabulary semantic segmentation with 3DGS surface uncertainty, prioritizing view sampling in regions with high instance uncertainty to reduce redundant view generation, which eliminates the need for significant effort in pre-establishing a model database and avoids the time cost associated with point cloud completion.

To summarize, our contributions are:

- We propose the first centralized multi-robot autonomous 3D reconstruction framework utilizing 3DGS, which comprises three modules: Perception, Task generation, and Hierarchical planning.
- We incorporate online open-vocabulary semantic segmentation and surface uncertainty of 3DGS, focusing view sampling on areas with high instance uncertainty, which reduces the generation of redundant viewpoints, thereby reducing time for task generation.
- We propose an efficient multi-robot collaboration strategy: global-local tasks for mode assignment of robots and improved K-means for task assignment, improving reconstruction quality while ensuring planning efficiency.

## II. RELATED WORK

## A. Autonomous 3D reconstruction

Neural radiance field (NeRF) has become a highly effective method for 3D scene reconstruction due to its remarkable ability to render photorealistic images [15] and represent scene geometry [21] effectively. However, achieving high visual quality still necessitates using neural networks, which are expensive to train and render. More recently, 3DGS [9] has demonstrated comparable or superior rendering performance to NeRF [15], achieving faster rendering and optimization speeds with an order of magnitude.

View path planning is essential in autonomous reconstruction, focusing on optimizing the sequence of viewpoints to reconstruct a 3D scene efficiently. Ran et al. [17] propose the first autonomous 3D reconstruction system using an implicit neural representation. To enhance global coverage capabilities in complex environments while avoiding local minima, Zeng et al. [24] incorporate frontier-based exploration tasks with surface-based reconstruction tasks to improve the efficacy of reconstruction. Jin et al. [8] leverages the advantageous features of 3DGS to incorporate real-time quality and completeness assessment to guide the robot's reconstruction process.

### B. Multiple-robot collaboration

Recently, multi-task assignment methods [11], [1], [6] have been utilized in multi-robot active reconstruction to ensure efficient global coverage, precise scene reconstruction, and load balancing among multiple robots. Lauria et al. [11] considers a similar matroid-constrained submodular maximization problem for multisensor NBV planning. Dong et al. [1] formulate task assignments based on Optimal mass transport (OMT) and propose efficient solutions based on a divide-and-conquer scheme. Hardouin et al. [6] introduce multiagent NBV planners to route robots to viewpoint configurations.

# C. Semantic guidance in planning

Due to the focused nature of semantic information and corresponding segmentation attributes for each object in the scene, some methods [12], [25], [4] have attempted to use semantic information to guide robots in scanning the scene. Liu et al. [12] proposes a model-driven objectness metric to measure the similarity and completeness of segmented components from objects in the 3D model database. Zheng et al. [25]estimate the next best view based on the uncertainty in scene reconstruction and understanding. Furthermore, Guo et al. [4] introduces semantic information into the multi-robot system, guiding task generation through incompleteness score based on point cloud completion.

## III. METHOD

# A. Problem Statement and System Overview

In this paper, we address the challenge of multi-robot autonomous scene reconstruction with 3DGS utilizing a centralized framework. This study aims to achieve efficient exploration and scanning of an unknown indoor scene. By starting from the initial positions of multiple robots, the goal is to maximize scanning coverage and reconstruction quality while minimizing scanning effort [1], [4].

Under the framework of the centralized strategy, our pipeline comprises three components, as illustrated in Fig. 2. The method overview is shown in Fig. 3. **The Perception module** receives images rendered from the Unity at specified viewpoints similar to [17], [23] and maintains four representations (Section III-B). A volumetric representation (occupancy grid map) is adopted for exploration tasks. 3DGS [22] and Open-vocabulary 3D instance retrieval (OVIR-3D) [13] are adopted for reconstruction tasks. A modern framework for Parallel spatial hashing (ASH) [2] is used for scene geometry acquisition and noise point pruning of 3DGS. In **the Task generation module** (Section III-C), surface uncertainty is obtained through 3DGS, and reconstruction instances are acquired through OVIR-3D segmentation. These processes allow

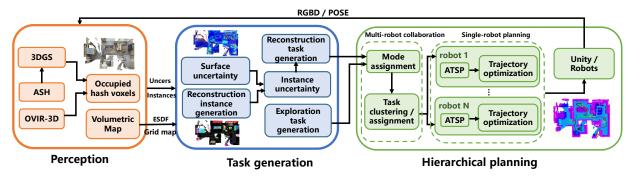


Fig. 2: The pipeline of our proposed method.

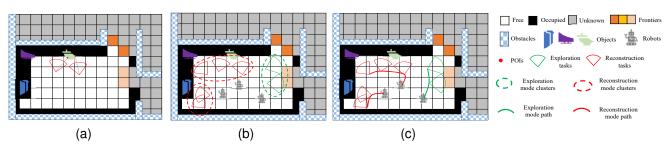


Fig. 3: Overview of our method. (a) Semantic-guided reconstruction task generation. (b) Multi-robot collaboration for mode and task assignments. (c) Single-robot view path planning.

for determining uncertainty for each point in every segmented instance, thereby guiding the generation of reconstruction tasks. Exploration tasks are generated from frontiers extracted from the volumetric map [26]. In **the Hierarchical planning Module** (Section III-D), Firstly, we determine the number of robots assigned to exploration and reconstruction modes based on the ratio between global whole exploration and reconstruction tasks. Each robot's mode is determined based on the number of tasks within its local range. Subsequently, for each robot, a global viewpoint path is planned using the Asymmetric traveling salesman problem (ATSP) solver, and B-spline trajectory optimization [19] is used to refine the path for smooth navigation. The images captured along this path are then processed by the perception system until the entire autonomous reconstruction process is completed.

# B. Perception

In this section, we first introduce the 3DGS representation and design a surface uncertainty retrieval method to evaluate its reconstruction quality. Then, we introduce an ASH) [2] method to obtain scene geometry and prune noisy Gaussians. Finally, we employ open-vocabulary online 3D semantic segmentation to focus view sampling on instances, thereby reducing the generation of redundant viewpoints.

1) 3D Gaussian splatting with surface uncertainty: In this study, we use Gaussian-SLAM [22] with online sub-map optimization as the base framework for our research. However, in our autonomous reconstruction system, the merging and quality evaluation of sub-maps require considerable computation time. Therefore, we maintain only a global 3DGS map during training, accelerating the quality evaluation process.

Our goal is to optimize a scene representation that captures the appearance of the scene, leading to a detailed, dense map and high-quality novel view synthesis. To achieve this, we represent the scene as a set of 3D Gaussians  $G^g$  as follows:

$$\mathbf{G}^g = \{ G_i^g : (\mu_i, SC_i, R_i, o_i, SH_i) \mid i = 1, \dots, N \}$$
 (1)

where each Gaussian  $G_i^g$  is defined by a mean  $\mu_i \in R^3$ , scale  $SC_i \in R^3$ , rotation  $R_i \in R^4$ , opacity  $o_i \in R$  and spherical harmonics  $SH_i \in R^{48}$ .

To evaluate the reconstruction quality of 3DGS, we adopt a loss caching strategy similar to GS-Planner [8]. This process requires projecting the loss  $L_{proj}$  from the image space to the world space and caching the loss into occupied voxels. Specifically,  $L_{proj}$  is a weighted sum of color loss  $L_{color}$  and depth loss  $L_{depth}$ :

$$L_{proj} = L_{color} + \lambda_d \cdot L_{depth} \tag{2}$$

where  $\lambda_d=0.5$  is the weight coefficient. To reduce computational overhead and memory usage, we maintain occupied hash voxels  $H_{occ}$  with resolution 0.05 m to store the occupied voxels, which are updated by adding new images. The losses will be averaged for each projected voxel v if it is already in  $H_{occ}$ . Otherwise, the new voxel with the position of projecting point p and associated loss  $L_p$  is added to  $H_{occ}$ . We then obtain uniform surface points S from  $H_{occ}$  and and surface uncertainty  $U_{qs}$  with uncertainty  $L_p$  of each point  $p \in S$ :

2) ASH for scene geometry and pruning noisy Gaussians: It is widely recognized that 3DGS has achieved impressive results in novel view synthesis. However, for complex scenes, such as multi-room indoor environments, the geometric representation capability of 3DGS is still limited. Therefore, we introduce an ASH [2] method to represent scene geometry. ASH implements spatially varying operations, transitioning seamlessly from volumetric geometry reconstruction to differentiable appearance reconstruction. ASH compensates for the limitations of 3DGS and can also prune noisy Gaussians

generated during the training process of 3DGS. Specifically, ASH extracts geometric surface points at intervals of 30 frames and generates a KD-Tree. For each Gaussian in 3DGS, the nearest point in surface points of ASH is found, and if the distance is greater than 0.05 m, the Gaussian is considered noisy and is pruned.

3) Open-vocabulary online 3D semantic segmentation: As the goal is to obtain high-quality reconstruction of 3D scenes, more scans should be conducted around surface points with high uncertainty, similar to AIISRFE [24]. However, due to the scattered distribution of surface uncertainty  $U_{gs}$ , excessive redundant views are generated in regions with high uncertainty points, thus reducing the scanning efficiency.

OVIR-3D [13] aims to return a set of 3D instance segments given a 3D point cloud and the corresponding RGBD images with poses. This is achieved by a multi-view fusion of textaligned 2D region proposals [27], [10] into 3D space. To obtain the aforementioned 3D point cloud online and align it with the surface uncertainty of 3DGS, we employ surface points S from the occupied voxel  $H_{occ}$  as the online point cloud input for OVIR-3D segmentation. The 2D region proposals of each frame are then projected to the surface points S given the camera pose. The projected 3D regions are either matched to existing 3D object instances  $O = \{O_1, \dots, O_b\}$ with 3D features  $F^{3D} = \{f_1^{3D}, \dots, f_b^{3D}\}$ , or added as a new instance if not matched with anything. Instance points of all objects are  $S_O = \{S_1, \dots, S_b\}, S_O \in S$ , where b is the number of segmented 3D instances. Viewpoints are generated for incomplete object instances to overcome the local minima caused by the scattered distribution of surface uncertainty.

**Algorithm 1:** Semantic-guided Reconstruction Task Generation

```
 \begin{array}{c} \textbf{Input} : \text{Downsampling ratio } N_{down}, \text{ POI distance threshold} \\ d_{POI}, \text{ 3DGS surface uncertainty } U_{gs}, \text{ segmented instances} \\ O \text{ with features } F^{3D}, \text{ Scannet200 vocabulary } F^{vocab}, \\ \text{points } S_O, \text{ volumetric map } V \\ \textbf{Output: Updated reconstruction tasks } \mathcal{T}^{rec} \\ \textbf{1} \quad O^{rec} \leftarrow RecInstanceGenerating(O, S_O, F^{3D}, F^{vocab}); \\ \textbf{2} \quad U^{rec} \leftarrow InstanceUncertaintyRetrieving(U_{gs}, O^{rec}); \\ \textbf{3} \quad \textbf{foreach } O_i^{rec} \in O^{rec} \quad \textbf{do} \\ & | // \text{ Process each instance} \\ \textbf{4} \quad & | S_i^d, U_i^d \leftarrow SurfaceDownsampling(S_i, U_i, N_{down}); \\ \textbf{5} \quad & | P_i^d \leftarrow POIRetrieving(S_i^d, U_i^d, d_{POI}); \\ \textbf{6} \quad & | V_i^{rec} \leftarrow Viewsampling(V, S_i^d, U_i^d, P_i^d); \\ \textbf{7} \quad & | v_i^{rec} \leftarrow ViewsofMaxGain(VP_i^{rec}); \\ \textbf{8} \quad & | \mathcal{T}_i^{rec} \leftarrow v_i^{rec}; \\ \textbf{9} \quad \textbf{end} \\ \textbf{10} \quad & \mathcal{T}^{rec} \leftarrow \{\mathcal{T}_1^{rec}, \mathcal{T}_2^{rec}, \dots, \mathcal{T}_{N_{rec}}^{rec}\}; \\ \end{array}
```

# C. Task generation

Within our approach, scanning tasks are categorized into two types: exploration tasks aimed at achieving rapid coverage and reconstruction tasks focused on ensuring high-quality reconstruction.

1) Exploration task generation: The exploration tasks  $\mathcal{T}^{exp} = \{\mathcal{T}_1^{exp}, \mathcal{T}_2^{exp}, ..., \mathcal{T}_{N^{exp}}^{exp}\}$  are intended to cover more unknown regions, with  $N^{exp}$  denoting the total number of tasks. We firstly update incremental frontiers and Euclidean signed distance field (ESDF) map E [5] from the maintained volumetric map  $V = V_o \cup V_e \cup V_u$  similar to Fuel [26], where  $V_o, V_e, V_u$  represent occupied, empty and unknown voxels.

Subsequently, we select viewpoints that can provide superior coverage of the frontiers in exploration tasks, similar to our previous method AIISRFE [24].

2) Semantic-guided reconstruction task generation: The reconstruction tasks  $\mathcal{T}^{rec} = \{\mathcal{T}^{rec}_1, \mathcal{T}^{rec}_2, \dots, \mathcal{T}^{rec}_{N^{rec}}\}$  are designed to refine areas with low reconstruction quality, and  $N^{rec}$  denotes the total number of tasks. Generally, objects with complex textures and structures are more difficult to reconstruct, necessitating focused attention during scanning. Therefore, reconstruction tasks can be transformed into scanning the low-quality regions of objects with low completeness. Algorithm 1 describes the generation process of reconstruction tasks.

Reconstruction instance generation To evaluate the objectness scores [12] of 3D object instances O, which represents the probability that an instance belongs to a certain category and reflects the completeness of the object, we need a vocabulary to match each instance with its corresponding category. We choose ScanNet200 vocabulary [18] since it contains most indoor object categories. The similarity scores between the feature  $f_i^{3D}$  of each instance  $O_i \in O$  and ScanNet200 vocabulary features  $F^{vocab}$  are computed and sorted in descending order as  $\Lambda_i = \{\Lambda_i^1, \dots, \Lambda_i^{200}\}$ . To enhance inter-class separability, a scaling factor  $\lambda_e = 50$  is applied to these scores. Finally, for each instance  $O_i \in O$ , the classification probabilities  $P_i$  for the 200 categories are obtained by normalizing the scores using a softmax function as  $C_i = \{C_i^1, \dots, C_i^{200}\} = softmax(\lambda_e \Lambda_i)$ . We select the category with the highest probability score as the label  $L_i$  for instance  $O_i$ . The objectness score of this instance is defined as  $C_i^1$ .

If  $C_i^1 < C_{min}$ , it is considered that the objectness score calculation of the instance may be affected by noise from CLIP labels. If  $C_i^1 > C_{max}$ , it is deemed that the instance has already been well reconstructed. We exclude these two types of instances and denote the remaining instances as reconstruction instances  $O^{rec}$ :

$$O^{rec} = \{O_1^{rec}, \dots, O_m^{rec}\} \tag{3}$$

where  $O^{rec}$  represents the instances used for generating reconstruction tasks, m is the number of reconstruction instances. For each instance  $O^{rec}_i = \{L^{rec}_i, C^{1,rec}_i, S^{rec}_i\} \in O^{rec}$ , it includes the label  $L^{rec}_i$ , objectness score  $C^{1,rec}_i$ , and surface points  $S^{rec}_i$ , where  $S^{rec}_i \in S_O$ .

Instance uncertainty For each instance  $O_i^{rec} \in O^{rec}$  in (3), we need to evaluate instance uncertainty to generate viewpoints that cover its low-quality areas. We query uncertainty of surface points  $S_i^{rec}$  in instance  $O_i^{rec}$  from surface uncertainty  $U_{gs}$  and denotes it as instance uncertainty  $U_i^{rec}$ . The instance uncertainty for each instance within the reconstruction instances  $O^{rec}$  can be obtained as follows

$$U^{rec} = \{U_1^{rec}, \dots, U_m^{rec}\} \tag{4}$$

**View sampling** For each instance  $O_i^{rec} \in O^{rec}$ , we downsample its surface points by a factor of  $N_{down} = 5$ , obtaining the downsampled surface points  $S_i^d$  and uncertainty  $U_i^d$ . To generate task viewpoints, we need to identify the points of interest (POIs) for each instance. Specifically, we initially select the point with the highest uncertainty. Next, from

the remaining points, we choose iteratively the point with the highest uncertainty that is at least  $d_{POI}$  away from all previously selected points. This process continues until no more points can be added. The set of selected points is denoted as  $P_i^d = \{P_{i,1}^d, \dots, P_{i,k_i}^d\}$ , where  $k_i$  is the number of sampled POIs from instance  $O_i^{rec}$ .

For each sampled POI  $P_{i,j}^d \in P_i^d$ , we set the center as  $P_{i,j}^d$  and generate a series of candidate viewpoints  $VP_{i,j}^d$  oriented towards the center within the empty space Ve, similar to AIISRFE [24]. This yields all candidate viewpoints  $VP_i^d = \{VP_{i,1}^d, \dots, VP_{i,k_i}^d\}$  for the instance  $O_i^{rec}$ .

Instance uncertainty based Information gain To select reconstruction tasks from these candidate viewpoints, we define information gain of viewpoint v as:

$$g(v) = \sum_{k=1}^{N_{vis}} \sigma_k e^{-0.5d_{v,k}}$$
 (5)

where  $N_{vis}$  represents the number of visible instance surface points,  $\sigma_k \in U_i^d$  represents the uncertainty of each visible surface point  $s_k \in S_i^d$ , and  $d_{v,k}$  is the distance from the viewpoint v to the surface point  $s_k$ .

We then choose the viewpoint  $v_{i,j}^d$  with the highest information gain as the reconstruction task for the sampled POI  $P_{i,j}^d \in P_i^d$  in instance  $O_i^{rec}$ . Reconstruction tasks  $\mathcal{T}_i^{rec} = v_i^{rec}$  for instance  $O_i^{rec}$  are generated for all sampled POIs from instance  $O_i^{rec}$ . This allows us to obtain the final reconstruction tasks  $\mathcal{T}^{rec} = \{\mathcal{T}_1^{rec}, \mathcal{T}_2^{rec}, \dots, \mathcal{T}_{N_{rec}}^{rec}\}$ .

#### D. Hierarchical planning

Once the new exploration tasks  $\mathcal{T}^{exp}$  and reconstruction tasks  $\mathcal{T}^{rec}$  are generated, they need to be assigned to the robots for execution. To find the best task assignment, we first construct a weighted graph  $G^{\mathcal{T}} = (\mathcal{T}^{exp} \cup \mathcal{T}^{rec}, \mathcal{E})$  to represent the spatial relationships between the tasks and robots similar to ACAMS [4], where  $\mathcal{E}$  consists of edges connecting each pair of tasks with travel cost. We optimize the travel costs during task assignment and formulate it as an MDMTSP problem. The objective is to find a set of disjoint paths  $\{\mathcal{T}^*_r\}_{r=1}^{N_R}$  that covers the entire set  $G^{\mathcal{T}}$  such that the total tour costs for all robots are minimized:

$$E_d = \sum_{r=1}^{N_R} \sum_{\mathcal{T}_k \in \mathcal{T}^*} t_{lb}(\mathcal{T}_k, \mathcal{T}_{k+1})$$
 (6)

where  $N_R$  is the number of robots,  $t_{lb}(\mathcal{T}_k, \mathcal{T}_{k+1})$  is the travel cost with the time lower bound from task  $\mathcal{T}_k$  to  $\mathcal{T}_{k+1}$  similar to Fuel [26], in which we consider path length cost, yaw cost and pitch cost. Moreover, directly merging exploration and reconstruction tasks is susceptible to falling into local minima, as demonstrated in [24], making it essential to assign modes for the robots.

To solve the complex NP-hard MDMTSP problem as denoted in (6), we adopt a hierarchical planning strategy. In multi-robot collaboration, the robots are assigned to different modes and tasks to reduce the number of robot and task assignments in the MDMTSP problem. In single-robot planning, after the clusters with tasks are assigned to robots, the task execution order is determined for each robot.

1) Multi-robot collaboration: In this section, we first determine the mode of each robot based on the global-local task distribution to reduce the number of robots and tasks in the MDMTSP problem. Subsequently, we employ an improved K-means clustering algorithm to assign tasks to the robots to reduce the number of task assignments within the problem.

**Mode assignment** The mode of the robot should be determined based on the distribution of nearby tasks due to the execution path cost. Besides, the number of robots assigned to exploration and reconstruction modes should also take into account the global task distribution, as demonstrated in [4]. To this end, we designed a global-local task distribution-guided method for robot mode assignment.

Firstly, we compute the proportion of two modes of all global tasks and assign the number of robots for each scanning mode according to this proportion. Secondly, we improve the assignment efficiency through a local task count statistic to avoid the iterative optimization process in ACAMS [4]. Specifically, for each robot  $R_r$ , we count the number of exploration and reconstruction tasks  $N_{local,r}^{exp}$ ,  $N_{local,r}^{rec}$ , within its local range not exceeding  $d_{local}$ . The overall score is denoted as  $N_{local,r}^{ove} = N_{local,r}^{exp} - N_{local,r}^{rec}$ . We then rank the robots based on their overall scores, with higher scores indicating a priority for exploration mode. This prioritization means that robots with more exploration tasks in their local area will be assigned to further exploration. Consequently, the modes for all robots are determined, considering the number of robots in each mode.

Task clustering and assignment After assign different modes to robots, we need to assign tasks to robots within each mode. Although the number of robots and tasks is reduced through mode assignment, it is still an MDMTSP problem. To further reduce the complexity of solving this problem, we propose an improved K-means algorithm for task clustering and assignment.

For robots in task mode  $M \in \{exp, rec\}$ , we cluster tasks  $\mathcal{T}^M$  into  $N_R^M$  classes using an improved K-means method and design the following objective function for task clustering:

$$\min_{\mathcal{T}^M} \sum_{r=1}^{N_R^M} \underbrace{D_r + \gamma(R_r^M, \omega_r)}_{\text{moving cost}} + \underbrace{\left(N_r^M - \overline{N}\right)^2 + |D_r - \overline{D}|}_{\text{robot capacity}}$$
(7)

where  $D_r = \sum_{\mathcal{T}_k \in \mathcal{T}_r^M} \gamma(\mathcal{T}_k, \omega_r)$  is the distance sum from assigned tasks to the task centroid  $\omega_r$  for robot  $R_r^M$  and  $\overline{D} = \sum_{D_r}/N_R^M$  is the average distance of  $N_R^M$  robots.  $\mathcal{T}_r^M \in \mathcal{T}^M$  is the tasks assigned to robot  $R_r^M$ .  $\gamma(R_r^M, \omega_r)$  is the distance from robot  $R_r^M$  to the task centroid  $\omega_r$ .  $N_r^M$  is the number of tasks assigned to robot  $R_r^M$  and  $\overline{N} = N^M/N_R^M$  is average number of tasks for  $N_R^M$  robots. Finally, we obtain the clusters of tasks assigned to all  $N_R$  robots, denoted as  $\{\mathcal{T}_r\}_{r=1}^{N_R}$ .

As above in (7), the task clustering and assignment are determined by the following two aspects:

Movement cost. Similar to [1], We approximate the moving cost by considering two components: 1) the Euclidean distance from the tasks to the centroid of the set of assigned tasks, and 2) the Euclidean distance from the robot to the centroid of the set of assigned tasks.

Robot capacity. Different from [1], [4], to ensure load balancing constraints on different robots, we consider not only

TABLE I: Evaluations of the effectiveness and efficiency of view path for 3DGS and ASH representations.

		Method	1			Qa	ıLdn			N	icut		Oyens			
Variant	Pitch	$\mathcal{T}^{rec}$	MA	IKM	PSNR↑	$Acc\!\!\downarrow$	$Comp \!\!\downarrow$	Recall↑	PSNR↑	$Acc\!\!\downarrow$	$Comp \!\!\downarrow$	Recall↑	PSNR↑	Acc↓	$Comp \!\!\downarrow$	Recall↑
V1(Fuel swarm <sup>[26]</sup> )					19.30	10.15	2.89	0.87	18.10	3.06	2.51	0.90	18.71	3.54	4.65	0.81
V2	✓				21.78	10.25	2.87	0.87	22.12	3.02	2.15	0.92	23.25	3.53	3.62	0.84
V3(AIISRFE swarm <sup>[24]</sup> )	✓	Surface			19.09	10.64	11.10	0.74	18.00	3.03	10.71	0.78	21.85	3.56	14.79	0.73
V4	✓	Semantics			20.33	10.35	5.38	0.79	21.51	3.03	3.43	0.89	22.56	3.52	10.24	0.75
V5	✓	Semantics	$\checkmark$		21.89	10.25	2.41	0.89	22.84	3.05	1.91	0.94	23.35	3.54	4.17	0.83
V6(Ours full)	✓	Semantics	$\checkmark$	$\checkmark$	22.75	10.04	2.25	0.91	24.28	3.01	1.62	0.95	24.81	3.50	3.36	0.86
Variant	Pitch	$\mathcal{T}^{rec}$	MA	IKM	$T_{task}$	$T_{colla}$	$T_{GP}$	P.L.	$T_{task}$	$T_{colla}$	$T_{GP}$	P.L.	T <sub>task</sub>	$T_{colla}$	$T_{GP}$	P.L.
V1(Fuel swarm <sup>[26]</sup> )					0.003	0.026	0.56	107.61	0.005	0.007	0.36	89.80	0.005	0.007	0.38	84.83
V2	<b>✓</b>				0.004	0.034	1.00	121.44	0.003	0.015	1.51	103.23	0.003	0.014	1.09	90.97
V3(AIISRFE swarm <sup>[24]</sup> )	✓	Surface			2.735	0.110	18.77	98.10	2.232	0.103	18.53	93.39	2.785	0.550	18.44	78.68
V4	<b>✓</b>	Semantics			0.799	0.353	10.15	97.19	0.477	0.109	4.12	107.95	0.432	0.231	5.03	108.60
V5	✓	Semantics	$\checkmark$		0.701	0.241	9.16	144.49	0.421	0.147	4.10	129.46	0.441	0.263	5.56	113.33
V6(Ours full)	✓	Semantics	$\checkmark$	$\checkmark$	0.718	0.546	10.92	167.77	0.471	0.253	4.91	131.18	0.437	0.280	5.85	126.46

TABLE II: Evaluations of the effectiveness and efficiency with existing planning methods.

	QaLdn							Nicut							Oyens					
Method	PSNR↑	$Acc\!\downarrow$	Comp↓	Recall↑	$T_{GP}$	P.L.	PSNR↑	$Acc\!\!\downarrow$	$Comp \!\!\downarrow$	Recall↑	$T_{GP}$	P.L.	PSNR↑	$Acc\!\downarrow$	$Comp \!\!\downarrow$	$Recall \uparrow$	$T_{GP}$	P.L.		
ACAMS <sup>[4]</sup>	22.36	10.14	2.83	0.88	87.10	175.94	22.96	3.07	1.88	0.94	41.62	170.42	23.36	3.53	3.53	0.85	65.53	146.45		
MS3DSR <sup>[6]</sup>	21.05	10.19	2.96	0.87	24.87	121.90	20.75	3.09	2.46	0.91	13.62	120.16	19.80	3.53	25.67	0.65	8.78	125.51		
Ours	22.75	10.04	2.25	0.91	10.92	167.77	24.28	3.01	1.62	0.95	4.91	131.18	24.81	3.50	3.36	0.86	5.85	126.46		

the balance of the number of tasks assigned to each robot but also the balance of the internal distances among the tasks assigned to each robot.

2) Sing-robot view path planning: In this section, we plan the view path for each robot with assigned tasks. For robot  $R_r$  with assigned tasks  $\mathcal{T}_r$ , we plan a global shortest path based on ATSP similar to [24] as follows:

$$\mathcal{T}_r^* = \arg\min \sum_{\mathcal{T}_k \in \mathcal{T}_r} t_{lb}(\mathcal{T}_k, \mathcal{T}_{k+1})$$
 (8)

where  $\mathcal{T}_r^*$  is planed path of robot  $R_r$ . For smooth navigation of the robot, we adopt B-spline optimization to generate smooth, safe, and dynamically feasible B-spline trajectories, similar to Fuel [26]. Viewpoints  $V_r^s$  are sampled along the trajectory at time intervals of t=0.4s, ensuring that the execution path length does not exceed  $L_{exec}$ . The sampled viewpoints for all robots are executed for map updates until a certain number of images are collected.

#### IV. RESULTS

## A. Implementation details

- 1) Data: The experiments are conducted on three virtual scenes:  $QaLdn~(187m^2)$  from HM3D [16],  $Nicut~(120m^2)$  and  $Oyens~(80m^2)$  from Gibson [20]. These scenes are reconstructed from real-world environments. We maintain the same field of view parameters, as specified in prior works such as [23], [24].
- 2) Implementation: Our method runs on two GPUs. The 3DGS reconstruction and surface uncertainty updates are performed on one A6000 GPU; OVIR-3D and view path planning are on the other RTX3090 GPU, where view path planning runs in the ROS environment. We set the maximum planned views to be 600 views for QaLdn, 600 views for Nicut, and 500 views for Oyens. In reconstruction instance generation, objectness score thresholds are  $C_{min}=0.2$ ,  $C_{max}=0.6$ . In view sampling, the POI distance threshold is  $d_{POI}=1.2m$ .

In mode assignment, the local range is  $d_{local} = 6m$ . In singrobot planning, the execution path length is  $L_{exec} = 6m$ . The experiments are conducted using four robots.

3) Metric: Similar to [23], [24], we evaluate our method from two aspects including effectiveness and efficiency. The effectiveness is measured in two parts: the quality of the rendered images and the quality of the geometry of the reconstructed surface. We adopt metrics: Accuracy (cm), Completion (cm), Recall. The number of sampled points and completion distance threshold are the same as AIISRFE [24].

For efficiency, we evaluate the path length (meter) and the planning time (second). The total path length is denoted as P.L. and the planning time as  $T_{GP}$ . For a more detailed comparison, we further divide the time required for view path planning at each step of the reconstruction process into three components: 1) the task generation time  $T_{task}$ , 2) the time  $T_{colla}$  for the multi-robot collaboration, 3) the time  $T_{sr}$  for single-robot planning.  $T_{sr}$  for all the variants in Table I is about 0.04s to 0.4s.

# B. Efficacy of the Method

Similar to [24], the efficacy of the method is evaluated regarding both the effectiveness and efficiency of our contributions. We design variants V1-V5 of our method based on 3DGS representation. V1 (Fuel swarm [26]) considers only exploration tasks and extends Fuel to the multi-robot system. V2 adds a degree of pitch from V1. V3 (AIISRFE swarm [24]) considers exploration tasks and surface uncertainty-based reconstruction tasks and extends AIISRFE to the multi-robot system. V4 merges exploration and semantic-guided reconstruction tasks for assignment. V5 considers mode assignments for robots but uses moving cost and robot capacity in [1] for task clustering and assignment. Our method is best for effectiveness and better than the method with surface-based reconstruction tasks for efficiency.

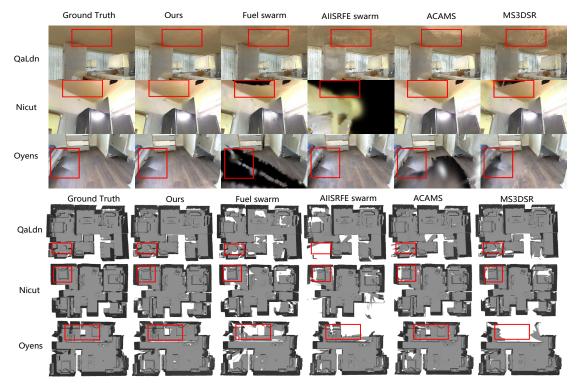


Fig. 4: Comparison with different methods. Top: novel view synthesis from 3DGS; Bottom: reconstructed meshes from ASH.



Fig. 5: Comparison of trajectories with different methods.

1) Combination of exploration and semantic-guided reconstruction tasks: We make V1 and V2 as our baselines to verify the efficacy. To verify the efficacy of different reconstruction tasks, we make V2, V3, and V4 as our baselines.

The metrics for V2 in Table I show that the pitch angle enhances reconstruction quality compared to Fuel swarm, with only a minor increase in planning time. Concentrating exclusively on exploration tasks (V2) neglects the detailed scanning of intricate features, which in turn diminishes the reconstruction quality. When exploration and surface-based reconstruction tasks (V3) are accounted for, the continuous updating of the KD-tree during task generation, coupled with the generation of numerous redundant viewpoints due to surface uncertainty error, further diminishes planning efficiency. The reconstruction quality becomes very poor because it falls into the local optimum and cannot cover the entire scene. V4 reduces task generation time due to efficiency of reconstruction instance generation.

2) Multi-robot collaboration: To validate the effectiveness of mode assignment (MA) and improved K-means (IKM) based task assignment of robots, we use V4 and V5 as our baselines. The merger of these exploration and reconstruction tasks (V4) significantly slows the pace of scene exploration and can even result in getting trapped in a local optimum,

thereby reducing the reconstruction quality. V5 does not consider balancing internal distances in (7), which can result in long execution paths for some robots. Due to the limited energy of robots, path truncation is required for execution, which reduces exploration efficiency and, consequently, the quality of the reconstruction. (Ours) ensures the speed of exploration and detailed scanning while avoiding the local optima often encountered when scanning areas with high instance uncertainty.

## C. Comparison with existing planning methods

We select two recent works on multi-robot autonomous reconstruction: ACAMS [4] and MS3DSR [6], as our baselines. ACAMS proposes a modified MDMTSP and corresponding approximate solver to optimize each robot's task assignment and execution order. MS3DSR generates viewpoints-based surface elements and introduces multiagent NBV planners to route robots to viewpoint configurations. The metrics presented in Table II indicate that our method excels compared to these baselines, showing superior reconstruction quality and planning efficiency. This is because the two aforementioned methods require continuous computation of the TSP when assigning modes and tasks, which increases time costs. Additionally, MS3DSR based on surface elements is susceptible

to environmental factors such as holes, making it challenging to achieve rapid scene coverage and thereby reducing reconstruction quality.

Fig. 4 shows our method provides better reconstruction results in novel views and geometry. For more visual comparisons and results, we refer readers to the supplementary video. Fig. 5 demonstrates that the trajectory of our method expands in scene QaLdn that of other methods.

## D. Robot experiments in real scene

We implement our proposed method on two Turtlebots equipped with an Azure Kinect camera $^1$  and a Realsense T265 camera $^2$  to perform room reconstruction, specifically targeting an exhibition room with dimensions of  $12m \times 8m \times 4m$ . The pose of the Kinect camera is provided by the T265 camera. For this scene, turtlebots take about 8 minutes to explore and reconstruct the room. The exploration and reconstruction results will be presented in the supplementary video.

### V. CONCLUSION

In this paper, we propose the first centralized multi-robot autonomous 3D reconstruction framework utilizing 3DGS. Subsequently, we incorporate online open-vocabulary semantic segmentation and surface uncertainty of 3DGS, focusing view sampling on areas with high instance uncertainty to reduce the time cost of task generation. Finally, we propose an efficient multi-robot collaboration strategy to improve reconstruction quality while ensuring planning efficiency. Comprehensive experiments demonstrate the superior performance of our method.

In the future, we plan to investigate the deployment of our research on the distributed system to alleviate the communication load and enhance system robustness.

## REFERENCES

- Siyan Dong, Kai Xu, Qiang Zhou, Andrea Tagliasacchi, Shiqing Xin, Matthias Nießner, and Baoquan Chen. Multi-robot collaborative dense scene reconstruction. ACM Transactions on Graphics, 38(4):1–16, 2019.
- [2] Wei Dong, Yixing Lao, Michael Kaess, and Vladlen Koltun. ASH: A modern framework for parallel spatial hashing in 3D perception. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5417–5435, 2022.
- [3] Qiao Gu, Alihusein Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs: Openvocabulary 3d scene graphs for perception and planning. arXiv preprint arXiv:2309.16650, 2023.
- [4] Junfu Guo, Changhao Li, Xi Xia, Ruizhen Hu, and Ligang Liu. Asynchronous collaborative autoscanning with mode switching for multirobot scene reconstruction. ACM Transactions on Graphics, 41(6):1–13, 2022.
- [5] Luxin Han, Fei Gao, Boyu Zhou, and Shaojie Shen. Fiesta: Fast incremental euclidean distance fields for online motion planning of aerial robots. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 4423–4430. IEEE, 2019.
- [6] Guillaume Hardouin, Julien Moras, Fabio Morbidi, Julien Marzat, and El Mustapha Mouaddib. A multirobot system for 3-d surface reconstruction with centralized and distributed architectures. *IEEE Transactions* on Robotics, 39(4):2623–2638, 2023.
- [7] Chenxing Jiang, Hanwen Zhang, Peize Liu, Zehuan Yu, Hui Cheng, Boyu Zhou, and Shaojie Shen. H \_{2}-mapping: Real-time dense mapping using hierarchical hybrid representation. *IEEE Robotics and Automation Letters*, 2023.
- <sup>1</sup>https://azure.microsoft.com/services/kinect-dk/
- <sup>2</sup>https://www.intelrealsense.com/tracking-camera-t265/

- [8] Rui Jin, Yuman Gao, Haojian Lu, and Fei Gao. Gs-planner: A gaussiansplatting-based planning framework for active high-fidelity reconstruction. arXiv preprint arXiv:2405.10142, 2024.
- [9] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. ACM Trans. Graph., 42(4):139–1, 2023.
- [10] 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.
- [11] Mikko Lauri, Joni Pajarinen, Jan Peters, and Simone Frintrop. Multisensor next-best-view planning as matroid-constrained submodular maximization. *IEEE Robotics and Automation Letters*, 5(4):5323–5330, 2020.
- [12] Ligang Liu, Xi Xia, Han Sun, Qi Shen, Juzhan Xu, Bin Chen, Hui Huang, and Kai Xu. Object-aware guidance for autonomous scene reconstruction. ACM Transactions on Graphics, 37(4):1–12, 2018.
- [13] 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 *Conference on Robot Learning*, pages 1610– 1620. PMLR, 2023.
- [14] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison.
  Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18039–18048, 2024.
  [15] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi,
- [15] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In European conference on computer vision, 2020.
- [16] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2), 2021.
- [17] Yunlong Ran, Jing Zeng, Shibo He, Jiming Chen, Lincheng Li, Yingfeng Chen, Gimhee Lee, and Qi Ye. Neurar: Neural uncertainty for autonomous 3d reconstruction with implicit neural representations. *IEEE Robotics and Automation Letters*, 8(2):1125–1132, 2023.
- [18] David Rozenberszki, Or Litany, and Angela Dai. Language-grounded indoor 3d semantic segmentation in the wild. In *European Conference* on Computer Vision, pages 125–141. Springer, 2022.
- [19] Jesus Tordesillas, Brett T Lopez, and Jonathan P How. Faster: Fast and safe trajectory planner for flights in unknown environments. In 2019 IEEE/RSJ international conference on intelligent robots and systems (IROS), pages 1934–1940. IEEE, 2019.
- [20] Fei Xia, Amir R. Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson Env. real-world perception for embodied agents. In Computer Vision and Pattern Recognition (CVPR), 2018 IEEE Conference on. IEEE, 2018.
- [21] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. Advances in neural information processing systems, 35:25018–25032, 2022.
- processing systems, 35:25018–25032, 2022.
  [22] Vladimir Yugay, Yue Li, Theo Gevers, and Martin R Oswald. Gaussian-slam: Photo-realistic dense slam with gaussian splatting. arXiv preprint arXiv:2312.10070, 2023.
- [23] Jing Zeng, Yanxu Li, Yunlong Ran, Shuo Li, Fei Gao, Lincheng Li, Shibo He, Jiming Chen, and Qi Ye. Efficient view path planning for autonomous implicit reconstruction. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 4063–4069. IEEE, 2023.
- [24] Jing Zeng, Yanxu Li, Jiahao Sun, Qi Ye, Yunlong Ran, and Jiming Chen. Autonomous implicit indoor scene reconstruction with frontier exploration. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 18041–18047. IEEE, 2024.
- [25] Lintao Zheng, Chenyang Zhu, Jiazhao Zhang, Hang Zhao, Hui Huang, Matthias Niessner, and Kai Xu. Active scene understanding via online semantic reconstruction. In *Computer Graphics Forum*, volume 38, pages 103–114. Wiley Online Library, 2019.
- [26] Boyu Zhou, Yichen Zhang, Xinyi Chen, and Shaojie Shen. Fuel: Fast UAV exploration using incremental frontier structure and hierarchical planning. *IEEE Robotics and Automation Letters*, 6(2):779–786, 2021.
  [27] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and
- [27] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In European Conference on Computer Vision, pages 350– 368. Springer, 2022.