

# OrionNav: Online Planning for Robot Autonomy with Context-Aware LLM and Open-Vocabulary Semantic Scene Graphs

Venkata Naren Devarakonda<sup>\*1</sup>, Raktim Gautam Goswami<sup>\*1</sup>, Ali Umut Kaypak<sup>\*1</sup>, Naman Patel<sup>\*1</sup>,  
Rooholla Khorrambakht<sup>1</sup>, Prashanth Krishnamurthy<sup>1</sup>, Farshad Khorrami<sup>1</sup>

**Abstract**—Enabling robots to autonomously navigate unknown, complex, dynamic environments and perform diverse tasks remains a fundamental challenge in developing robust autonomous physical agents. These agents must effectively perceive their surroundings while leveraging world knowledge for decision-making. Although recent approaches utilize vision-language and large language models for scene understanding and planning, they often rely on offline processing, offboard compute, make simplifying assumptions about the environment and perception, limiting real-world applicability. We present a novel framework for real-time onboard autonomous navigation in unknown environments that change over time by integrating multi-level abstraction in both perception and planning pipelines. Our system fuses data from multiple onboard sensors for localization and mapping and integrates it with open-vocabulary semantics to generate hierarchical scene graphs from continuously updated semantic object map. The LLM-based planner uses these graphs to create multi-step plans that guide low-level controllers in executing navigation tasks specified in natural language. The system’s real-time operation enables the LLM to adjust its plans based on updates to the scene graph and task execution status, ensuring continuous adaptation to new situations or when the current plan cannot accomplish the task, a key advantage over static or rule-based systems. We demonstrate our system’s efficacy on a quadruped navigating dynamic environments, showcasing its adaptability and robustness in diverse scenarios. Project video: <https://www.youtube.com/watch?v=gqWCSjLLG1E>.

## I. INTRODUCTION

Autonomous physical agents operating in real-world environments require robust and scalable scene representations capable of adapting to continuously evolving conditions, enabling them to perform complex tasks such as navigation, manipulation, and interaction [1]–[4]. Moreover, these agents must be capable of swiftly recognizing when a plan’s execution is unsuccessful and dynamically adjusting their strategies, ensuring resilience in the face of unexpected changes in the environment. This necessitates integrated approaches to both perception and planning, where agents must recognize and classify objects, construct structured models of their surroundings, and formulate multi-step plans to make informed decisions and execute tasks. A fundamental challenge lies in developing these representations and plans in real-time using

onboard sensors while ensuring adaptability to the dynamically changing nature of real-world settings.

Simultaneous Localization and Mapping (SLAM) has long been the cornerstone for enabling robots to build maps and localize themselves within them [5]–[7]. While traditional SLAM techniques primarily focus on geometric and spatial information, they lack the semantic understanding that is essential for higher-level reasoning and task execution. Object detection and segmentation pipelines have been previously integrated into these systems to bridge this gap, allowing them to identify and categorize objects within their environment while also refining the map and odometry [8]–[10]. Recent advancements in open-vocabulary detection and segmentation [11] are now being integrated into the mapping process, enhancing the robot’s environmental understanding by identifying and classifying objects within varying environments without requiring retraining [12]. However, these approaches often rely on pre-built semantic maps and struggle with real-time operation, limiting their applicability in dynamically changing, large-scale cluttered environments. Effective navigation and task performance in such settings demand scene representations organized hierarchically across multiple abstraction levels. Scene graphs have emerged as a powerful framework to meet this need, capturing objects and their interrelationships in a structured manner [13]–[16]. These graphs not only encode geometric information but also capture semantic attributes and relations, providing a concise depiction of the scene to facilitate complex planning using large language models (LLMs) to achieve a desired goal [16]–[20].

The advent of LLMs has transformed autonomous planning due to their ability to understand and reason about complex, context-rich scenarios [17]–[26]. Thus, our system integrates LLMs with hierarchical scene graph representations, enabling real-time adjustments and error corrections based on the agent’s evolving understanding of its environment. This enables the generation of context-aware high-level plans for autonomous task execution using information from the scene graphs. The execution of these high-level plans is aided by lower-level planners and controllers. We demonstrate this on a quadruped robot, focusing on navigation tasks that involve reaching specific objects as directed by the user, in unknown environments that change over time. This requires robust coordination between high-level planning and lower-level motion control, including obstacle avoidance. This coordination is crucial for effectively translating simple natural language instructions into complex plans and then into precise physical actions, enabling the robot to navigate autonomously to specific objects or rooms based on conversational commands.

<sup>\*</sup>Equal Contribution

<sup>1</sup>Control/Robotics Research Laboratory (CRRL), Department of Electrical and Computer Engineering, NYU Tandon School of Engineering, Brooklyn, NY, 11201. E-mails: {d.naren, rgg9769, ak10531, nkp269, rk4342, prashanth.krishnamurthy, khorrami}@nyu.edu

This work was supported in part by Google and in part by ARO under Grant W911NF-22-1-0028, and in part by the New York University Abu Dhabi (NYUAD) Center for Artificial Intelligence and Robotics (CAIR), funded by Tamkeen under the NYUAD Research Institute Award CG010.

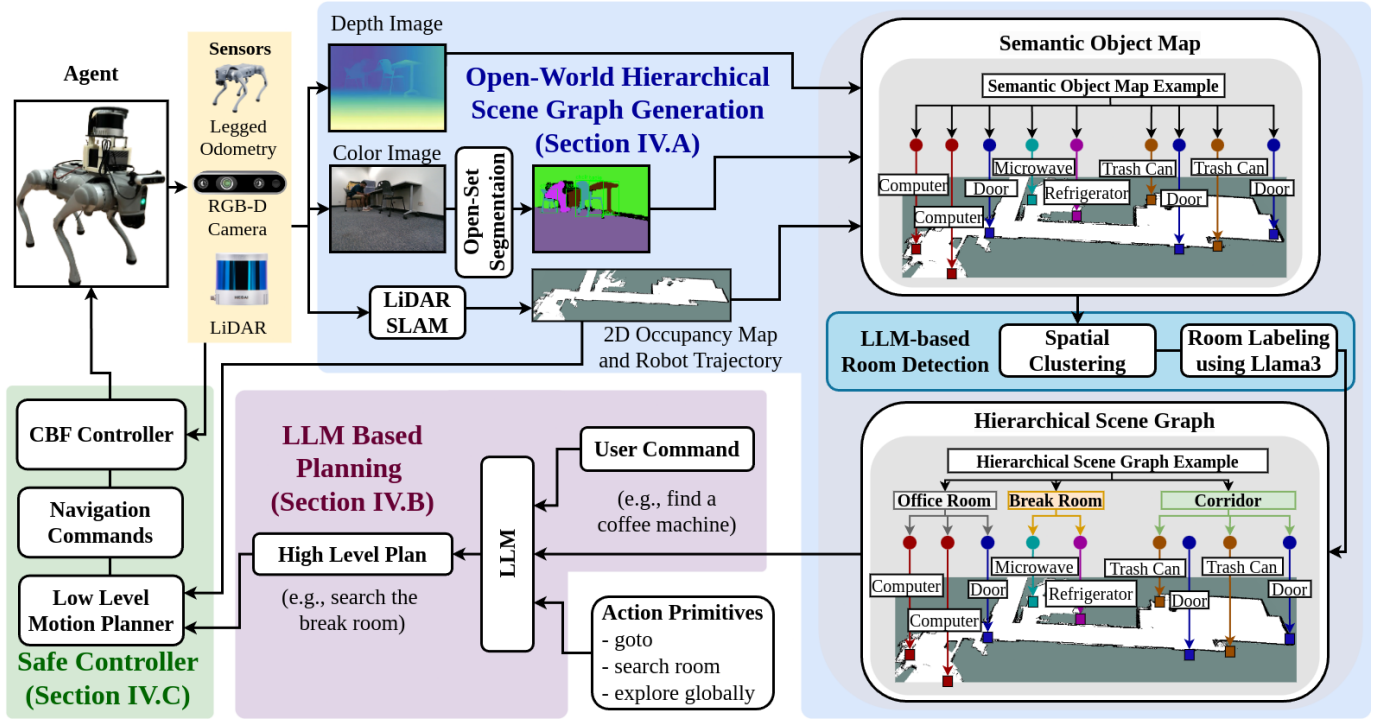


Fig. 1. **Overview of our OrionNav framework.** The OrionNav system fuses data from onboard LiDAR and odometry sensors for robust localization and mapping, while integrating open-world semantics to produce a **semantic object map** of the environment. This map is then clustered into distinct rooms, and room labels are assigned using the Llama3 LLM, generating a **hierarchical scene graph**. An **LLM-based planner** utilizes this scene graph, along with user commands, to create a high-level task execution plan, guiding low-level controllers to safely and efficiently achieve designated goals.

Our proposed OrionNav framework utilizes a quadruped robot equipped with a combination of RGB-D and LiDAR sensors to autonomously navigate an unknown environment to execute a remote user’s command. It generates an object map through integrated open-vocabulary semantic segmentation and SLAM that are continuously updated in real time as the robot navigates around the scene. The objects in the map are spatially clustered and associated with specific rooms using an onboard LLM to construct a hierarchical scene graph. This graph representation enables another LLM to formulate a contextually appropriate plan for executing user-defined navigation tasks. A key innovation of OrionNav lies in the system’s ability to continuously update the scene representations and plan online during robot operation. This allows the robot to swiftly respond to environmental changes and correct planning errors on the fly. This capability is essential for maintaining effectiveness in dynamic, unpredictable real-world settings where static pre-programmed behaviors often fail. The high-level plan guides the low-level motion planner to generate goal-directed navigation commands, while an obstacle-avoiding control barrier function (CBF) controller reactively adjusts the robot’s path when obstacles are detected. These adjustments ensure safe and efficient navigation toward the designated goals. The overall flow of the proposed OrionNav framework is illustrated in Fig. 1. All modules except for the LLM used for planning operate onboard the robot in real-time while demonstrating autonomy in complex environments to perform diverse tasks. Our key contributions are:

- A real-time, language-driven autonomous navigation that employs multi-level abstractions of perception and planning to execute diverse tasks while handling changes in environments and recovery from task failure.
- A robust and efficient open-vocabulary 3D semantic mapping approach for real-time generation and updating of hierarchical scene graphs during robot operation.
- A LLM-based planner that uses hierarchical scene graphs to create multi-step plans for navigation tasks with ability to adjust plan based on environmental changes, task progress, and potential task failure.
- An enhanced exploration strategy integrated with low-level motion planners and CBF-based controller to enable a robot to operate safely in dynamically changing environments while executing diverse navigation tasks.
- Real-world implementation of the framework on a quadruped robot, autonomously performing various real-time navigation and search tasks in environments using onboard sensors and computing resources.

The rest of the paper is organized as follows: Section II covers related work, Section III presents the problem and Section IV describes our framework, Section V shows experimental results, and Section VI concludes the paper.

## II. RELATED WORKS

### A. Scene Graph

*Scene Graphs* have emerged as an efficient and compact method for representing environments. They employ a tree

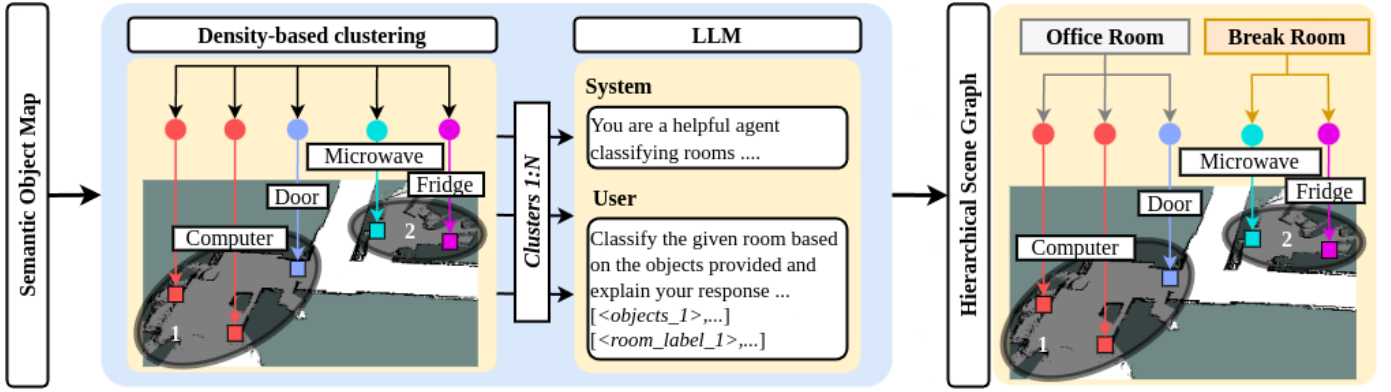


Fig. 2. Our method makes use of the semantic constructs of indoor environments to generate a hierarchical scene graph from the semantic object map. The objects are then clustered based on the density of the objects in the scene. An LLM is then queried with the labels of the objects in each cluster and a set of candidate room labels. Finally, the generated hierarchical graph is passed to the LLM-based planner.

structure consisting of ‘nodes’ that represent environmental features, connected by ‘edges’ that denote relative spatial or semantic relationships between node pairs. These graphs facilitate in understanding both global and local characteristics of the environment, thereby enabling more informed and accurate planning for executing a given task. The scene information represented by the nodes can be image and object instances conveyed through 3D point clouds [13], [27], [28], 2D masks [18], [20], [24], or language-guided 3D representations using radiance fields [29]–[32]. Edges are encoded to represent geometric or semantic connections between nodes, establishing relationships between these features. For example, an edge might indicate that an *apple* is *on a table* [33]. Such structures have been commonly used for image retrieval [34], [35], object retrieval [16], and task and motion planning [36].

Traditionally, Conditional Random Fields [37] were used to model graph nodes and edges but advancements in deep learning have enabled the use of graph neural networks to generate scene graphs from image instances and object masks [38], [39]. The advent of Large Visual-Language Models has facilitated the creation of open-vocabulary annotations of nodes and the use of language-aligned features to improve performance on downstream tasks. Foundational models [40]–[43] have been used to enrich scene graph nodes with semantic features like Contrastive Language-Image Pre-training (CLIP) [40] embeddings and linguistic annotations like labels, captions, and descriptions. Next, object nodes associated with corresponding CLIP embeddings obtained by fusing multiple 2D views, masked views, or re-projected 3D features [13]–[16]. However, these approaches are typically computationally expensive, preventing online scene graph generation. On the other hand, OrionNav uses a lightweight scene graph generation process, enabling it to work online. ConceptGraphs [24] augment the nodes with language captions to enable the use of generated scene graphs with LLM-based planners. Along with node enhancement, foundational models are also capable of generating semantic relations between nodes. Large Vision-Language models (VLMs) and LLMs such as LLaVa [42], BLIP [44] or GPT [41] are used either to directly generate spatial semantic relations between objects using SLAM-

derived poses and 2D camera views [24], or indirectly to train graph networks to populate the nodes and edges [13]. Within the domain of scene graph generation, two approaches closely align with our work: HOV-SG [16] and Clio [15].

HOV-SG [16] obtains 2D objects masks from image instances using the Segment Anything Model [45]. Their fused CLIP features are aggregated for 3D point corresponding to the objects, and later processed to identify a single feature for each object node. These are added to a hierarchical scene graph [46], [47], where heuristics are used to segment floors and rooms and their corresponding nodes are characterized using fused CLIP features. However, their computationally expensive clustering and classification methods make it difficult to implement them in real time. On the other hand, Clio [15] extends the Hydra framework [48] to generate scene graphs with open vocabulary labels, while introducing a novel approach to determine appropriate object representation granularity. It employs an Information Bottleneck formulation to cluster 3D primitives into task-relevant semantic concepts, thus enhancing the scene graph’s utility for specific tasks while generating it efficiently. However, such scene graphs are constrained by their reliance on predefined task sets, limiting their adaptability to unforeseen scenarios and tasks. Our method mitigates these shortcomings by facilitating real-time scene graph generation and allowing the execution of a diverse range of tasks, including those not originally anticipated.

### B. Semantic Mapping

Scene graph generation relies on information about the environment acquired through mapping and semantic labeling process. These processes integrate image or point cloud segmentation techniques with spatial mapping data generated from SLAM frameworks. SLAM provides crucial geometric and structural information about the environment, complementing the semantic understanding derived from segmentation. This synergy enables the creation of robust and contextually rich representations through semantic object maps, enabling better scene understanding for autonomous robot navigation and decision-making. [2], [4], [49]–[53]. Segmentation techniques are fundamental to semantic labeling and scene graph

generation, enabling detailed parsing of visual information to improve the quality and detail of scene representations. The field encompasses various approaches, each contributing unique capabilities to the interpretation of complex visual scenes. Semantic segmentation infers high-level category concepts, while instance segmentation groups foreground pixels into distinct object instances. Panoptic segmentation combines these two approaches, and entity segmentation infers masks of object instances without classification [11], [54], [55]. Panoptic segmentation, in particular, predicts both class and instance labels for each pixel, assigning instance labels to foreground objects ('things') and single labels to background parts ('stuff'). The advent of foundation VLMs like CLIP [40] has enabled more sophisticated approaches to fine-grained recognition and segmentation of objects in scenes in an open vocabulary setting [12]. These models leverage large-scale pre-training on diverse image-text pairs, allowing them to generalize to a wide range of visual concepts and enabling zero-shot learning capabilities which are made robust with out-of-distribution detection [56]. Recent developments have further extended these capabilities to open-vocabulary semantic and instance segmentation tasks, demonstrating impressive zero-shot generalization to novel object categories [11], [57]–[63].

Semantic mapping techniques, vital for scene graph generation, extend beyond 2D segmentation to encompass complex 3D structures. Reconstruction of 3D structures with object recognition and grouping remains a fundamental challenge in computer vision, especially in the context of 3D scene understanding. These methods help a large language model understand the scene through reasoning [64], [65] for building grounded 3D large language models [66]–[72]. The semantic information aids in both understanding the scene as well as improving the accuracy of the map and odometry for semantic SLAM [8], [73]–[81]. End-to-end 3D methods that predict 3D semantic information directly from sensor modalities like camera and depth, or from intermediate representations like reconstructed point clouds or voxel-based representations, can capture complex 3D relationships. Recent works have significantly improved the performance of open vocabulary 3D segmentation using point clouds from LiDAR or depth data from RGB-D cameras [82]–[90]. In addition, there have also been approaches fusing embeddings from open world VLMs like CLIP with neural radian fields, Gaussian splats or implicit neural network based representation [29]–[32] for generating a consistent segmentation of the scene. The open vocabulary allows for greater flexibility in real-world applications and can help an LLM to understand the scene with multiple objects, which multi-modal LLMs hallucinate or fail to perceive [91]. However, these methods do not scale well in larger scenes making them computationally intensive and memory-consuming, rendering them impractical for real-time segmentation in robotic applications.

OrionNav draws inspiration from prior works while prioritizing efficiency, striking a balance between adaptability and real-time performance. The proposed method leverages the strengths of both 2D and 3D approaches, utilizing the efficiency and generalization capabilities of 2D CLIP-based models while incorporating 3D geometric information to enhance

spatial understanding. It integrates a CLIP-based open vocabulary 2D segmentation model with 3D geometric information from maps and odometry for improved scene understanding. We utilize the FC-CLIP [92] model with a ConvNeXt [93] CLIP [40] backbone, enabling efficient processing of high-resolution images where transformer-based models struggle. This architecture maintains translation equivariance and allows multi-scale feature fusion, crucial for detailed environmental segmentation in robotic applications. The system performs lightweight updates to the semantic object map leveraging spatial and semantic properties of the scene and observed objects in it. This map can be further utilized by a scene graph to generate hierarchical object clusters and relationships. This hybrid approach offers a robust, efficient solution that addresses trade-offs between computational efficiency and comprehensive 3D scene understanding.

### C. Multi-Modal Language Models for Robotic Task Planning

In recent years, robotic task planning has benefitted from the use of LLMs [17]–[23], [25], [26]. Early studies in this field primarily focused on manipulation tasks and humanoid simulations [21]–[23]. Recent works have shifted toward generating plans for mobile robots operating in larger environments [17]–[20], [26], [94]. Strategies such as environment feedback through the use of the Python assert keyword in the plans [21], affordance functions [25], and success detection and scene description [23] have been used to ground the LLM planners in the environment. In addition, robot navigation frameworks have also employed VLMs to build semantic maps for embodied question answering, replacing frontier-based search by leveraging semantic knowledge of the environment. OK-Robot [95] proposed an end-to-end system for open-ended pick-and-place tasks using pre-trained models. Specifically, they utilize CLIP [40] to convert language queries into semantic vectors, which is then used to search for the semantically closest voxel in a semantic point cloud. PixelNav [94] proposes a direct RGB-based solution that utilizes an LLM planner for waypoint generation in RGB space for object navigation tasks. However, these methods either create an environment representation that cannot be used for subsequent tasks in the same environment or require a prebuilt scene representation.

Studies focusing on mobile robots operating in large environments have often utilized scene graphs to ground LLMs in their environment [16]–[20], [24], a strategy we also adopt in OrionNav. ConceptGraphs [24], while focusing primarily on scene graph generation, employed an LLM-based planner to retrieve target objects from a pre-built scene graph as a downstream task. Similarly, HOV-SG [16] emphasized scene graph generation and utilized an LLM to parse complex natural language queries for object retrieval from the pre-built scene graph. In contrast, we introduce an LLM planner with several action primitives that function as a high-level controller, actively participating in the scene graph generation process. While SayPlan [17], MoMa-LLM [20], and SayNav [18] also leverage LLMs as high-level planners, they face constraints that limit their practical application. SayPlan [17] requires pre-built scene graphs, whereas MoMa-LLM [20] and SayNav

[18] rely on known robot poses and semantic object masks. In contrast, OrionNav overcomes these constraints by performing perception, planning, and control for navigation tasks using onboard sensors in real-time, enabling autonomous navigation in diverse, unknown, and dynamic real-world environments.

### III. PROBLEM FORMULATION

We address the challenge of developing an autonomous language-driven physical agent capable of operating in dynamic, unknown environments. The agent, a quadruped robot, is equipped with an RGB-D camera, LiDAR, and leg odometry sensors (IMU, joint encoders, and contact sensors). The primary objective is to enable the robot to autonomously search for and navigate to specific objects or rooms in response to natural language queries issued by a remote operator. The robot must operate without prior knowledge of the environment, which consists of multiple interconnected rooms containing various objects. The environment is subject to change over time, with no predetermined types of objects, allowing for arbitrary categories and adding complexity to the search and navigation tasks. The robot must interpret queries, explore the environment, and locate targets in real time while adapting its strategies when initial attempts are unsuccessful. Performance is evaluated based on the robot’s ability to successfully locate and reach queried targets, considering time and path efficiency. The problem encompasses challenges in perception, navigation, natural language understanding, and decision-making in unstructured, real-world settings. The desired outcome is a system that can build and maintain a memory-efficient representation of the environment and objects within it, which can be utilized to handle subsequent queries more efficiently.

### IV. METHODOLOGY

This section introduces OrionNav, our novel framework designed to address the challenges outlined in Section III. The framework efficiently explores the environment, constructing a hierarchical scene graph that captures object semantics and relationships (Section IV-A). To interpret human queries, OrionNav employs the LLM-based high-level planner that translates natural language into structured commands, which are then executed using the scene graph information (Section IV-B). It employs safe motion primitives for exploration and navigation, reducing accident risk (Section IV-C). By combining open-vocabulary semantic mapping, hierarchical scene graph generation, and LLM-based query interpretation with a low-level motion planner and controller, OrionNav enables safe autonomous operations that directly address the problem outlined in Section III. Importantly, these modules, which are further elaborated in the following subsections, are designed with modularity in mind, allowing for easy substitution as technological advancements occur or to meet varying user requirements, provided that new components adhere to the established communication protocols.

#### A. Open-World Hierarchical Scene Graph

1) *Open-Vocabulary Semantic Segmentation*: The agent uses color frames from the RGB-D camera to semantically

segment the scene. We use an off-the-shelf open-vocabulary semantic segmentation framework, FC-CLIP [92], that comprises a class agnostic mask generator, in-vocabulary classifier, and out-of-vocabulary classifier. All modules leverage multi-scale feature maps extracted from a frozen ConvNeXt-L [93] CLIP [40] backbone. The mask generator passes these feature maps through a pixel decoder with multi-scale deformable attention, as described in Mask2Former [96], to enhance the pixel features. These enhanced features, along with object query embeddings, are then passed through a mask decoder with each layer in the decoder consisting of a masked cross-attention layer, a self-attention layer, and a feed-forward network. The object query embeddings are a fixed number (250) of learnable embeddings representing each potential object the model can predict. The mask decoder predicts a fixed set of candidate masks (100 masks) representing the segmentation of each object in the scene. The number of object queries and masks are independent, allowing flexible prediction of varying object counts while improving computational efficiency by generating masks only for valid predictions.

Category Text Prompt Templates	
• <label>.	• There is a small <label> in the scene.
• a photo of a <label>.	• There is a medium <label> in the scene.
• a dark photo of a <label>.	• There is a large <label> in the scene.
• This is a photo of a <label>	• a photo of a hard to see <label>.
• There is a <label> in the scene	• a bright photo of a <label>.
• There is the <label> in the scene	• a close-up photo of a <label>.
• a photo of a <label> in the scene	• a black and white photo of <label>.
• a photo of a small <label>.	• a bright photo of the <label>.
• a photo of a medium <label>.	• a blurry photo of the <label>.
• a photo of a large <label>.	• a photo of the <label>.
• This is a photo of a small <label>.	• a close-up photo of the <label>.
• This is a photo of a medium <label>.	• a low resolution photo of a <label>.
• This is a photo of a large <label>.	• a blurry photo of a <label>.

Fig. 3. Prompt templates used to generate text embeddings for each category. Each instance of ‘<label>’ is replaced with corresponding category names.

The candidate masks in FC-CLIP require category text embeddings for each class to enable zero-shot classification. These embeddings of dimension 768, are generated using a method similar to CLIP, employing a set of predefined template prompts (e.g., ‘a photo of a <label>’, ‘an image of a <label>’) for each class as shown in Fig. 3. For each class, the class name and its closely related categories (e.g., cabinet, cupboard, bureau, china cabinet) are individually inserted into these templates, tokenized [97], and processed through a CLIP text encoder. Importantly, each synonym generates its own unique embedding, allowing the model to capture subtle semantic differences. However, during the classification process, these separate embeddings are all associated with a single, general class name (e.g., ‘cabinet’) for output, which aids in LLM room classification and scene graph generation. This approach allows for on-the-fly computation of embeddings for new objects without additional training and enhances the model’s robustness to different object naming conventions.

The classification process involves both out-of-vocabulary and in-vocabulary classifiers. The out-of-vocabulary classifier, which applies mask pooling to the frozen CLIP backbone’s feature maps, is designed to maintain CLIP’s open-vocabulary recognition ability for unseen classes. The in-



vocabulary classifier, which pools the pixel decoder’s output feature maps, is optimized for better performance on classes seen during training. Both compute classification scores using cosine similarity between the pooled features and each of the synonym embeddings, followed by softmax with a temperature of 0.07. The highest score among synonyms determines the class prediction. These scores from both classifiers are then combined using geometric ensembling, leveraging the strengths of both classifiers to improve overall performance on seen and unseen classes. Importantly, experiments have shown that using both classifiers also improves performance in out-of-vocabulary classes, demonstrating the synergistic effect of this approach. Candidate masks are filtered out based on classification scores or if they belong to the background class, with further filtering out of ambiguous or inconsistent masks where the percentage of pixels belonging to majority classes falls below a specified threshold. Remaining masks  $\mathcal{M}_o$  are associated with their corresponding object class  $o$ , using the general class name. Objects like ceilings, floors, and walls are segmented but excluded from object mapping to improve robustness against outlier masks. The classes for the classifiers consist of classes from COCO Panoptic dataset [98]–[100] (used to train FC-CLIP) and additional class labels relevant to indoor environments, along with their synonyms. All class labels are mapped to appropriate general class names, with provisions for adding new classes as needed. This process enables FC-CLIP to perform effective zero-shot segmentation across a wide range of object classes, including unseen ones while adapting to various naming conventions and ensuring consistent class names for semantic object mapping.

2) *Semantic Object Mapping*: The 2D segmentation mask,  $\mathcal{M}_o$  (for object  $o$ ), obtained from the model described in Section IV-A1 are used to estimate the object’s center  $(u_o, v_o)$  in the camera plane. These center points are then projected from the camera plane to 3D coordinates  $P_o^c = (X_o^c, Y_o^c, Z_o^c)$  in the camera frame, using depth frame from RGB-D camera, and camera intrinsics. This approach was selected over alternatives such as using 3D bounding boxes or computing the mean of the point cloud, considering simplicity and computational complexity. These object observations are subsequently converted from the camera frame to the world frame using the pose from the robot’s odometry. We employ an off-the-shelf SLAM toolbox [101] that utilizes the robot’s onboard LiDAR and legged-inertial odometry to build a 2D map of the environment as well as estimate the pose of the robot.

The detected objects ( $P_o^c$  in the camera frame) are transformed to the world frame  $P_o^w$  using the estimated pose and associated with previous map objects. The map objects are observed in prior frames and frame objects are objects observed in the current frame. All frame objects observed in the first frame are considered as map objects and if there are multiple map objects belonging to the same class, they are given a different instance id. The instance id is the instance count of a particular class in the scene and is thus initialized at 1 and incremented with new observations. For the subsequent frames, the frame object and map object are associated using density-based clustering (DB-SCAN) [102] based assignment as it offers a good balance between accuracy

and computational speed compared to more complex methods like Hungarian matching [103]. Frame objects are projected into the world coordinate frame and DB-SCAN is applied to group frame and map objects of the same class, using a local radius threshold of  $\epsilon$  meters. For each resulting cluster, only one map object is retained, while the others are discarded. DBSCAN is effectively able to handle grouping without using heuristics, such as when a large object is detected as several smaller bounding boxes across consecutive frames.

Three cases for grouping objects are considered: for larger objects like tables, doors, and noticeboards, we set  $\epsilon = 2$  meters to account for their size; for smaller objects, such as books, a reduced threshold of  $\epsilon = 0.5$  meters is used; and for all other objects, a threshold of  $\epsilon = 1$  meter is applied. All the objects without any associations are assigned as new map objects with unique instance id, thus forming a semantic object map  $\mathcal{S}$ . The system handles updates in object locations following loop closure and local bundle adjustment using a pose graph-based approach. Each object in the semantic object map  $\mathcal{S}$  is associated with the keyframe in which it was first observed. A relative transform  $T_o^k$  between the object and its associated keyframe in the pose graph is maintained. When loop closure or local bundle adjustment occurs and keyframe poses are updated, the location of each object is adjusted by applying the updated keyframe pose to the stored relative transform:  $T_o^w = T_k^w \cdot T_o^k$ , where  $T_o^w$  is the updated world pose of the object,  $T_k^w$  is the updated pose of the keyframe in the world coordinates, and  $T_o^k$  is the relative transform between the object and its keyframe. An alternative approach involving bundle adjustment using all keyframes that observe an object was considered but not chosen due to its higher complexity and computational cost, favoring simplicity and efficiency in our current system. Our method ensures that object locations remain consistent with the refined map and robot trajectory, maintaining the integrity of the semantic object map even as the underlying SLAM system refines its estimates.

3) *Hierarchical Scene Graph*: Indoor environments are inherently structured into distinct semantic components, primarily floors and rooms. These spaces are typically categorized using familiar natural language terms such as ‘office,’ ‘kitchen,’ and ‘bedroom,’ based on their intended function and the objects they contain. The consistency of these labels across various indoor settings makes them highly informative semantic descriptors. For instance, a room labeled ‘kitchen’ has a significantly higher likelihood of containing a microwave compared to one designated as an ‘office’. This semantic structuring of indoor spaces lends itself to efficient representation through 3D scene graphs with room hierarchies, as proposed in prior works [15], [16], [46], [47]. Such representations not only capture the physical layout of indoor environments but also encode the semantic relationships between spaces and objects, facilitating more intuitive and context-aware spatial understanding and navigation. The indoor spaces consist of several rooms connected by walkways called corridors. Hence, it is natural to find multiple clusters of objects while exploring such environments. Making use of this observed property, we propose classifying the scene into multiple areas by applying density-based clustering (DBSCAN [102]) on the coordinates

of the map objects. For each identified cluster, we leverage LLMs to determine the most appropriate room label. Our approach involves inputting the labels of all objects present in each area cluster to the LLM, along with a curated list of candidate room labels and general guidelines to ensure consistency across the environment. The generation of these labels occurs in real-time and on-device, utilizing a memory-efficient 4-bit quantized version of the LLaMA 3 [104] model. While prior works rely on language-aligned visual embedding models to find the best match [15], [16], [20], our experiments indicate that modeling the influence of certain combinations of objects toward accurate room classification is fairly difficult. Hence, we leverage an LLM to identify labels that are consistent with the general notion of such semantic concepts.

The efficient storage of the scene graph is achieved through a tree structure. Each node in the tree contains properties such as its hierarchy level, text label, class-specific unique ID relative to the parent node, location in the map, and parent and child relationships as shown in Fig. 2. This structure is easily converted to a Python dictionary format, allowing direct input as a prompt to the LLM-based planner, as detailed in Section IV-B. Our approach combines density-based clustering for spatial organization, LLM-powered semantic labeling for accurate room classification, and an efficient tree-based representation for the scene graph. This integrated method offers a robust and flexible approach to understanding and representing complex indoor environments, overcoming the challenges of traditional visual embedding models in capturing the nuanced relationships between objects and room types. By leveraging the semantic understanding capabilities of LLMs, we achieve a more intuitive and context-aware representation of indoor spaces, which can significantly enhance the agent’s ability to navigate, localize objects, and comprehend the overall layout of the indoor scene for autonomous task execution.

### B. LLM Planner

The LLM planner leverages GPT-4-Turbo [41] to navigate and interact with the environment through three core action primitives: *goto(<room name> or <object name>)*, *search\_room(<room name>)*, and *explore\_globally()*. The *goto* command directs the robot to the specified room or object and queries a visual agent to verify task completion. The *search\_room* command moves the robot to the specified room and then rotates to detect nearby objects. The *explore\_globally* command initiates frontier-based exploration of the environment. The documentation of these commands with detailed explanations and use-case scenarios, is provided to the LLM (see Fig. 4 for prompting details). The LLM’s grounding in the environment is achieved through a hierarchical scene graph formatted as JSON, which stores information about rooms and the objects within them (as described in Section IV-A3). Leveraging this scene graph, the LLM selects an appropriate action and fills in the necessary arguments to accomplish the user-defined task, while also providing reasoning for its decision. The resulting output is directly executed as a function call to the robot API, implemented using the ‘Safe Robot Controller’ described in Section IV-C.

---

### Algorithm 1 Task Execution

---

```

1: Input: user_command
2: action_history  $\leftarrow$  []
3: feedback  $\leftarrow$  None
4: task_accomplished  $\leftarrow$  False
5: while not task_accomplished do
6:   scene_graph  $\leftarrow$  HierarchicalSceneGraph.get_latest()
7:   action  $\leftarrow$  LLMPlanner.call(user_command,
   scene_graph, feedback, action_history)
8:   action_history.append(action)
9:   action_success  $\leftarrow$  RobotApi.send(action)
10:  if action_success then
11:    if action = goto(<arg>) then
12:      task_accomplished  $\leftarrow$  ask_agent('Is task done?')
13:      if not task_accomplished then
14:        feedback  $\leftarrow$  'Task has not been accomplished.'
15:      else
16:        feedback  $\leftarrow$  None
17:      end if
18:    else
19:      task_accomplished  $\leftarrow$  False
20:      feedback  $\leftarrow$  'Task has not been accomplished.'
21:    end if
22:  else
23:    task_accomplished  $\leftarrow$  False
24:    feedback  $\leftarrow$  RobotApi.get_recent_action_logs()
25:  end if
26: end while

```

---

The plan generation framework initiates with a 360° environmental scan if the scene graph is empty. The LLM then generates an action based on the given task and hierarchical scene graph. Following action execution, if the task remains incomplete, the LLM receives an updated scene graph and generates a subsequent action. Error scenarios—stemming from incorrect syntax, low-level controller issues, or environmental constraints like inaccessible target points—trigger feedback messages explaining the error, prompting the LLM to produce alternative actions. After the initial action, a command history is maintained to prevent action repetition. This iterative action generation process is formalized in Algorithm 1. The action primitives used by the LLM are elaborated as follows:

1) *goto(<room name> or <object name>)*: The *goto* primitive enables the robot to navigate to a specified room or approach a particular object within the environment. When tasks involve navigating to an object, the LLM planner is instructed via system prompt to provide both the object’s name and the room it resides in as arguments. For example, to locate a ‘printer in the office,’ the command would be *goto(office-2, printer-1)*. Conversely, for tasks requiring navigation to a specific room, only the room name is specified, such as *goto(office-1)* for ‘an office with three computers.’ In each command, the name consists of a label and a unique instance ID, where object labels are assigned during semantic object mapping, and instance IDs and room labels are generated during hierarchical scene graph creation. This primitive also supports task completion verification by interacting with a

visual agent upon arrival at the destination.

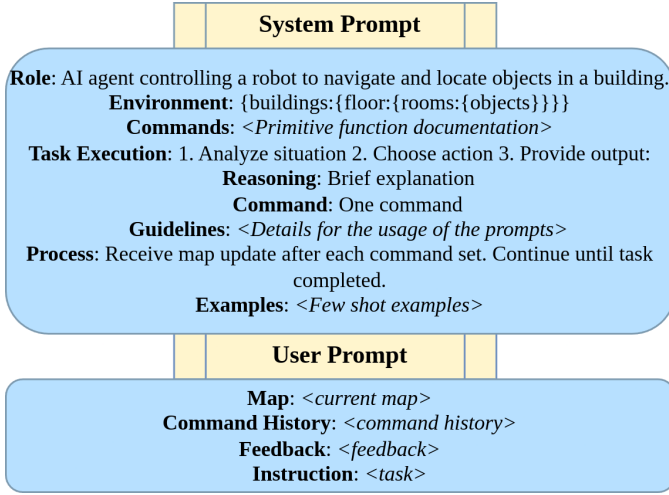


Fig. 4. **Overview of the system and user prompts for the LLM planner.** The system prompt explains the agent’s role, action primitives, and map format. The user prompt provides the map, command history, feedback, and task details. In the initial call to the LLM, command history and feedback are absent, as there is no prior interaction. Feedback includes task status and error messages from previous command executions. ‘<text>’ represents a placeholder for the information summarized in *text*.

2) *search\_room(<room name>)*: This primitive commands the robot to navigate to the specified room and perform a 360° rotation within it. This action is designed to explore a partially observed room, enabling the robot to detect and recognize objects in the scene and update its scene graph. Even if the robot has previously visited the room, limitations in perception or environmental changes may have caused it to miss or misidentify objects. By re-observing the room, the robot gathers more complete and up-to-date information, which is critical for accurate planning and task execution. While this command might be redundant in simulated environments with ground truth perception and segmentation, it is essential in real-world settings where the processed data from the sensors can be imperfect or incomplete. The ability to update the scene graph with newly observed objects as well as updated locations ensures that the LLM has access to the most current environmental information for decision-making.

3) *explore\_globally()*: When the scene graph is empty or lacks sufficient information for the LLM to utilize the other commands effectively, the *explore\_globally()* primitive is employed. This action initiates an autonomous exploration phase, typically lasting one minute, using a frontier-based exploration strategy as detailed in Section IV-C3 to enable identification of unexplored areas. The robot continuously updates its scene graph using new observations from onboard sensors during the exploration phase. Previously unobserved aspects of its surroundings are updated by rotating the quadruped around its vertical axis. The enhanced scene representation enables the LLM planner to generate more specific and task-relevant commands. The duration of the exploration is selected to balance thoroughness with operational efficiency, ensuring that the robot gathers enough information to proceed with task-specific actions without unnecessary delay. This primitive is

vital when the robot starts in an unknown environment or when significant portions of the environment remain unexplored, providing the foundational knowledge required for effective planning and execution of subsequent tasks.

### C. Safe Robot Controller

1) *Navigation*: The ROS2 navigation stack is used to execute the *goto* commands generated by the LLM planner. The agent begins by retrieving the locations of the specified objects or rooms from the scene graph. This scene graph is continuously updated as the robot explores and gathers new sensory data. The retrieved location is used to generate a target goal location for global trajectory planning. A cost-aware  $A^*$  search method is used for global trajectory planning. This method calculates the optimal path to the target location by searching through a grid-based representation of the environment. The global cost map required for this search is computed from an occupancy grid map generated from LiDAR mapping.

To enhance safety and navigation efficiency, obstacles in the occupancy grid map are inflated radially. This inflation accounts for the robot’s physical dimensions and adds a safety buffer to prevent collisions. The inflation radius and cost scaling factors are manually tuned parameters that define how close the robot can approach obstacles and how the proximity to obstacles affects the path cost. Considering the non-linear dynamics inherent in the robot’s motion model and the necessity of maintaining a forward-looking planning horizon for improved safety, we integrate a sampling-based Model Predictive Path Integral (MPPI) controller for local trajectory execution. The MPPI controller outputs the desired velocity commands to the locomotion controller on the robot.

2) *CBF-QP Action Filter*: While the MPPI controller accounts for local obstacles through the cost map, it does not inherently guarantee collision-free actions with surrounding obstacles. The MPPI optimizes trajectories based on a predictive model and cost functions but may not react swiftly enough to sudden changes or unmodeled dynamics that could lead to unsafe situations. To further enhance the safety of our motion planning pipeline, we incorporate a Control Barrier Function Quadratic Programming (CBF-QP) safety filter as proposed in [105]. This safety filter is positioned between the MPPI controller’s output and the robot’s low-level velocity tracker, acting as an additional layer of safety oversight.

The CBF-QP safety filter operates analogously to the reactive instincts observed in biological organisms, establishing a low-latency connection between perception and action. Specifically, it directly utilizes raw LiDAR point cloud data to assess the immediate environment and modifies the control commands to avoid collisions proactively. Integrating this filter into our motion planning pipeline significantly enhances the robot’s ability to navigate safely in complex and dynamic environments. It provides a robust safeguard against potential collisions, complementing the predictive capabilities of the MPPI controller with real-time reactive adjustments. This strategy ensures that the robot not only follows optimal paths but also adapts instantaneously to unforeseen obstacles, thereby improving overall navigation safety and reliability.



3) *Exploration*: To facilitate efficient and comprehensive exploration of unknown environments, we utilize the *m-explore* ROS2 package<sup>1</sup>, which is based on the exploration strategy described in [106]. We have implemented several key modifications to adapt the algorithm to our specific requirements and to enhance its robustness in practical applications.

The exploration process relies on the robot’s current 2D position and the occupancy grid map generated using SLAM, as detailed in Section IV-A2. The occupancy grid map represents the environment in a discretized form, where each grid cell is classified into one of three categories based on sensor data: unexplored, occupied, or explored. Unexplored cells are those that the robot has not yet observed, occupied cells contain obstacles or impassable terrain, and explored cells represent free space that has been observed and mapped. To identify areas on the boundary between known and unknown space, we cluster unexplored grid cells that are adjacent to explored cells. This clustering process groups contiguous unexplored cells into distinct regions, and the centroid of each cluster is designated as a *frontier point*. Frontiers represent gateways to unexplored regions and are critical targets for the robot to expand its knowledge of the environment.

Each frontier  $\mathcal{F}$  is characterized by several key parameters, including its 2D coordinates  $(\mathcal{F}_x, \mathcal{F}_y)$  in the map reference frame, the Euclidean distance  $\mathcal{F}_d$  from the robot’s current position to the frontier point, and the number of grid cells  $\mathcal{F}_n$  in the cluster, which indicates the potential area that can be explored by visiting that frontier. We compute a weight  $\mathcal{F}_w = \alpha \mathcal{F}_d + \mathcal{F}_n$  for each frontier to prioritize them for exploration. We set  $\alpha = 50$  to prioritize closer frontiers and balance the trade-off between minimizing travel distance and maximizing exploration gain. A higher value of  $\alpha$  increases the influence of proximity, encouraging the robot to select closer frontiers to reduce navigation time and energy consumption, while the term  $\mathcal{F}_n$  promotes frontiers that lead to larger unexplored areas, maximizing the informational gain.

The robot selects the frontier with the highest weight  $\mathcal{F}_w$  as the next navigation goal, ensuring efficient expansion of the map by choosing frontiers that offer the best trade-off between proximity and exploration potential. However, not all frontiers are navigable due to environmental constraints such as obstacles blocking the path, narrow passages that the robot cannot traverse, or inaccuracies in the map stemming from sensor noise or localization errors. These non-navigable frontiers are identified through path planning attempts; if the robot cannot find a viable path to a frontier using its navigation stack, the frontier is marked as *blacklisted* and excluded from future consideration. Unlike the original implementation of *m-explore*, which terminates the exploration process when all frontiers are either explored or blacklisted, our modified algorithm continues exploration by assigning random goals within the already explored free space when no navigable frontiers are available. This strategy serves multiple important functions. First, by moving to different locations within the explored area, the robot may gain new perspectives on previously blacklisted frontiers. Changes in sensor viewpoints or

accumulated mapping data can reveal new paths or correct mapping errors, potentially converting blacklisted frontiers into navigable ones. Second, the robot may uncover new frontiers that were not previously detected due to occlusions, limited sensor range, or initial mapping inaccuracies. Continuous movement increases the likelihood of observing unexplored areas that were not visible from previous positions. Third, revisiting explored areas allows the robot to refine its occupancy grid map, correcting any errors or inconsistencies caused by sensor noise or environmental changes, leading to a more accurate and reliable representation of the environment. By integrating this strategy, we enhance the robustness of the exploration algorithm, ensuring that the robot continues to make progress even when initial exploration efforts stall. This approach mitigates the impact of measurement errors and uncertainties inherent in real-world environments. The robot remains active, persistently seeking opportunities to expand its map and improve its understanding of the surroundings.

In summary, our modifications to the *m-explore* package enhance autonomous exploration in scenarios involving measurement errors in mapping due to perceptual limitations or error, changes in environment, and complex obstacle configurations. The improved algorithm balances exploration efficiency and thoroughness through weighted frontier selection, frontier blacklisting, and a novel strategy of random goal assignment in explored space when viable frontiers are unavailable. This approach enables continuous exploration even when initial efforts stall due to environmental constraints or mapping inaccuracies. The enhanced capabilities enable the construction of accurate scene graphs to empower the LLM planner to generate more informed and contextually appropriate commands.

## V. EXPERIMENTS

### A. Experiment Setup

Our experiments aim to achieve three primary objectives: (a) demonstrate the real-world applicability of our pipeline through various end-to-end experiments (Section V-B), (b) showcase the advantages of utilizing an LLM as a high-level planner (Section V-C1), and (c) highlight the benefits of incorporating hierarchical scene graphs and their integration with the LLM planner (Section V-C2). We organized the experiments into three scenarios to evaluate different aspects of our system. Scenarios A and B focus on object retrieval tasks, where the system must locate and navigate to a specific object. Scenario C addresses room navigation tasks, involving detecting and navigating to a particular room within the environment. In total, we conducted 96 end-to-end experiments, as detailed in Table I. In Scenario A (short-range object retrieval), the target objects were placed within the same room or area as the robot’s starting position, with no occlusions between the robot and the object. This scenario tests the system’s ability to quickly identify and approach objects in an environment with minimal navigation complexity and visual obstructions. Scenario B (long-range object retrieval) presented a more complex challenge, requiring the robot to navigate through corridors and multiple rooms to locate the target objects, thereby testing the system’s mapping, navigation and planning

<sup>1</sup><https://github.com/robo-friends/m-explore-ros2>

TABLE I  
EXPERIMENT SCENARIOS

Task Type	Example Prompts	No. of Exp.
<b>Scenario A:</b> Short Range Object Retrieval	find a potted plant, go to a monitor, find a chair	36
<b>Scenario B:</b> Long Range Object Retrieval		45
<b>Scenario C:</b> Room Navigation	find a break room	15
<b>Total</b>		<b>96</b>

capabilities in a more realistic setting. Scenario C (room navigation) assessed the robot’s ability to interpret higher-level navigation goals by detecting and moving to a specified room, independent of specific objects. This highlights the integration of hierarchical scene understanding with LLM planning.

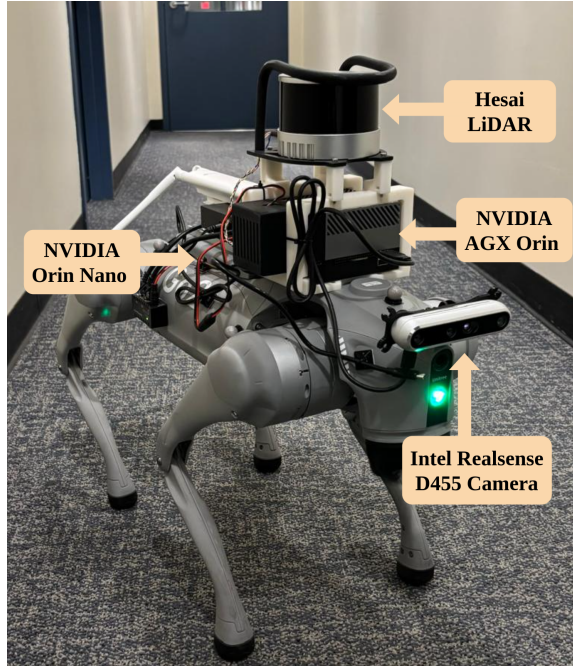


Fig. 5. **Robot Setup:** OrionNav’s capabilities are demonstrated on a Unitree Go2 quadrupedal robot equipped with onboard LiDAR sensor, stereo camera, and embedded computers equipped with a graphics processing unit (GPU).

The experiments were conducted in four distinct office environments, each featuring multiple rooms and corridors, as depicted in Fig. 5. These environments contain a variety of objects, such as computers, chairs, microwaves, coffee makers, and potted plants, alongside different room types including offices, break rooms, and conference rooms, providing diverse scenarios for the robot’s tasks. All experiments utilized the Unitree Go2 quadrupedal robot as the physical agent, depicted in Fig. 5. It is equipped with a top-mounted Hesai LiDAR for high-resolution 3D mapping and a Realsense D455 stereo camera for visual perception. All modules in our framework, except for the LLM planner, run onboard the robot, while the LLM planner (GPT-4-Turbo [41]) is accessed via a cloud API. Onboard computation is handled by two key modules: the NVIDIA Jetson AGX Orin (32 GB RAM) and the Jetson Orin

Nano (8 GB RAM). The AGX Orin handles computationally intensive tasks like open-vocabulary semantic segmentation, semantic object mapping, and scene graph generation. The Jetson Orin Nano handles mapping and odometry at 10 Hz, matching the LiDAR’s update rate, by utilizing data from LiDAR and legged-inertial sensors and employing the SLAM Toolbox [101]. It also handles the low-level control for navigation, comprising the ROS2 navigation stack, MPPI controller with CBF-QP action filter, and exploration module.

During the experiments, the stereo camera outputs color and depth frames at 30 Hz, each with a resolution of  $640 \times 480$ . The open-vocabulary semantic segmentation model uses the color frames without resizing to generate object masks at around 10 Hz. The model was optimized for memory and computational efficiency through vectorized operations and TensorRT with BF16 quantization, enabling it to recognize 560 distinct object classes and process 1306 text categories, with the flexibility to dynamically add new classes as needed. The semantic object map is refreshed at 2 Hz by integrating the odometry output, depth frames from the stereo camera, semantic information from the segmentation model, and pose graph corrections from bundle adjustment. The hierarchical scene graph is generated from the semantic object map using spatial clustering and room hierarchies, updating every two seconds. Room label

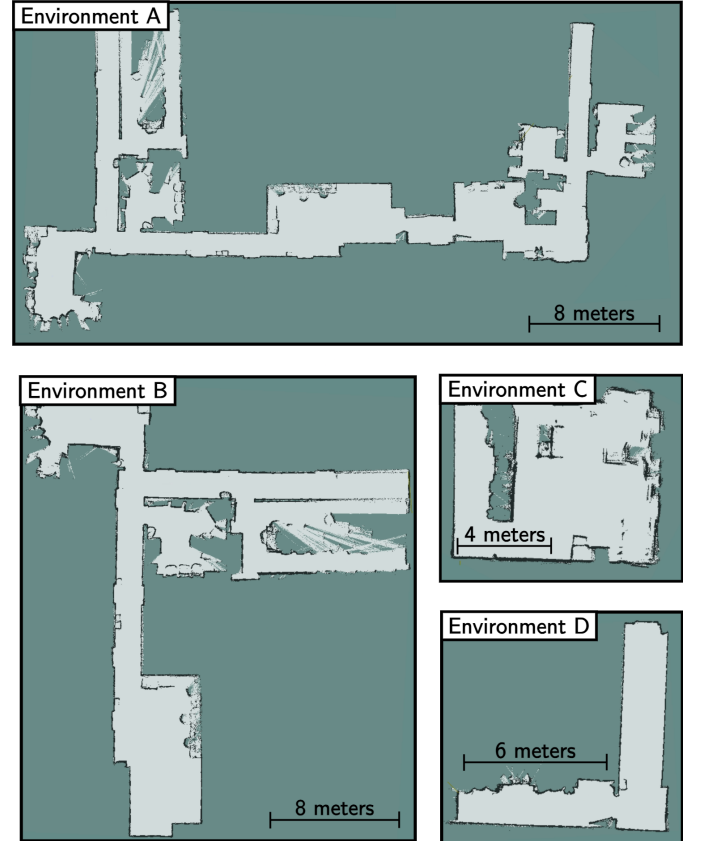


Fig. 6. **Experiment Environments:** The experiments were performed in four distinct environments. Environment A is a large, multi-room space with interconnected corridors. Environment B is a subset of Environment A, while Environment C is a smaller, single-room setting with diverse objects. Environment D comprises two perpendicular corridors, forming an L-shaped layout. The 2D LiDAR maps created using SLAM are shown in the figures.

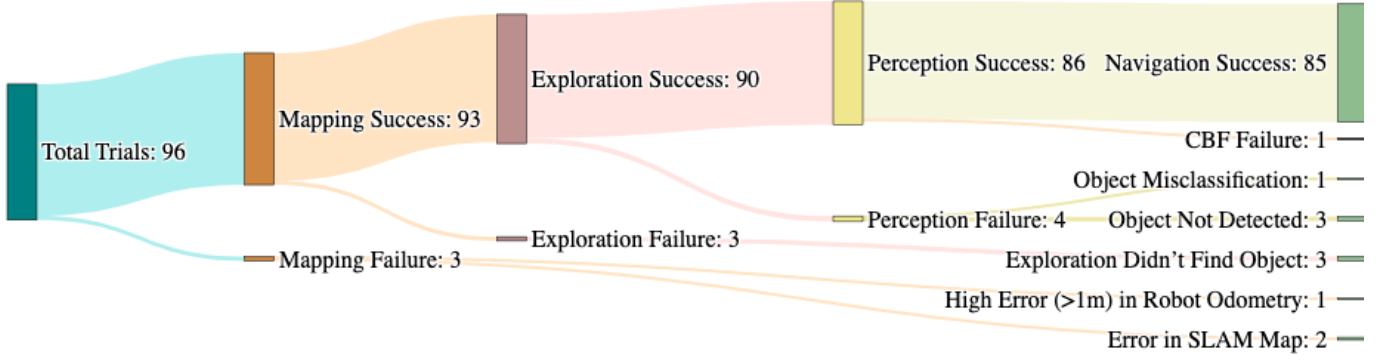


Fig. 7. **Success and failure cases of OrionNav across all experiments.** Visualization of our task execution framework with breakdown of different observed failures that occur due to perception, and navigation failures. A detailed breakdown of failures provided on the right.

predictions are handled by a quantized LLaMA 3 model, with system prompts containing 21 tokens and user prompts averaging 175 tokens, depending on the input. The LLaMA 3 model processes these predictions at a rate of about 24 tokens per second. For the LLM planner, system prompts comprises 566 tokens, while user prompts averages around 150 tokens.

## B. Results

We visualize the results of all the experiments conducted on the robot using OrionNav in Fig. 7. It presents the overall results of our experiments, including individual failure cases and their corresponding causes. Out of a total of 96 trials, OrionNav successfully completed 85, resulting in a success rate exceeding 88%. Failures due to erroneous SLAM maps occurred in two trials, and a high odometry error caused failure in one trial. Excluding these three instances, the remaining 93 trials proceeded without localization issues. In this subset, the exploration algorithm failed to guide the robot to the location containing the queried objects in three cases. Among the 90 trials with successful exploration, four experienced semantic segmentation errors, where the target object was either undetected or misclassified by the perception system. Additionally, there was one instance of failure in the CBF-based navigation system, resulting in the robot colliding with an obstacle. Importantly, no failures were attributed to the LLM planner or scene graph generator, highlighting their robustness and reliability. A detailed breakdown of the success rates across various scenarios along with the average number of steps taken by the LLM planner is provided in Tab. II. Next, we offer an in-depth analysis of the system’s performance under varying conditions for three different task types, shedding light on the strengths and limitations observed during the experiments.

1) *Scenario A:* The robot and the target object are placed within the same region without occlusions between them. We conducted 36 experiments to evaluate our framework’s performance, using four distinct objects placed strategically in various locations. To account for the stochastic nature of the LLM’s outputs, we perform three trials for each target location, starting from the same initial positions. Our framework achieves a high success rate of approximately 88% in tasks related to this scenario. During the robot’s initial 360° scan, the system is able to detect most objects in the environment

and assign room labels. If the target object is observed during this initial scan, the LLM planner promptly issues a *goto*

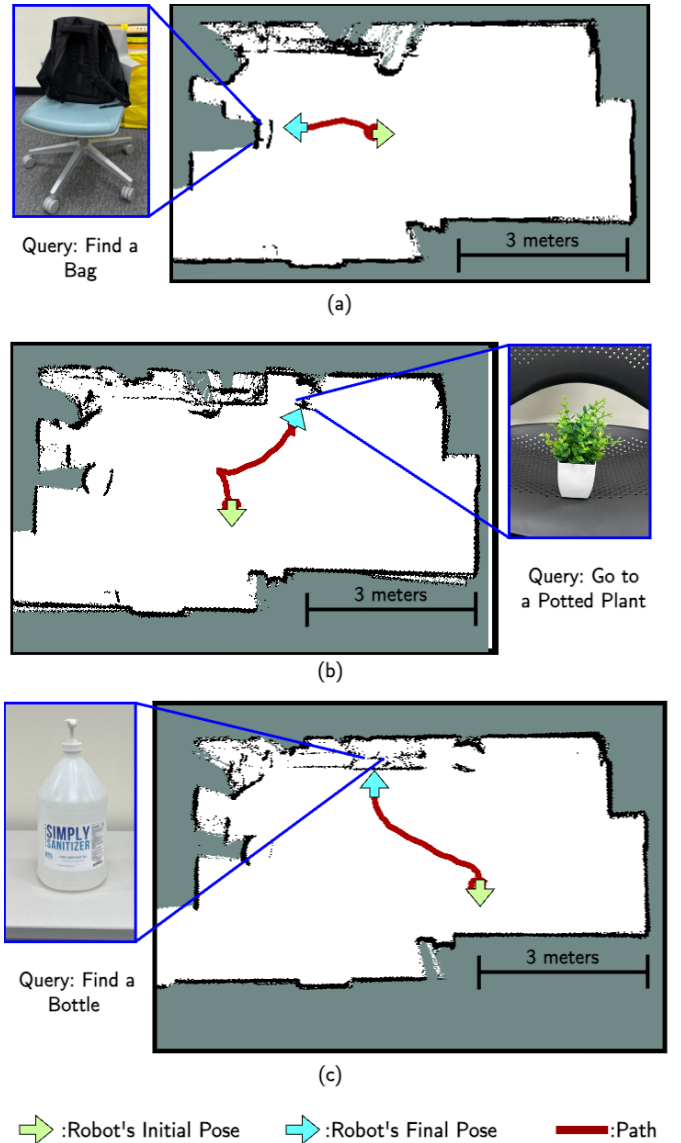


Fig. 8. **Short range object navigation task:** The robot observes the queried objects during the first 360° rotation and then navigates toward them.



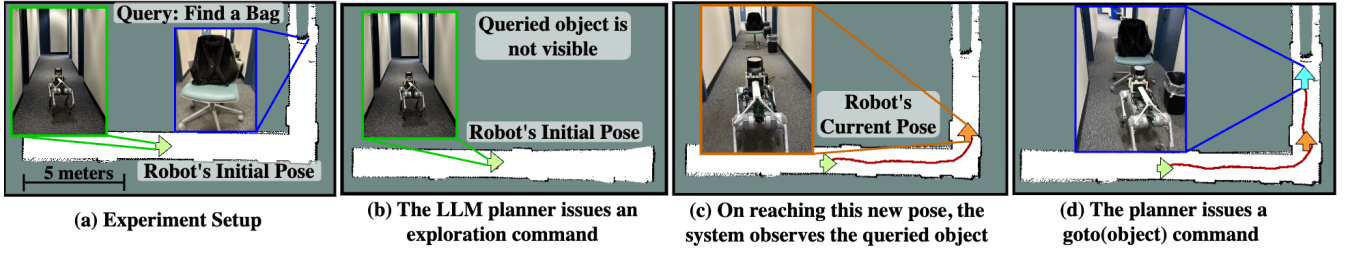


Fig. 9. **Long range object navigation task:** The system is tasked with locating a bag in a corridor environment. Initially, the bag is not visible, prompting the LLM planner to issue an exploration command. The robot explores until it detects the bag, then successfully navigates to it.

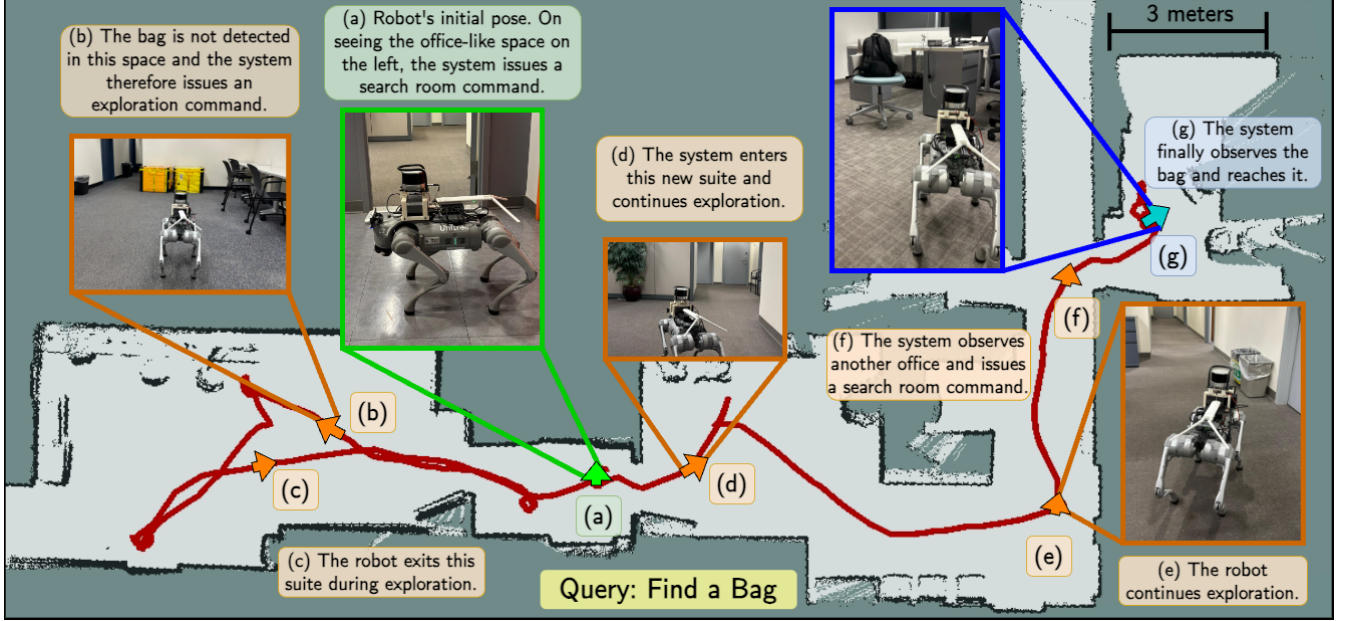


Fig. 10. **Long range object navigation task:** The system is tasked with locating a bag on a floor of a building. It begins at the starting position (a) and initially searches the office-like space (b). After failing to find the bag, the system initiates global exploration (c, d, e). Upon detecting another office space (f), the system enters and searches it (g), successfully locating the bag.

command, allowing the robot to navigate directly to the object. If the target object is not detected during the initial scan, the LLM planner issues a *search\_room* command. This adaptive approach is effective, as the target object is usually detected during the subsequent search phase, demonstrating the robustness of our exploration strategy. Three examples of this scenario are demonstrated in Fig. 8, where the system is tasked with finding a bag, a potted plant, and a bottle, respectively. In each case, the system is successfully able to locate the queried object during the 360° LiDAR scan and then navigate towards, away from the target object.

2) *Scenario B:* For this scenario, the target object is placed at a greater distance from the robot, often in separate rooms or corridors. Without prior knowledge of the environment, the robot needs to perform real-time exploration, making this scenario more challenging than Scenario A. We select different objects and position them at various locations, resulting in 15 distinct target locations. As before, we perform three trials from similar starting positions for each target location to account for the stochastic nature of the LLM’s outputs. Despite the increased complexity, our framework achieves a

high success rate of over 91%. The target object was placed at a larger distance from the robot, often located in separate rooms or corridors. Since the robot lacks prior knowledge of the environment, the LLM planner typically issued global exploration commands at the start, unlike in Scenario A. Through this exploration, the robot successfully navigated the environment to detect the target objects. Once detected, the LLM issued *goto* commands to guide the robot to the target locations. Fig. 9 depicts an example of the robot’s state during Scenario B, where it discovers a bag that is not initially visible from its starting position. After conducting a 360° scan, the LLM planner issues a global exploration command, as illustrated in Fig. 9b. The robot then follows exploration frontiers and eventually detects the bag at the position shown in Fig. 9c. At this point, the LLM planner issues a *goto* command, enabling the robot to successfully navigate to the bag. An additional example of Scenario B, conducted in the larger ‘Environment A’ from Fig. 6, is presented in Fig. 10.

3) *Scenario C:* Scenario C focuses on locating specific rooms, such as offices or break rooms. As described in Section IV-A3, our system identifies rooms using the LLM-

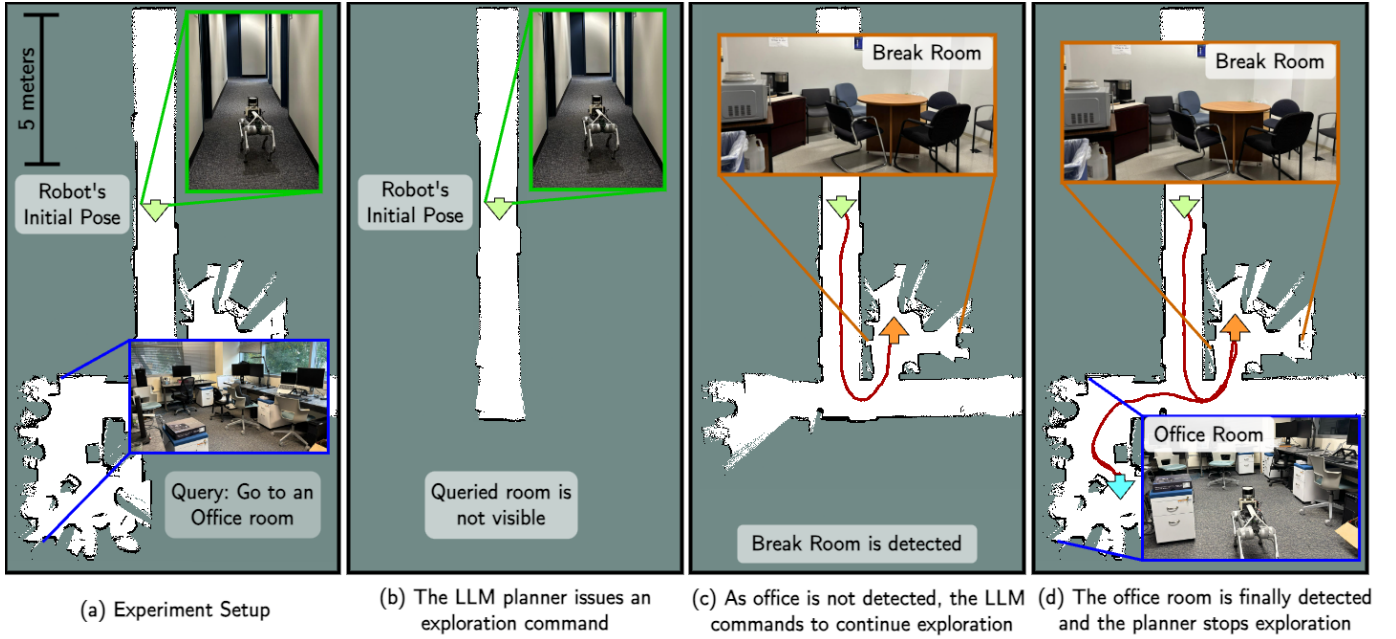


Fig. 11. **Room navigation task:** The system is tasked with locating an office room on a building floor. The LLM planner initiates the task by issuing an exploration command, during which the robot first encounters a break room. Continuing its search, the robot eventually reaches an office room which is correctly detected and classified during scene graph generation. The LLM planner then halts exploration, successfully completing the task.

based room detection module, which relies on the objects present within each room. For this set of experiments, we select four target rooms and initialize the robot from different positions. Each experiment is repeated three times with slight variations in the robot’s starting position, resulting in a total of 15 trials and 12 successes. Fig. 11 illustrates a real-world example of Scenario C, where the system is tasked with finding an office room that is not initially visible from the robot’s starting position. Consequently, the LLM planner issues an exploration command (Fig. 11b). During exploration, the robot encounters a room containing a coffee maker, microwave, refrigerator, table, and several chairs. The room detection module classifies this as a break room (Fig. 11c). Since the target has not been found, the LLM planner instructs the robot to continue exploring. The robot then enters a room equipped with multiple computers, monitors, tables, cabinets, and chairs (Fig. 11d). Upon identifying this room as an office, the system halts exploration, successfully completing the task.

TABLE II  
EXPERIMENT RESULTS: NUMBER OF SUCCESSFUL TRIALS AND THE AVERAGE NUMBER OF LLM PLANNER CALLS

Task Type	No. of Exp.	No. of Success	Avg. No. of LLM Planner Steps
<b>Scenario A:</b> Short Range Object Retrieval	36	32	1.25
<b>Scenario B:</b> Long Range Object Retrieval	45	41	2.97
<b>Scenario C:</b> Room Navigation	15	12	3.71
<b>Total</b>	<b>96</b>	<b>85</b>	

### C. Comparative Studies

We compare baseline methods and evaluate the design choices of the proposed OrionNav system, particularly in scenarios involving dynamic environmental changes. In each experiment, the robot is tasked with searching for an object that is initially absent from the environment and introduced at a later stage, for e.g., when an object is placed in a room after the robot has already mapped the area. To assess performance, we design two scenarios: one using a pre-built map and another where the robot builds the map from scratch. In both scenarios, the queried object is added to the environment after the mapping phase is complete. These setups allow us to evaluate the system’s ability to adapt to changes in the environment, both with and without prior knowledge of the map. The goal is to locate the object while minimizing the time required to complete the task. We conduct nine trials to evaluate our method and compare it with baseline approaches.

We consider *Frontier-Search* and *Object-Map-Search* baselines to demonstrate the efficacy of our framework. *Frontier-Search* explores the environment by searching through frontiers and transitions to random search once all frontiers are exhausted. This method leverages the ROS2 exploration package, as described in Section IV-C3. Since this baseline lacks an integrated mechanism to recognize when the target object has been detected, an operator manually inspects the generated object map for the presence of the queried object after the exploration concludes. The second baseline, *Object-Map-Search*, utilizes a semantic object map without room labels to locate the target object. A new action primitive, `search_object(<object name>)` is introduced for this task. It instructs the robot to navigate to the specified object in the map and perform a 360° scan to search for other nearby objects.



Fig. 12 presents the comparative results of task completion times for each experiment as bars, with tasks classified as failures when exceeding 10 minutes.

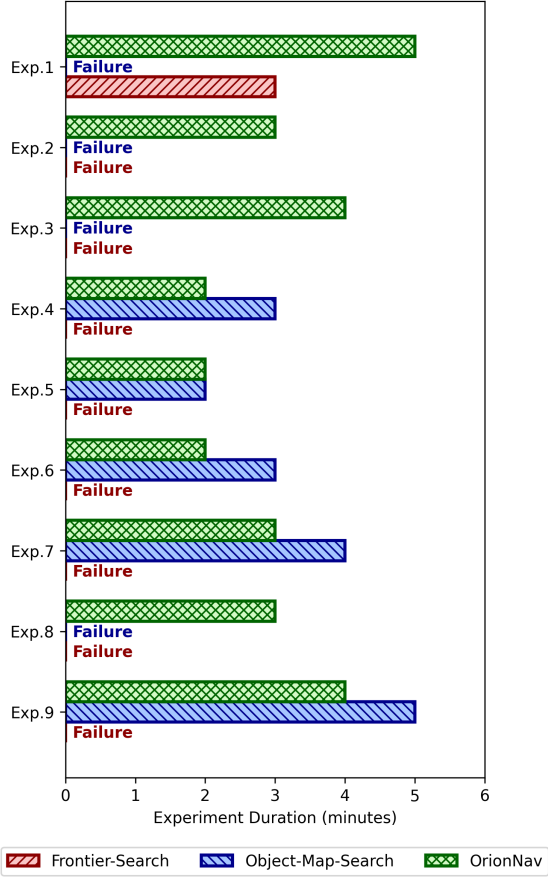


Fig. 12. **Comparisons with Baseline Methods:** The bar chart illustrates the task completion times (in minutes) for each method across various experiments. Methods that exceed 10 minutes are considered as failures.

1) *Frontier-Search*: In our experiments, the *Frontier-Search* baseline demonstrated limited effectiveness, succeeding in only one out of nine trials while failing in the other eight attempts. The single successful trial required approximately three minutes to complete. In stark contrast, OrionNav achieved success in all trials with significantly reduced search times. This superior performance of our framework can be attributed to two key components: the open-world hierarchical scene graph and the LLM-based planner. OrionNav leverages the hierarchical scene graph to maintain a continuously updated understanding of the environment, including probable object locations based on semantic information. The LLM planner analyzes this scene graph to identify the room most likely to contain the queried object and efficiently directs the system to search within that specific area, leading to consistent and rapid object localization. Conversely, the *Frontier-Search* baseline lacks semantic understanding and relies solely on frontier-based or random exploration strategies to locate objects. This approach fails to adapt its exploration strategy based on the likelihood of object locations, resulting in inefficient searches and significantly lower success rates. Importantly, even when environmental changes prevent object localization

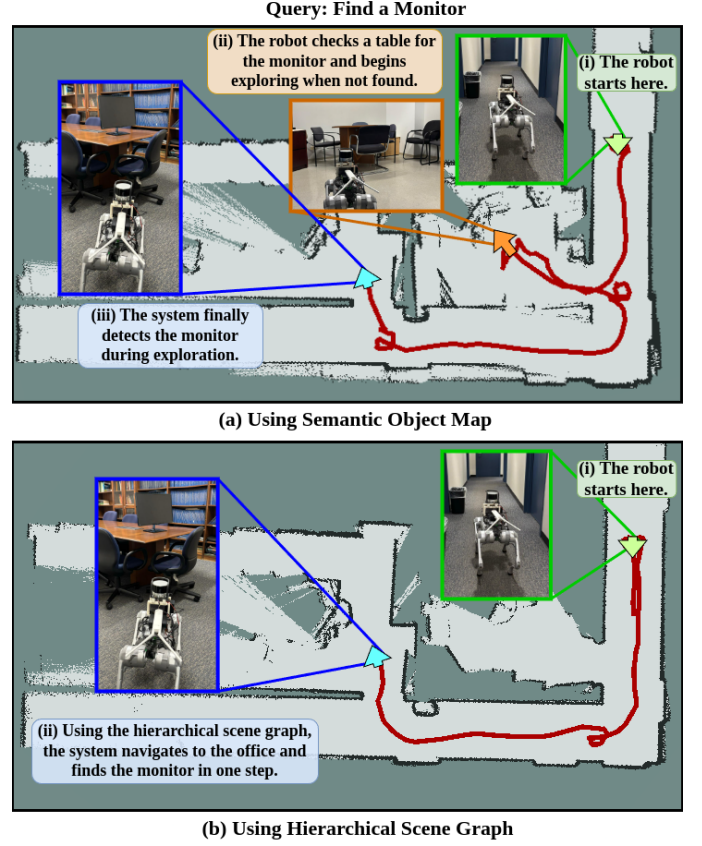


Fig. 13. **Autonomous object navigation in changing environment:** The system is provided with a partial map of the environment containing an office room among other locations. A monitor is later placed in the office room and the system is tasked with finding it. In (a), using only the semantic object map without room labels, the system searched near the table in the break room but did not find the monitor. It then explored the environment and eventually located it in the office. In (b), with the proposed hierarchical scene graph, the system identified the office and directly found the monitor inside it.

during the robot’s initial mapping phase, OrionNav’s approach enables the system to successfully locate the object in subsequent attempts. This adaptability further underscores the robustness of OrionNav in complex, changing environments.

2) *Object-Map-Search*: The key distinction between the *Object-Map-Search* baseline and OrionNav lies in the former’s use of a semantic object map without room labeling for object search. This approach results in a success rate of about 56%, significantly lower than OrionNav’s 100% success rate, underscoring the importance of the hierarchical scene graph in the proposed method. Fig. 13 illustrates this comparison, presenting a scenario where both systems are provided with a partial map of an environment comprising an office room, a break room, and a long L-shaped corridor. The task involves locating a monitor that is temporarily removed from the office room for the initial mapping phase and then reintroduced to its original location. As shown in Fig. 13a, the baseline system assumes that a monitor would likely be located on a table and navigates to the nearest one. However, this table, situated in the break room, lacks a monitor. Consequently, it initiates a global exploration command to eventually reach the office room where it detects the monitor. In contrast, OrionNav (Fig.

13b) leverages its hierarchical scene graph to deduce that the monitor is likely in an office, navigating directly to it and detecting the monitor in a single step. As shown in Fig. 12, even in successful cases, the baseline consistently required more time to complete the task compared to OrionNav.

## VI. CONCLUSION

We propose OrionNav, an autonomous system that integrates open-vocabulary hierarchical scene graphs with an LLM-based planner for real-time language-driven autonomous navigation. The system builds a map of the environment using SLAM and fuses open-world semantics, room detection, and LLM-based labeling to generate a hierarchical scene graph. The LLM planner leverages this scene graph, along with user commands, to formulate high-level plans, which are executed by a lower-level safe controller. The multi-level abstractions of perception and planning with continuously updated scene graphs and adaptive planner enables robust execution of diverse navigation tasks in large-scale complex, changing environments with the ability to recover from task failures. Extensive experiments in real-world environments demonstrated the system's high success rates and practical applicability. Future research will explore mobile robots equipped with onboard manipulation and interaction capabilities, enabling tasks such as the opening of doors and pick-and-place.

## REFERENCES

- [1] T. Zhang, X. Hu, J. Xiao, and G. Zhang, "A survey of visual navigation: From geometry to embodied AI," *Engineering Applications of Artificial Intelligence*, vol. 114, p. 105036, 2022.
- [2] J. Duan, S. Yu, H. L. Tan, H. Zhu, and C. Tan, "A survey of embodied ai: From simulators to research tasks," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 6, no. 2, pp. 230–244, 2022.
- [3] F. Safavi, P. Olikkal, D. Pei, S. Kamal, H. Meyerson, V. Penumalee, and R. Vinjamuri, "Emerging frontiers in human-robot interaction," *Journal of Intelligent and Robotic Systems*, vol. 110, no. 2, p. 45, 2024.
- [4] D.-S. Jang, D.-H. Cho, W.-C. Lee, S.-K. Ryu, B. Jeong, M. Hong, M. Jung, M. Kim, M. Lee, S. Lee, *et al.*, "Unlocking robotic autonomy: A survey on the applications of foundation models," *International Journal of Control, Automation and Systems*, vol. 22, no. 8, pp. 2341–2384, 2024.
- [5] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [6] H. Wang and M. Li, "A new era of indoor scene reconstruction: A survey," *IEEE Access*, vol. 12, pp. 110 160–110 192, 2024.
- [7] H. Yin, X. Xu, S. Lu, X. Chen, R. Xiong, S. Shen, C. Stachniss, and Y. Wang, "A survey on global lidar localization: Challenges, advances and open problems," *International Journal of Computer Vision*, vol. 132, no. 8, pp. 3139–3171, 2024.
- [8] A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone, "Kimera: From SLAM to spatial perception with 3d dynamic scene graphs," *International Journal of Robotics Research*, vol. 40, no. 12-14, pp. 1510–1546, 2021.
- [9] Y. Tao, X. Liu, I. Spasojevic, S. Agarwal, and V. Kumar, "3d active metric-semantic SLAM," *IEEE Robotics and Automation Letters*, vol. 9, no. 3, pp. 2989–2996, 2024.
- [10] Y. Wang, Y. Tian, J. Chen, K. Xu, and X. Ding, "A survey of visual SLAM in dynamic environment: The evolution from geometric to semantic approaches," *IEEE Transactions on Instrumentation and Measurement*, vol. 73, pp. 1–21, 2024.
- [11] C. Zhu and L. Chen, "A survey on open-vocabulary detection and segmentation: Past, present, and future," *CoRR*, vol. abs/2307.09220, 2023.
- [12] J. Wu, X. Li, S. Xu, H. Yuan, H. Ding, Y. Yang, X. Li, J. Zhang, Y. Tong, X. Jiang, B. Ghanem, and D. Tao, "Towards open vocabulary learning: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 7, pp. 5092–5113, 2024.
- [13] S. Koch, N. Vaskevicius, M. Colosi, P. Hermosilla, and T. Ropinski, "Open3DSG: Open-vocabulary 3d scene graphs from point clouds with queryable objects and open-set relationships," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 14 183–14 193.
- [14] A. Takmaz, E. Fedele, R. W. Sumner, M. Pollefeys, F. Tombari, and F. Engelmann, "Openmask3d: Open-vocabulary 3d instance segmentation," *arXiv preprint arXiv:2306.13631*, 2023.
- [15] D. Maggio, Y. Chang, N. Hughes, M. Trang, D. Griffith, C. Dougherty, E. Cristofalo, L. Schmid, and L. Carlone, "Clio: Real-time task-driven open-set 3d scene graphs," *IEEE Robotics and Automation Letters*, vol. 9, no. 10, pp. 8921–8928, 2024.
- [16] A. Werby, C. Huang, M. Büchner, A. Valada, and W. Burgard, "Hierarchical Open-Vocabulary 3D Scene Graphs for Language-Grounded Robot Navigation," in *Proceedings of the Robotics: Science and Systems*, Delft, Netherlands, July 2024.
- [17] K. Rana, J. Haviland, S. Garg, J. Abou-Chakra, I. Reid, and N. Suen-derhauf, "Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning," in *Proceedings of the Conference on Robot Learning*, Atlanta, GA, USA, November 2023, pp. 23–72.
- [18] A. Rajvanshi, K. Sikka, X. Lin, B. Lee, H.-P. Chiu, and A. Velasquez, "Saynav: Grounding large language models for dynamic planning to navigation in new environments," in *Proceedings of the International Conference on Automated Planning and Scheduling*, Banff, AB, Canada, June 2024, pp. 464–474.
- [19] P. Wu, Y. Mu, B. Wu, Y. Hou, J. Ma, S. Zhang, and C. Liu, "Voronav: Voronoi-based zero-shot object navigation with large language model," *arXiv preprint arXiv:2401.02695*, 2024.
- [20] D. Honerkamp, M. Büchner, F. Despinoy, T. Welschehold, and A. Valada, "Language-grounded dynamic scene graphs for interactive object search with mobile manipulation," *IEEE Robotics and Automation Letters*, vol. 9, no. 10, pp. 8298–8305, 2024.
- [21] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, "ProgPrompt: Generating situated robot task plans using large language models," in *Proceedings of the International Conference on Robotics and Automation*, London, United Kingdom, May 2023, pp. 11 523–11 530.
- [22] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, "Code as policies: Language model programs for embodied control," in *Proceedings of the International Conference on Robotics and Automation*, London, United Kingdom, May 2023, pp. 9493–9500.
- [23] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar, *et al.*, "Inner monologue: Embodied reasoning through planning with language models," in *Proceedings of the Conference on Robot Learning*, Atlanta, GA, USA, December 2023, pp. 1769–1782.
- [24] Q. Gu, A. Kuwajerwala, S. Morin, K. M. Jatavallabhula, B. Sen, A. Agarwal, C. Rivera, W. Paul, K. Ellis, R. Chellappa, *et al.*, "Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning," in *Proceedings of the International Conference on Robotics and Automation*, Yokohama, Japan, May 2024, pp. 5021–5028.
- [25] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, *et al.*, "Do as I can, not as I say: Grounding language in robotic affordances," in *Proceedings of the Conference on Robot Learning*, Auckland, New Zealand, December 2022, pp. 287–318.
- [26] C. Huang, O. Mees, A. Zeng, and W. Burgard, "Visual language maps for robot navigation," in *Proceedings of the International Conference on Robotics and Automation*, London, United Kingdom, May 2023, pp. 10 608–10 615.
- [27] S.-C. Wu, J. Wald, K. Tateno, N. Navab, and F. Tombari, "SceneGraph-Fusion: Incremental 3d scene graph prediction from rgb-d sequences," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, June 2021, pp. 7511–7521.
- [28] Y. Deng, J. Wang, J. Zhao, X. Tian, G. Chen, Y. Yang, and Y. Yue, "OpenGraph: Open-vocabulary hierarchical 3d graph representation in large-scale outdoor environments," *IEEE Robotics and Automation Letters*, vol. 9, no. 10, pp. 8402–8409, 2024.
- [29] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik, "LERF: language embedded radiance fields," in *Proceedings of the*

- International Conference on Computer Vision*, Paris, France, October 2023, pp. 19672–19682.
- [30] H. Zhang, F. Li, and N. Ahuja, “Open-nerf: Towards open vocabulary nerf decomposition,” in *Proceedings of the Winter Conference on Applications of Computer Vision*, Waikoloa, HI, USA, January 2024, pp. 3444–3453.
  - [31] Y. Wang, H. Chen, and G. H. Lee, “Gov-nesf: Generalizable open-vocabulary neural semantic fields,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 20443–20453.
  - [32] T. Nguyen, A. Bourki, M. Macudzinski, A. Brunel, and M. Bannamoun, “Semantically-aware neural radiance fields for visual scene understanding: A comprehensive review,” *CoRR*, vol. abs/2402.11141, 2024.
  - [33] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al., “Visual genome: Connecting language and vision using crowdsourced dense image annotations,” *International Journal of Computer Vision*, vol. 123, pp. 32–73, 2017.
  - [34] B. Schroeder and S. Tripathi, “Structured query-based image retrieval using scene graphs,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition Workshops*, Seattle, WA, USA, June 2020, pp. 178–179.
  - [35] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei, “Image retrieval using scene graphs,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, Boston, MA, USA, June 2015, pp. 3668–3678.
  - [36] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, “Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs,” in *Proceedings of the International Conference on Robotics and Automation*, Xi’an, China, May 2021, pp. 6541–6548.
  - [37] M. Xu, M. Qu, B. Ni, and J. Tang, “Joint modeling of visual objects and relations for scene graph generation,” *Proceedings of the Advances in Neural Information Processing Systems*, pp. 7689–7702, December 2021.
  - [38] M. Qi, W. Li, Z. Yang, Y. Wang, and J. Luo, “Attentive relational networks for mapping images to scene graphs,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, Long Beach, CA, USA, June 2019, pp. 3957–3966.
  - [39] J. Yang, J. Lu, S. Lee, D. Batra, and D. Parikh, “Graph r-cnn for scene graph generation,” in *Proceedings of the European Conference on Computer Vision*, Munich, Germany, September 2018, pp. 670–685.
  - [40] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning transferable visual models from natural language supervision,” in *Proceedings of the International Conference on Machine Learning*, vol. 139, Vienna, Austria, July 2021, pp. 8748–8763.
  - [41] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat, et al., “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
  - [42] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” in *Proceedings of the Advances in Neural Information Processing Systems*, New Orleans, LA, USA, December 2023.
  - [43] H. Liu, C. Li, Y. Li, and Y. J. Lee, “Improved baselines with visual instruction tuning,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 26296–26306.
  - [44] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” in *Proceedings of the International Conference on Machine Learning*, Honolulu, HI, USA, July 2023, pp. 19730–19742.
  - [45] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W. Lo, P. Dollár, and R. B. Girshick, “Segment anything,” in *Proceedings of the International Conference on Computer Vision*, Paris, France, October 2023, pp. 3992–4003.
  - [46] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese, “3D scene graph: A structure for unified semantics, 3d space, and camera,” in *Proceedings of the International Conference on Computer Vision*, Seoul, Korea, October 2019, pp. 5664–5673.
  - [47] A. Rosinol, A. Gupta, M. Abate, J. Shi, and L. Carlone, “3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans,” in *Proceedings of the Robotics: Science and Systems*, Corvallis, Oregon, USA, July 2020.
  - [48] N. Hughes, Y. Chang, and L. Carlone, “Hydra: A Real-time Spatial Perception System for 3D Scene Graph Construction and Optimization,” in *Proceedings of the Robotics: Science and Systems*, New York City, NY, USA, June 2022.
  - [49] J. Huang, S. Yong, X. Ma, X. Linghu, P. Li, Y. Wang, Q. Li, S. Zhu, B. Jia, and S. Huang, “An embodied generalist agent in 3d world,” in *Proceedings of the International Conference on Machine Learning*, Vienna, Austria, July 2024.
  - [50] L. Xia, J. Cui, R. Shen, X. Xu, Y. Gao, and X. Li, “A survey of image semantics-based visual simultaneous localization and mapping: Application-oriented solutions to autonomous navigation of mobile robots,” *International Journal of Advanced Robotic Systems*, vol. 17, no. 3, 2020.
  - [51] W. Xu, M. Liu, O. Sokolsky, I. Lee, and F. Kong, “Llm-enabled cyber-physical systems: Survey, research opportunities, and challenges,” in *IEEE International Workshop on Foundation Models for Cyber-Physical Systems & Internet of Things (FMSys)*, Hong Kong, China, May 2024, pp. 50–55.
  - [52] K. Bousmalis, G. Vezzani, D. Rao, C. Devin, A. X. Lee, M. Bauza, T. Davchev, Y. Zhou, A. Gupta, A. Raju, et al., “Robocat: A self-improving foundation agent for robotic manipulation,” *arXiv preprint arXiv:2306.11706*, 2023.
  - [53] H. Huang, Z. Wang, R. Huang, L. Liu, X. Cheng, Y. Zhao, T. Jin, and Z. Zhao, “Chat-scene: Bridging 3d scene and large language models with object identifiers,” in *Proceedings of the Advances in Neural Information Processing Systems*, Vancouver, BC, Canada, December 2024.
  - [54] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, “Object detection in 20 years: A survey,” *Proceedings of the IEEE*, vol. 111, no. 3, pp. 257–276, 2023.
  - [55] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos, “Image segmentation using deep learning: A survey,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 7, pp. 3523–3542, 2021.
  - [56] H. Fu, N. Patel, P. Krishnamurthy, and F. Khorrami, “Clipscope: Enhancing zero-shot ood detection with bayesian scoring,” *arXiv e-prints*, pp. arXiv–2405, 2024.
  - [57] M. Xu, Z. Zhang, F. Wei, Y. Lin, Y. Cao, H. Hu, and X. Bai, “A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model,” in *Proceedings of the European Conference on Computer Vision*, Tel Aviv, Israel, October 2022, pp. 736–753.
  - [58] H. Yuan, X. Li, C. Zhou, Y. Li, K. Chen, and C. C. Loy, “Open-vocabulary SAM: segment and recognize twenty-thousand classes interactively,” in *Proceedings of the European Conference on Computer Vision*, Milan, Italy, September 2024.
  - [59] J. Chen, Q. Yu, X. Shen, A. L. Yuille, and L. Chen, “ViTamin: Designing scalable vision models in the vision-language era,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 12954–12966.
  - [60] B. Xie, J. Cao, J. Xie, F. S. Khan, and Y. Pang, “Sed: A simple encoder-decoder for open-vocabulary semantic segmentation,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 3426–3436.
  - [61] S. Cho, H. Shin, S. Hong, A. Arnab, P. H. Seo, and S. Kim, “CAT-seg: Cost aggregation for open-vocabulary semantic segmentation,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 4113–4123.
  - [62] W. Kang, G. Liu, M. Shah, and Y. Yan, “Segvg: Transferring object bounding box to segmentation for visual grounding,” in *Proceedings of the European Conference on Computer Vision*, Milan, Italy, September 2024.
  - [63] Y. Wang, R. Sun, N. Luo, Y. Pan, and T. Zhang, “Image-to-image matching via foundation models: A new perspective for open-vocabulary semantic segmentation,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 3952–3963.
  - [64] J. Wang and L. Ke, “Llm-seg: Bridging image segmentation and large language model reasoning,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 1765–1774.
  - [65] H. Shi, S. D. Dao, and J. Cai, “Llmformer: Large language model for open-vocabulary semantic segmentation,” *International Journal of Computer Vision*, August 2024.
  - [66] X. Ma, S. Yong, Z. Zheng, Q. Li, Y. Liang, S. Zhu, and S. Huang, “SQA3D: situated question answering in 3d scenes,” in *Proceedings of the International Conference on Learning Representations*, Kigali, Rwanda, May 2023.
  - [67] D. Azuma, T. Miyanishi, S. Kurita, and M. Kawanabe, “Scanqa: 3d question answering for spatial scene understanding,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, New Orleans, LA, USA, June 2022, pp. 19107–19117.
  - [68] P. Achlioptas, A. Abdelreheem, F. Xia, M. Elhoseiny, and L. J. Guibas, “Referit3d: Neural listeners for fine-grained 3d object identification

- in real-world scenes,” in *Proceedings of the European Conference in Computer Vision*, Glasgow, UK, August 2020, pp. 422–440.
- [69] Y. Chen, S. Yang, H. Huang, T. Wang, R. Lyu, R. Xu, D. Lin, and J. Pang, “Grounded 3d-llm with referent tokens,” *arXiv preprint arXiv:2405.10370*, 2024.
- [70] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan, “3d-llm: Injecting the 3d world into large language models,” in *Proceedings of the Advances in Neural Information Processing Systems*, New Orleans, LA, USA, December 2023.
- [71] S. Chen, X. Chen, C. Zhang, M. Li, G. Yu, H. Fei, H. Zhu, J. Fan, and T. Chen, “LI3da: Visual interactive instruction tuning for omni-3d understanding reasoning and planning,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, June 2024, pp. 26 428–26 438.
- [72] J. Yang, X. Chen, S. Qian, N. Madaan, M. Iyengar, D. F. Fouhey, and J. Chai, “LLM-Grounder: Open-vocabulary 3d visual grounding with large language model as an agent,” in *Proceedings of the International Conference on Robotics and Automation*, Yokohama, Japan, May 2024, pp. 7694–7701.
- [73] B. Xu, W. Li, D. Tzoumanikas, M. Bloesch, A. J. Davison, and S. Leutenegger, “Mid-fusion: Octree-based object-level multi-instance dynamic SLAM,” in *Proceedings of the International Conference on Robotics and Automation*, Montreal, QC, Canada, May 2019, pp. 5231–5237.
- [74] M. Rünz, M. Buffier, and L. Agapito, “Maskfusion: Real-time recognition, tracking and reconstruction of multiple moving objects,” in *Proceedings of the International Symposium on Mixed and Augmented Reality*, D. Chu, J. L. Gabbard, J. Grubert, and H. Regenbrecht, Eds., Munich, Germany, October 2018, pp. 10–20.
- [75] L. Nicholson, M. Milford, and N. Sünderhauf, “Quadricslam: Dual quadrics from object detections as landmarks in object-oriented SLAM,” *IEEE Robotics and Automation Letters*, vol. 4, no. 1, pp. 1–8, 2019.
- [76] N. Patel, P. Krishnamurthy, and F. Khorrami, “Semantic segmentation guided slam using vision and lidar,” in *Proceedings of the International Symposium on Robotics*, Munich, German, June 2018, pp. 1–7.
- [77] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. J. Kelly, and A. J. Davison, “SLAM++: simultaneous localisation and mapping at the level of objects,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Portland, OR, USA, June 2013, pp. 1352–1359.
- [78] X. Kong, S. Liu, M. Taher, and A. J. Davison, “vmap: Vectorised object mapping for neural field SLAM,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Vancouver, BC, Canada, June 2023, pp. 952–961.
- [79] N. Patel, F. Khorrami, P. Krishnamurthy, and A. Tzes, “Tightly coupled semantic RGB-D inertial odometry for accurate long-term localization and mapping,” in *Proceedings of the International Conference on Advanced Robotics*, Belo Horizonte, Brazil, December 2019, pp. 523–528.
- [80] X. Han, H. Liu, Y. Ding, and L. Yang, “RO-MAP: real-time multi-object mapping with neural radiance fields,” *IEEE Robotics and Automation Letters*, vol. 8, no. 9, pp. 5950–5957, 2023.
- [81] J. McCormac, R. Clark, M. Bloesch, A. J. Davison, and S. Leutenegger, “Fusion++: Volumetric object-level SLAM,” in *Proceedings of the International Conference on 3D Vision*, Verona, Italy, September 2018, pp. 32–41.
- [82] R. Ding, J. Yang, C. Xue, W. Zhang, S. Bai, and X. Qi, “PLA: language-driven open-vocabulary 3d scene understanding,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Vancouver, BC, Canada, June 2023, pp. 7010–7019.
- [83] Z. Jin, M. Hayat, Y. Yang, Y. Guo, and Y. Lei, “Context-aware alignment and mutual masking for 3d-language pre-training,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Vancouver, BC, Canada, June 2023, pp. 10984–10994.
- [84] Z. Zhu, X. Ma, Y. Chen, Z. Deng, S. Huang, and Q. Li, “3d-vista: Pre-trained transformer for 3d vision and text alignment,” in *Proceedings of the International Conference on Computer Vision*, Paris, France, October 2023, pp. 2911–2921.
- [85] S. Chen, H. Zhu, M. Li, X. Chen, P. Guo, Y. Lei, Y. Gang, T. Li, and T. Chen, “Vote2cap-detr++: Decoupling localization and describing for end-to-end 3d dense captioning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 11, pp. 7331–7347, 2024.
- [86] D. Z. Chen, A. Gholami, M. Nießner, and A. X. Chang, “Scan2cap: Context-aware dense captioning in RGB-D scans,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Nashville, TN, USA, June 2021, pp. 3193–3203.
- [87] D. Z. Chen, A. X. Chang, and M. Nießner, “Scanrefer: 3d object localization in RGB-D scans using natural language,” in *Proceedings of the European Conference on Computer Vision*, Glasgow, UK, August 2020, pp. 202–221.
- [88] J. Schult, F. Engelmann, A. Hermans, O. Litany, S. Tang, and B. Leibe, “Mask3d: Mask transformer for 3d semantic instance segmentation,” in *Proceedings of the International Conference on Robotics and Automation*, London, UK, May 2023, pp. 8216–8223.
- [89] Y. Yue, S. Mahadevan, J. Schult, F. Engelmann, B. Leibe, K. Schindler, and T. Kontogianni, “AGILE3D: attention guided interactive multi-object 3d segmentation,” in *Proceedings of the International Conference on Learning Representations*, Vienna, Austria, May 2024.
- [90] J. Zhou, J. Wang, B. Ma, Y.-S. Liu, T. Huang, and X. Wang, “Uni3d: Exploring unified 3d representation at scale,” in *Proceedings of the International Conference on Learning Representations*, Vienna, Austria, May 2024.
- [91] X. Chen, Z. Ma, X. Zhang, S. Xu, S. Qian, J. Yang, D. F. Fouhey, and J. Chai, “Multi-object hallucination in vision-language models,” in *Proceedings of the Advances in Neural Information Processing Systems*, Vancouver, BC, Canada, December 2024.
- [92] Q. Yu, J. He, X. Deng, X. Shen, and L. Chen, “Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional CLIP,” in *Proceedings of the Advances in Neural Information Processing Systems*, New Orleans, LA, USA, December 2023.
- [93] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, New Orleans, LA, USA, June 2022, pp. 11 966–11 976.
- [94] W. Cai, S. Huang, G. Cheng, Y. Long, P. Gao, C. Sun, and H. Dong, “Bridging zero-shot object navigation and foundation models through pixel-guided navigation skill,” in *Proceedings of the International Conference on Robotics and Automation*, Yokohama, Japan, May 2024, pp. 5228–5234.
- [95] P. Liu, Y. Orru, C. Paxton, N. M. M. Shafullah, and L. Pinto, “OK-Robot: What really matters in integrating open-knowledge models for robotics,” *arXiv preprint arXiv:2401.12202*, 2024.
- [96] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girshick, “Masked-attention mask transformer for universal image segmentation,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition*, New Orleans, LA, USA, June 2022, pp. 1280–1289.
- [97] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, Berlin, Germany, August 2016, pp. 1715–1725.
- [98] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: common objects in context,” in *Proceedings of the European Conference on Computer Vision*, vol. 8693, Zurich, Switzerland, September 2014, pp. 740–755.
- [99] H. Caesar, J. Uijlings, and V. Ferrari, “Coco-stuff: Thing and stuff classes in context,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, June 2018, pp. 1209–1218.
- [100] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, “Panoptic segmentation,” in *Proceedings of the International Conference on Computer Vision and Pattern Recognition*, Long Beach, CA, USA, June 2019, pp. 9404–9413.
- [101] S. Macenski and I. Jambrecic, “SLAM toolbox: SLAM for the dynamic world,” *Journal of Open Source Software*, vol. 6, no. 61, p. 2783, 2021. [Online]. Available: <https://doi.org/10.21105/joss.02783>
- [102] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al., “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *Proceedings of the International Conference on Knowledge Discovery and Data Mining*, Portland, OR, USA, August 1996, pp. 226–231.
- [103] H. W. Kuhn, “The hungarian method for the assignment problem,” *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [104] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al., “The llama 3 herd of models,” *arXiv preprint arXiv:2407.21783*, 2024.
- [105] B. Dai, R. Khorrambakht, P. Krishnamurthy, and F. Khorrami, “Sailing through point clouds: Safe navigation using point cloud based control barrier functions,” *IEEE Robotics and Automation Letters*, vol. 9, no. 9, pp. 7731–7738, 2024.
- [106] J. Hörner. (2016) Map-merging for multi-robot system. Prague. [Online]. Available: <https://is.cuni.cz/webapps/zzp/detail/174125/>