

STRIVE: Structured Representation Integrating VLM Reasoning for Efficient Object Navigation

Haokun Zhu^{1*}, Zongtai Li^{1*}, Zhixuan Liu¹,
Wenshan Wang¹, Ji Zhang¹, Jonathan Francis^{1,2}, Jean Oh¹
¹Carnegie Mellon University ²Bosch Center for AI

Abstract: Vision-Language Models (VLMs) have been increasingly integrated into object navigation tasks for their rich prior knowledge and strong reasoning abilities. However, applying VLMs to navigation poses two key challenges: effectively representing complex environment information and determining *when and how* to query VLMs. Insufficient environment understanding and over-reliance on VLMs (e.g. querying at every step) can lead to unnecessary backtracking and reduced navigation efficiency, especially in continuous environments. To address these challenges, we propose a novel framework that constructs a multi-layer representation of the environment during navigation. This representation consists of viewpoint, object nodes, and room nodes. Viewpoints and object nodes facilitate intra-room exploration and accurate target localization, while room nodes support efficient inter-room planning. Building on this representation, we propose a novel two-stage navigation policy, integrating high-level planning guided by VLM reasoning with low-level VLM-assisted exploration to efficiently locate a goal object. We evaluated our approach on three simulated benchmarks (HM3D, RoboTHOR, and MP3D), and achieved state-of-the-art performance on both the success rate ($\uparrow 7.1\%$) and navigation efficiency ($\uparrow 12.5\%$). We further validate our method on a real robot platform, demonstrating strong robustness across 15 object navigation tasks in 10 different indoor environments. Project page is available at [here](#).

Keywords: Robot Object Navigation, Vision-Language Models, Reasoning

1 Introduction

Object navigation is a fundamental task in robotics, where an agent must locate an instance of a given object category in unknown environments. This task is particularly challenging, as it requires the agent to understand complex visual information, reason about spatial relationships, and make decisions based on both current and past observations.

Advances in Vision-Language Models (VLMs) [1, 2, 3] have demonstrated strong capabilities in contextual visual understanding and common-sense reasoning. Building on this, recent works [4, 5, 6, 7, 8, 9] have integrated VLMs into object navigation tasks, utilizing their rich prior knowledge, visual understanding, and commonsense reasoning abilities to guide navigation. However, existing approaches often face two significant challenges: First, the input to VLMs typically lacks a structured representation of the environment and is often restricted to local observations. Without a coherent global view that integrates both current and previous observations, VLMs struggle to reason effectively about the environment and fail to make reasonable navigation decisions. Second, existing methods [6, 7, 9] typically rely on VLMs to select among all frontier viewpoints at each step, without utilizing navigation progress or environment layouts to effectively guide VLMs’ reasoning process. Besides, due to VLMs’ limited understanding of 3D spatial information [10, 11, 12], they cannot jointly reason about the spatial relationships and the navigation history when evaluating each viewpoint. As a result, their evaluation of viewpoints is largely based on viewpoints’ local semantic information, which often leads to redundant navigation behaviors such as backtracking or repeated exploration.

* Equal Contribution.

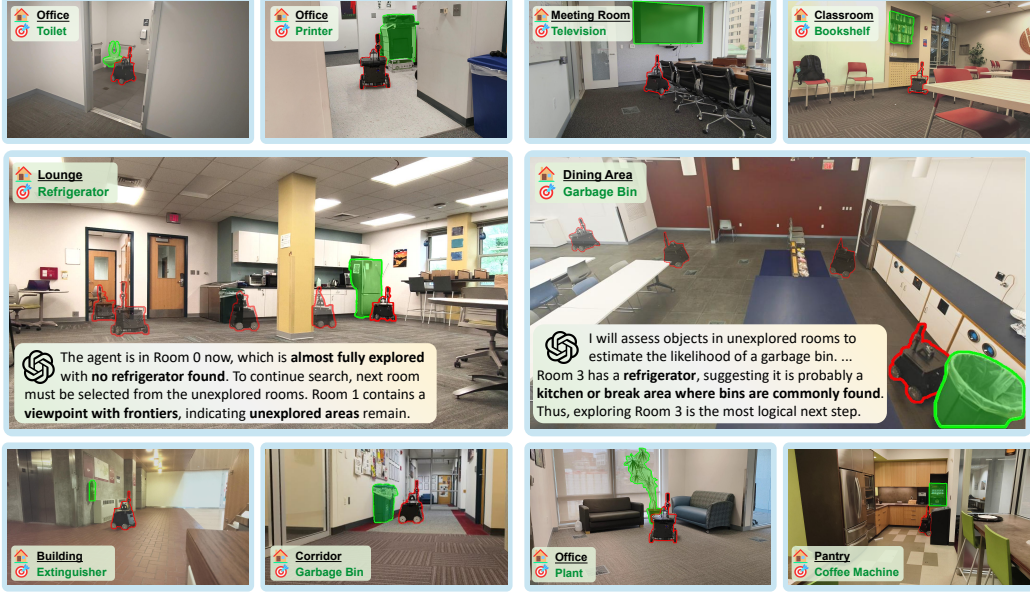


Figure 1: STRIVE can conduct zero-shot object navigation in diverse and complex real-world environments by leveraging our novel multi-layer representation and 2-stage efficient navigation policy.

To address these challenges, we propose STRIVE (S**TR**uctured R**E**presentation I**N**tegrating VLM R**E**asoning for E**FF**icient O**B**ject N**E**avigation), a novel framework that incrementally learns a structured representation of the environment and utilizes VLM’s reasoning abilities to guide the navigation. This representation consists of 3 layers: object nodes, viewpoint nodes, and room nodes. Object nodes represent all observed objects, provide comprehensive semantic information about the environment and assist in target localization; Viewpoint nodes discretize the environment into a set of key locations, enabling efficient intra-room exploration; Room nodes further segment the environment into distinct rooms and facilitate room-level reasoning by the VLM. This multi-layer representation enables a more comprehensive understanding of the environment, allowing VLM to better utilize their reasoning abilities for more effective decision-making. Furthermore, we design an efficient two-stage navigation policy based on this representation, combining high-level planning guided by the VLM’s reasoning and VLM-assisted low-level exploration. Specifically, for the high-level planning, instead of making step-by-step decisions among all viewpoint nodes, the VLM selects the next room to explore based on the spatial layout and semantic information of each room. For low-level exploration within rooms, we employ a traditional frontier-based algorithm for efficient exploration, while leveraging the VLM to decide whether continued exploration of the current room is worthwhile. Making high-level planning on rooms effectively mitigates the issue of VLMs’ insufficient 3D spatial understanding and prevents redundant actions, thereby enhancing navigation efficiency.

We evaluate our method on three widely-used simulated benchmarks: HM3D [13], RoboTHOR [14], and MP3D [15]. STRIVE achieves state-of-the-art (SOTA) results, significantly outperforming 13 existing methods in both Success Rate (SR) and navigation efficiency, measured by Success weighted by Path Length (SPL). This highlights the effectiveness of our proposed multi-layer representation and the VLM-guided reasoning policy in improving object navigation. Specifically, STRIVE achieves 79.6% SR and 38.7% SPL on HM3D, 68.1% SR and 36.3% SPL on RoboTHOR, and 52.3% SR and 23.1% SPL on MP3D. Besides, we also conduct 15 real-world experiments across 10 different indoor environments on a Mecanum wheel platform [16], demonstrating the effectiveness and robustness of our method in real-world scenarios.

2 Related Works

Object Navigation. Existing object navigation methods are typically categorized into end-to-end learning approaches and modular approaches. End-to-end methods [17, 18, 19, 20, 21, 22, 23] use reinforcement learning to directly map observations to actions, but often suffer from low sample

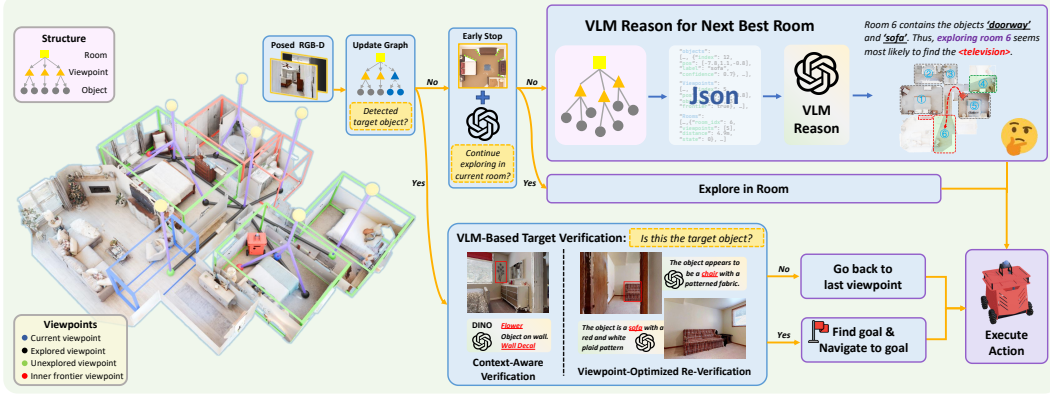


Figure 2: **Overview of STRIVE.** We construct a multi-layer representation \mathcal{R} (Sec. 3.1) on-the-fly, consisting of object, viewpoint, and room nodes, which serves as a structured input for VLM. Based on \mathcal{R} , we introduce a two-stage navigation policy, where the VLM reasons and plans at room-level (Sec. 3.2.2), while the agent explores in room at the viewpoint-level using a VLM-assisted frontier-based navigation strategy (Sec. 3.2.1) and VLM-based target verification (Sec. 3.2.3).

efficiency and poor generalization. In contrast, modular methods [24, 25, 26, 27, 28, 24, 6, 4] decompose navigation into steps such as mapping, planning, and action execution, and often build semantic maps in bird’s-eye view or 3D space to facilitate more interpretable and scalable navigation behavior. With the emergence of foundation models [1], object navigation has advanced towards zero-shot, open-vocabulary setting [27, 28, 29, 4, 6]. VLFM [29] aligns object goals with CLIP embeddings, while CogNav [6] further leverages LLMs to enable cognitive-like decision-making. We also leverage VLM’s reasoning abilities to improve zero-shot object navigation. but we employ a novel representation and a two-stage policy, enabling more efficient and effective VLM guidance.

VLM-guided Navigation. With internet-scale training data, Vision-Language Models (VLMs) [30, 1, 2, 3] have shown strong common-sense reasoning abilities and have been widely applied in Object Navigation tasks to guide the decision-making process. For example, InstructNav [5] leverages multi-sourced value maps to model key navigation elements. SG-Nav [4] constructs 3D scene graphs and prompts LLMs with structural relationships, while CogNav [6] utilizes LLMs to reason about the cognitive process of object navigation. However, due to VLMs’ limited 3D spatial understanding ability [10, 11, 12], over-reliance on VLMs for navigation can lead to inefficient behavior, such as frequent backtracking. To address this, we propose a two-stage navigation policy that combines VLM-guided high-level planning with VLM-assisted frontier-based low-level exploration strategies, leveraging the reasoning strength of VLMs while ensuring efficient and robust navigation behavior.

Scene Representation for Indoor Navigation. Scene representation is crucial for transforming raw observations into structured information for decision-making in navigation tasks. Frontier-based methods [25, 31, 32, 33] record frontiers on a grid map and integrate semantic information to guide exploration. In contrast, graph-based methods represent the environment as structured scene graphs to support navigation. Prior works [4, 9] use scene graphs to summarize semantic information and let VLMs to select among frontier locations. Others [34, 7, 6] explicitly construct viewpoints in the scene graph, enabling VLMs to reason over the graph and choose among viewpoints to guide navigation. Unlike traditional scene graphs [35, 7, 6], where viewpoints are typically derived from Voronoi partitions, we discretize the environment into semantically meaningful regions to select viewpoint nodes. As the middle layer, these viewpoints bridge the spatial structure (room nodes) and semantic content (object nodes), forming a structured representation facilitating VLM reasoning.

3 Method

Task Definition: In Object Navigation, the agent is required to find an instance of a given object category (e.g. Find the *bed*.) in an unknown environment. At each time step t , the agent receives a posed RGB-D observation $\mathbf{O}_t = \{I_t, D_t, P_t = \langle \mathbf{p}_t, \mathbf{R}_t \rangle\}$, where I_t is the RGB image,

D_t is the depth map, and P_t is the camera pose. The navigation policy then predicts an action $a_t \in \{\text{move_forward}, \text{turn_left}, \text{turn_right}, \text{stop}\}$. The task is considered successful if the agent stops within d_s meters of the target object in less than T steps.

Overview: Fig. 2 provides an overview of STRIVE, a framework that constructs a multi-layer environment representation and performs object navigation through a novel two-stage navigation policy. STRIVE enables the VLM to reason at the room-level while guiding the agent to explore within rooms at the viewpoint-level. The representation construction process is detailed in Sec. 3.1 and the two-stage navigation policy is presented in Sec. 3.2.

3.1 Multi-layer Environment Representation

We propose a framework that constructs a three-layer graph representation \mathcal{R} to model the environment, where each layer corresponds to a specific type of node: object nodes $V^{obj} = \{v_i^{obj}\}$, viewpoint nodes $V^{vp} = \{v_i^{vp}\}$, and room nodes $V^{room} = \{v_i^{room}\}$. Edges encode spatial and semantic relationships across nodes. We elaborate on the graph construction in following sections.

3.1.1 Viewpoint Nodes

Inspired by [36], we construct a skeleton graph as the viewpoint layer to discretize the environment. Importantly, the graph is incrementally built as the agent navigates—each time the agent reaches a new viewpoint, it updates the graph. We define a coverage range ζ_{cover} as the radius within which semantic information is associated with the center viewpoint. Each viewpoint node thus controls a local region determined by ζ_{cover} . Edges between viewpoint nodes indicate direct traversability. The maximum sensor range ζ_{max} denotes the effective measurable distance of the depth camera.

Node: To select the future viewpoint nodes, we begin by taking the input of agent’s position \mathbf{p} , coverage range ζ_{cover} , and accumulated point cloud (project from posed D). First, k rays are cast from \mathbf{p} in k uniformly sampled directions to intersect with point cloud and form a polygon P within ζ_{cover} , representing the controlled region of current viewpoint node. Next, we remove the polygon from the point cloud and divide the remaining point cloud into separate regions, which fall into two categories: regions with or without frontiers.

1) Regions with frontiers: Inspired by [37], we use frontiers to guide viewpoint selection. The frontiers are identified and clustered into frontier edge segments and construct a graph $G_{frontier}$ by connecting every pair of segments that are mutually visible. Then the *Maximum Clique* is iteratively removed from $G_{frontier}$. For each removed clique, its center is added as a new viewpoint node v_i^{vp} (green node in Fig. 3) to our representation \mathcal{R} .

2) Regions without frontiers: For these regions, the center of the region is directly added as a new viewpoint node v_i^{vp} (yellow nodes in Fig. 3) to our representation \mathcal{R} .

Edges between V^{vp} : We evaluate straight-line traversability between each pair of viewpoint nodes. An edge is added if the direct line between two nodes is free of obstacles.

3.1.2 Object Nodes

We leverage open-vocabulary detection and segmentation methods [38, 39] to obtain segmented 3D object instances. Specifically, given the observations at time step t , we reconstruct the 3D point cloud of each detected object using the predicted masks, depth map D_t and camera pose P_t . For each object, we instantiate an object node at its center, recording attributes such as 3D position,

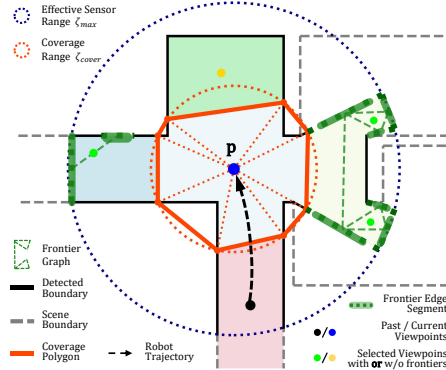


Figure 3: **Visualization of the viewpoint selection algorithm.** Green and yellow nodes are the selected viewpoints.

point cloud, predicted label, confidence score and 3D bounding box. Newly instantiated nodes are merged with previously observed nodes if they correspond to the same physical object.

Edges between V^{vp} and V^{obj} : An edge is added between v_i^{vp} and v_j^{obj} if v_j^{obj} is within the cover range ζ_{cover} of v_i^{vp} and is visible from v_i^{vp} . An object can be associated with multiple viewpoints. If an object isn't connected to any viewpoint, we connect it to the closest visible viewpoint.

3.1.3 Room Nodes

Following [40, 35], we identify all walls in the environment and iteratively dilate them to segment the environment into connected components. Then each connected component is added as a room node v_i^{room} to our representation \mathcal{R} . Finally, edges are added between each room node and the viewpoint nodes located within the corresponding room. Further details are provided in the App. A8.

3.2 Object Navigation Policy

In this section, we present our efficient two-stage navigation policy, where the VLM performs high-level reasoning and planning at the room level, while the agent conducts fine-grained exploration within each room at the viewpoint level, guided by a VLM-assisted frontier-based strategy and VLM-based target verification.

3.2.1 Explore in Room with Early Stop

For efficient low-level exploration within rooms, we introduce VLM-assisted early stop, combining VLM with traditional frontier-based algorithm. We first classify frontiers into two types: *True Frontiers*, which lie along room boundaries indicating incomplete exploration, and *Inner Frontiers*, resulting from objects' occlusions. The agent iteratively navigates to the nearest viewpoint with *True Frontiers* and explores until all *True Frontiers* are cleared. If *Inner Frontiers* still remain in the current room, we query the VLM to decide whether further exploration inside this room is necessary.

3.2.2 Next Best Room

In situations where exploration of the current room is completed *without* finding the target object, we must determine the next room to explore. To guide this decision, we leverage VLM's commonsense reasoning abilities by providing task-relevant context and general exploration heuristics.

The task-relevant context is consolidated into a combined *Prompt* as described in Fig. 4, which contains 1) target object instruction, 2) agent's current state, 3) agent's navigation history, and 4) the environment representation \mathcal{R} formatted as a JSON file. A detailed description of the JSON format is provided in the App. A3.

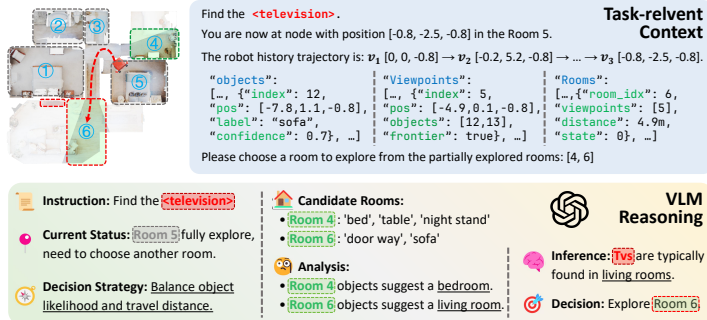


Figure 4: Visualization of the structured prompt and the VLM's reasoning process of selecting the next best room.

Besides the task-relevant context, we also provide the VLM with general exploration heuristics. Specifically, we explicitly instruct the VLM to evaluate two factors: 1) The semantic similarity between the objects in each room and the target object. 2) The distance from the agent's current position to each room, aiming to optimize the exploration path by minimizing unnecessary backtracking. Using a Chain-of-Thought reasoning strategy, the VLM selects the most suitable unexplored room for further exploration. Finally, the viewpoint closest to the current position in the selected room is chosen as the next action viewpoint. We detail the prompts used in the App. A2.

Table 1: Comparison with SOTA methods with different settings on HM3D, RoboTHOR, and MP3D datasets. We report the Success Rate (SR) and Success weighted by Path Length (SPL) metrics.

Method	Open-Set	Zero-Shot	HM3D		RoboTHOR		MP3D	
			SR(%) \uparrow	SPL(%) \uparrow	SR(%) \uparrow	SPL(%) \uparrow	SR(%) \uparrow	SPL(%) \uparrow
SemEXP [24]	\times	\times	-	-	-	-	36.0	14.4
PONI [25]	\times	\times	-	-	-	-	31.8	12.1
ZSON [41]	\checkmark	\times	25.5	12.6	-	-	15.3	4.8
L3MVN [27]	\times	\checkmark	54.2	25.5	41.2	22.5	34.9	14.5
ESC [28]	\checkmark	\checkmark	39.2	22.3	38.1	22.2	28.7	11.2
VoroNav [7]	\checkmark	\checkmark	42.0	26.0	-	-	-	-
VLFM [29]	\checkmark	\checkmark	52.5	30.4	-	-	36.4	17.5
SG-Nav [4]	\checkmark	\checkmark	54.0	24.9	47.5	24.0	40.2	16.0
OpenFMNav [42]	\checkmark	\checkmark	54.9	22.4	44.1	23.3	37.2	15.7
InstructNav [5]	\checkmark	\checkmark	58.0	20.9	-	-	-	-
TriHelper [43]	\checkmark	\checkmark	62.0	25.3	-	-	-	-
OSG [9]	\checkmark	\checkmark	69.3	28.3	-	-	-	-
CogNav [6]	\checkmark	\checkmark	72.5	26.2	54.6	24.3	46.6	16.1
STRIVE (ours)	\checkmark	\checkmark	79.6	38.7	68.1	36.3	52.3	23.1

Notably, in later navigation stage, continuing forward is more effective than backtracking, as remaining steps may not allow long detours. In light of this, we introduce a penalized distance that weights the geodesic distance by factors reflecting the steps already taken and the number of explored viewpoints along the path to each candidate room.

3.2.3 VLM-Based Target Verification

Accurate detection of the target object is crucial in object navigation. However, relying solely on detectin model [38] often results in false positives. To address this, we propose incorporating the VLM to verify detected target objects, leveraging its ability to reason about the contextual information of the surrounding environment.

Context-Aware Verification: When the agent detects a potential target object, we prompt the VLM with the detected object and its surrounding visual context for verification. The VLM leverages both the object’s appearance and its surrounding semantic information to determine whether it matches the target category, e.g. recognizing a painting of plant as a ‘decoration’ rather than ‘plant’.

Viewpoint-Optimized Re-Verification: The agent may initially observe and detect the target object from a suboptimal viewpoint (e.g., under occlusion or from a long distance), resulting in inaccurate detection. To address this, we perform a second observation from a better viewpoint. Unlike the baseline method [4, 6], which further observes the target object from multiple viewpoints, we compute the optimal viewpoint along the path from the current position to the target object and only perform one VLM re-verification at that viewpoint. This strategy improves detection accuracy without sacrificing navigation efficiency.

4 Experiment

4.1 Experiment Setup

We evaluate our method by comparing it with state-of-the-art methods on the Object Navigation task in Habitat [44] simulator. We also conduct real-world experiments across diverse environments.

Dataset: We perform simulated experiments on 3 datasets: 1) HM3D [13], a large-scale 3D indoor scene dataset comprising 20 high-fidelity scenes with 6 target object categories. 2) MP3D [15], another 3D scene dataset featuring 11 high-fidelity scenes with 21 target object categories. 3) RoboTHOR [14], a 3D indoor scene dataset containing 15 scenes with 12 target object categories.

Evaluation Metrics: Following [45], we use 4 metrics to evaluate the performance: 1) Success Rate (SR): the percentage of episodes in which the agent reaches the target object within a success distance. 2) Success weighted by Path Length (SPL): the success rate weighted by the ratio of the shortest path length to the actual path length. 3) Distance to Goal (DTG): the final distance to the

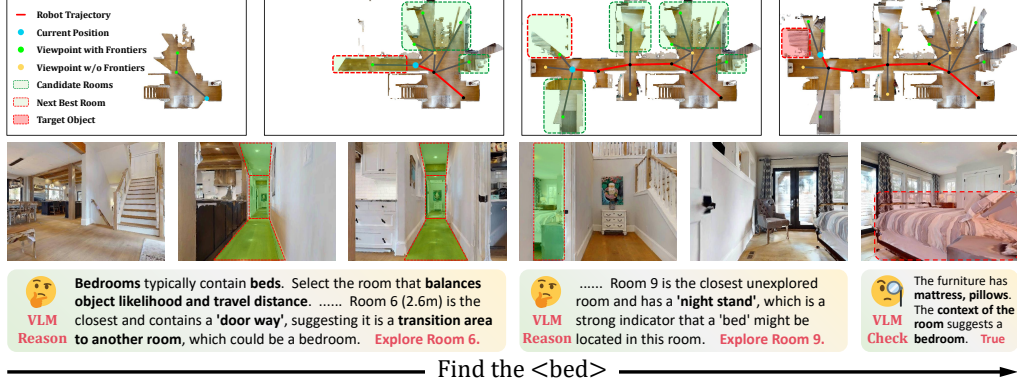


Figure 5: **Qualitative visualization of STRIVE.** The first and second steps show the VLM’s reasoning process, where it selects Room 6 and 9 by jointly considering room-layout (‘doorway’), semantic cues (‘nightstand’) and travel cost (penalized distance). The final step shows VLM-based verification, using contextual cues (e.g., mattress, pillows) to confirm the target object as a ‘bed’.

target object at the end of the episode. 4) SoftSPL: replacing the binary success term in SPL with a “soft” value that indicates the progress made by the agent towards the goal.

Implementation Details: Each episode allows a maximum of 500 steps and a success distance of 1.0m. The agent observes the environment using a 640×480 RGB-D image, with depth values from 0.5m to 5.0m and a horizontal field of view (HFOV) of 79° . The camera is mounted 0.88m above the ground. The agent moves forward by 0.25m per step and rotates by 30° . For VLM, we use Gemini [3] (gemini-2.0-flash), and for object detection and segmentation, we use MM-GroundingDINO [38] and SAM [39]. All experiments are conducted on RTX4090 GPUs.

For real-world experiments, we deploy STRIVE on the Mecanum wheel platform [16], which is equipped with a Ricoh Theta Z1 360-degree camera for RGB image capturing and a Livox Mid-360 LiDAR for 3D point cloud acquisition. To maintain compatibility with the input format in simulation, the collected point clouds are converted into depth maps when necessary.

4.2 Quantitative Results in Simulator

We compare STRIVE with state-of-the-art object navigation methods in different settings as shown in Tab. 1. STRIVE significantly outperforms all baselines across all benchmarks, with an increase of +7.1% SR, +12.5% SPL in HM3D, +13.5% SR, +12.0% SPL in RoboTHOR, +5.7% SR, +7.0% SPL in MP3D compared to CogNav [6]. The improvement in SPL is more significant than that in SR, indicating that our representation and navigation policy effectively improve navigation efficiency. Furthermore, increased navigation efficiency enables the agent to explore a larger area within a limited number of steps. The improvement in SR results from the combined effects of more efficient navigation, better utilization of the VLM’s reasoning capabilities, and more accurate VLM-based verification.

4.3 Qualitative Results in Simulator

We visualize the navigation process of STRIVE in Fig. 5. The results demonstrate that our structured representation enables the VLM to reason effectively about both spatial layout (e.g., a room with a “doorway” can lead to other rooms) and semantic cues (e.g., nightstands suggesting bedrooms), leading to improved room selection. Furthermore, the VLM balances the likelihood of finding the target object against travel distance cost when planning room-to-room exploration. It also leverages contextual information to re-verify detected objects and effectively reduces false positives.

4.4 Real-World Experiments

We conduct 15 real-world experiments across 10 different environments, including offices, meeting rooms, classrooms, lounges, dining areas, corridors, and kitchens. Part of the experiment environments and results are shown in Fig. 1. Compared to simulation, real-world deployment presents

additional challenges. Lidar-captured point clouds are much sparser than depth maps, and real environments are often more cluttered, introducing noise that affects both exploration and object detection. Despite these challenges, our agent demonstrates robust performance.

We elaborate on 2 difficult environments in Fig. 1. In the left scenario, the agent is initialized inside a small enclosed room connected to a larger lounge area. Despite the presence of inner frontiers—regions occluded by furniture—the agent correctly decides to abort exhaustive exploration of this room. Instead, it exits the room early and shifts its attention to unexplored rooms nearby, where it ultimately locates the target object. In the right scenario, the agent is initialized in a dining area and instructed to find a ‘Garbage bin’. STRIVE successfully uses the VLM to reason on semantic associations (‘refrigerator’ and ‘bins’) and find the target object efficiently. For other environments, we also present final frames where the agent successfully reaches the target object. From these results, we can conclude that STRIVE can handle diverse and complex real-world environments.

4.5 Ablation Study

Multi-layer Representation: We conduct an ablation study to evaluate the contributions of object nodes and room nodes in our multi-layer scene representation. Since room-level navigation depends on room nodes, it cannot be used when they are removed. To ensure consistency, we adopt a basic navigation policy that allows the VLM to guide navigation at the viewpoint level instead of the room level. As shown in Tab. 2, both object nodes and room nodes contribute significantly to performance improvement. Incorporating object nodes improves the agent’s ability to localize target objects, while adding room nodes enhances its understanding of the environment layout. Using both layers together leads to a substantial improvement over using either individually.

Navigation Policy: We conduct ablation studies to evaluate the effectiveness of our navigation policy, VLM-assisted early stopping, penalized distance, and VLM-based verification, as summarized in Tab. 3. The last two rows compare our room-level planning policy against a basic viewpoint-level approach. Results show that room-level planning enables the agent to better leverage the VLM’s reasoning capabilities, significantly boosting performance. We also report the average VLM token usage per episode. By querying the VLM only for room-level planning, our method significantly reduces token consumption compared to viewpoint-level planning. Finally, the VLM-assisted early stopping, penalized distance, and verification modules each contribute to further performance gains.

Table 2: **Ablation study of representation on HM3D.** We adopt a *viewpoint-level navigation policy* for experiment consistency.

V^{vp}	V^{obj}	V^{room}	SR \uparrow	SPL \uparrow	S-SPL \uparrow	DTG(m) \downarrow
✓	✗	✗	71.3	33.2	35.2	1.86
✓	✓	✗	72.4	34.0	36.1	1.95
✓	✗	✓	72.9	33.8	35.4	1.86
✓	✓	✓	75.0	34.9	36.5	1.80

Table 3: **Ablation study of navigation policy components on HM3D.** Viewpoint Policy stands for VLM planning on viewpoint-level.

	SR \uparrow	SPL \uparrow	S-SPL \uparrow	DTG(m) \downarrow	Tokens \downarrow
w/o Early Stop	74.8	34.8	36.4	1.62	-
w/o Penalized Dist	73.7	36.1	36.9	1.47	-
w/o VLM-Verify	72.1	32.7	34.1	1.83	-
Viewpoint Policy	75.0	34.9	36.5	1.80	22935
STRIVE	79.6	38.7	38.9	1.29	8068

5 Conclusion

In this paper, we introduce STRIVE, a novel framework that incrementally constructs a structured scene representation and leverages VLM’s reasoning capabilities to achieve efficient object navigation. STRIVE incrementally builds a multi-layer representation of the environment, consisting of room, viewpoint and object nodes. Based on this representation, we design an efficient two-stage VLM-guided navigation policy, which leverages VLM reasoning for room-level planning while using VLM together with traditional frontier-based methods for efficient exploration within rooms. To further improve robustness, we incorporate VLM-based target verification, utilizing VLMs’ contextual understanding to improve detection accuracy. Experiments across three simulated benchmarks demonstrate that STRIVE achieves state-of-the-art performance, significantly improving both success rate and navigation efficiency. Furthermore, our real-world experiments demonstrate the robustness and practicality of STRIVE in navigating complex and diverse real-world environments.

6 Limitations

Despite the promising results in object navigation, STRIVE still has several limitations:

Limiting Assumption: In simulation, the depth input from the depth camera is dense and accurate. However, in real-world settings, although we have accumulated LiDAR inputs along the trajectory, the resulting point clouds remain significantly sparser, which affects both object segmentation and traversability estimation. As a result, the agent must pause at each viewpoint for 2 seconds to accumulate denser point clouds. Incorporating lightweight point cloud completion models to improve perception quality without sacrificing efficiency might address this issue.

Failure Mode: In our experiments, we observed that even with VLM-based target verification, it is still challenging to avoid false positives in certain scenarios. For instance, when searching for a bed, if the agent encounters a sofa bed, the VLM may mistakenly identify it as a bed due to the functional and visual similarities. Future work could explore using more constrained prompts and stronger VLMs for target verification. Additionally, at the start of the task, employing image generation models to create a target image could help mitigate the influence of VLM’s common-sense knowledge on target verification since generative models typically produce very typical target images.

Other Limitations:

(1) STRIVE currently does not support real-time updates of the environment representation, which becomes especially evident on real-world experiments with limited onboard computation. The bottleneck lies primarily in the detection module: current 2D detection model, MM-Grounding-DINO [38], is both slow and error-prone. To improve detection reliability, we incorporate VLM-based target verification, which further increases computational overhead. A more efficient and accurate detection framework would significantly improve the system’s speed and responsiveness.

(2) In simulated environments, we found that many scenes have big holes in the mesh, which can lead to incorrect traversability estimation. This is primarily due to the data collection. Besides, we found that the dataset didn’t label all instances of the target categories in the scene, which cause some success episodes to be counted as failures, especially in MP3D [15]. More details is provides in the App. A7. We believe that the dataset should be improved to provide more accurate and complete annotations.

References

- [1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.
- [2] T. Wu, G. Yang, Z. Li, K. Zhang, Z. Liu, L. Guibas, D. Lin, and G. Wetzstein. Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22227–22238, 2024.
- [3] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [4] H. Yin, X. Xu, Z. Wu, J. Zhou, and J. Lu. Sg-nav: Online 3d scene graph prompting for llm-based zero-shot object navigation. *Advances in Neural Information Processing Systems*, 37: 5285–5307, 2024.
- [5] Y. Long, W. Cai, H. Wang, G. Zhan, and H. Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. In *8th Annual Conference on Robot Learning*.
- [6] Y. Cao, J. Zhang, Z. Yu, S. Liu, Z. Qin, Q. Zou, B. Du, and K. Xu. Cognav: Cognitive process modeling for object goal navigation with llms. *arXiv preprint arXiv:2412.10439*, 2024.

- [7] 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.
- [8] S. Saxena, B. Buchanan, C. Paxton, B. Chen, N. Vaskevicius, L. Palmieri, J. Francis, and O. Kroemer. Grapheqa: Using 3d semantic scene graphs for real-time embodied question answering. *arXiv preprint arXiv:2412.14480*, 2024.
- [9] J. Loo, Z. Wu, and D. Hsu. Open scene graphs for open-world object-goal navigation. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*.
- [10] Z. Zhang, F. Hu, J. Lee, F. Shi, P. Kordjamshidi, J. Chai, and Z. Ma. Do vision-language models represent space and how? evaluating spatial frame of reference under ambiguities. *arXiv preprint arXiv:2410.17385*, 2024.
- [11] B. Chen, Z. Xu, S. Kirmani, B. Ichter, D. Sadigh, L. Guibas, and F. Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024.
- [12] Z. Qi, Z. Zhang, Y. Fang, J. Wang, and H. Zhao. Gpt4scene: Understand 3d scenes from videos with vision-language models. *arXiv preprint arXiv:2501.01428*, 2025.
- [13] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. M. Turner, E. Undersander, W. Galuba, A. Westbury, A. X. Chang, M. Savva, Y. Zhao, and D. Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. URL <https://arxiv.org/abs/2109.08238>.
- [14] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford, L. Weihs, M. Yatskar, and A. Farhadi. Robothor: An open simulation-to-real embodied ai platform. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [15] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017.
- [16] J. Zhang. Autonomy stack for mecanum wheel platform. https://github.com/jizhang-cmu/autonomy_stack_mecanum_wheel_platform, 2024. Accessed: 2025-04-29.
- [17] R. Dang, L. Wang, Z. He, S. Su, J. Tang, C. Liu, and Q. Chen. Search for or navigate to? dual adaptive thinking for object navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8250–8259, 2023.
- [18] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. *arXiv preprint arXiv:1911.00357*, 2019.
- [19] J. Ye, D. Batra, A. Das, and E. Wijmans. Auxiliary tasks and exploration enable objectgoal navigation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16117–16126, 2021.
- [20] R. Ramrakhya, E. Undersander, D. Batra, and A. Das. Habitat-web: Learning embodied object-search strategies from human demonstrations at scale. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5173–5183, 2022.
- [21] R. Ramrakhya, D. Batra, E. Wijmans, and A. Das. Pirlnav: Pretraining with imitation and rl finetuning for objectnav. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17896–17906, 2023.

- [22] A. Mousavian, A. Toshev, M. Fišer, J. Košecká, A. Wahid, and J. Davidson. Visual representations for semantic target driven navigation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8846–8852. IEEE, 2019.
- [23] W. Yang, X. Wang, A. Farhadi, A. Gupta, and R. Mottaghi. Visual semantic navigation using scene priors. *arXiv preprint arXiv:1810.06543*, 2018.
- [24] D. S. Chaplot, D. P. Gandhi, A. Gupta, and R. R. Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. *Advances in Neural Information Processing Systems*, 33: 4247–4258, 2020.
- [25] S. K. Ramakrishnan, D. S. Chaplot, Z. Al-Halah, J. Malik, and K. Grauman. Poni: Potential functions for objectgoal navigation with interaction-free learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18890–18900, 2022.
- [26] J. Zhang, L. Dai, F. Meng, Q. Fan, X. Chen, K. Xu, and H. Wang. 3d-aware object goal navigation via simultaneous exploration and identification, 2023. URL <https://arxiv.org/abs/2212.00338>.
- [27] B. Yu, H. Kasaei, and M. Cao. L3myn: Leveraging large language models for visual target navigation. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, page 3554–3560. IEEE, Oct. 2023. doi:10.1109/iros55552.2023.10342512. URL <http://dx.doi.org/10.1109/IROS55552.2023.10342512>.
- [28] K. Zhou, K. Zheng, C. Pryor, Y. Shen, H. Jin, L. Getoor, and X. E. Wang. Esc: Exploration with soft commonsense constraints for zero-shot object navigation, 2023. URL <https://arxiv.org/abs/2301.13166>.
- [29] N. Yokoyama, S. Ha, D. Batra, J. Wang, and B. Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 42–48. IEEE, 2024.
- [30] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. URL <https://arxiv.org/abs/2201.12086>.
- [31] J. Chen, G. Li, S. Kumar, B. Ghanem, and F. Yu. How to not train your dragon: Training-free embodied object goal navigation with semantic frontiers. *arXiv preprint arXiv:2305.16925*, 2023.
- [32] S. Y. Gadre, M. Wortsman, G. Ilharco, L. Schmidt, and S. Song. Cows on pasture: Baselines and benchmarks for language-driven zero-shot object navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23171–23181, 2023.
- [33] T. Gervet, S. Chintala, D. Batra, J. Malik, and D. S. Chaplot. Navigating to objects in the real world. *Science Robotics*, 8(79):eadf6991, 2023.
- [34] D. An, H. Wang, W. Wang, Z. Wang, Y. Huang, K. He, and L. Wang. Etpnav: Evolving topological planning for vision-language navigation in continuous environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [35] N. Hughes, Y. Chang, and L. Carlone. Hydra: A real-time spatial perception system for 3D scene graph construction and optimization. 2022.
- [36] F. Yang, D.-H. Lee, J. Keller, and S. Scherer. Graph-based topological exploration planning in large-scale 3d environments. In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 12730–12736. IEEE, 2021.

- [37] B. Yamauchi. A frontier-based approach for autonomous exploration. In *Proceedings 1997 IEEE International Symposium on Computational Intelligence in Robotics and Automation CIRA'97: Towards New Computational Principles for Robotics and Automation*, pages 146–151. IEEE, 1997.
- [38] X. Zhao, Y. Chen, S. Xu, X. Li, X. Wang, Y. Li, and H. Huang. An open and comprehensive pipeline for unified object grounding and detection. *arXiv preprint arXiv:2401.02361*, 2024.
- [39] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023.
- [40] A. Werby, C. Huang, M. Büchner, A. Valada, and W. Burgard. Hierarchical open-vocabulary 3d scene graphs for language-grounded robot navigation. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024.
- [41] A. Majumdar, G. Aggarwal, B. Devnani, J. Hoffman, and D. Batra. Zson: Zero-shot object-goal navigation using multimodal goal embeddings. *Advances in Neural Information Processing Systems*, 35:32340–32352, 2022.
- [42] Y. Kuang, H. Lin, and M. Jiang. Openfmnav: Towards open-set zero-shot object navigation via vision-language foundation models. *arXiv preprint arXiv:2402.10670*, 2024.
- [43] L. Zhang, Q. Zhang, H. Wang, E. Xiao, Z. Jiang, H. Chen, and R. Xu. Trihelper: Zero-shot object navigation with dynamic assistance. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10035–10042. IEEE, 2024.
- [44] X. Puig, E. Undersander, A. Szot, M. D. Cote, R. Partsey, J. Yang, R. Desai, A. W. Clegg, M. Hlavac, T. Min, T. Gervet, V. Vondruš, V.-P. Berges, J. Turner, O. Maksymets, Z. Kira, M. Kalakrishnan, J. Malik, D. S. Chaplot, U. Jain, D. Batra, A. Rai, and R. Mottaghi. Habitat 3.0: A co-habitat for humans, avatars and robots, 2023.
- [45] K. Yadav, J. Krantz, R. Ramrakhya, S. K. Ramakrishnan, J. Yang, A. Wang, J. Turner, A. Gokaslan, V.-P. Berges, R. Mootaghi, O. Maksymets, A. X. Chang, M. Savva, A. Clegg, D. S. Chaplot, and D. Batra. Habitat challenge 2023. <https://aihabitat.org/challenge/2023/>, 2023.

STRIVE: Structured Representation Integrating VLM Reasoning for Efficient Object Navigation

(Appendix)

A1 Overview

In this supplementary material, more details about the proposed STRIVE and more experimental results are provided, including:

- **General Exploration Heuristic:** The general exploration heuristic prompt we provide to the VLM to help it make better decisions (App. A2).
- **Detailed Structure of Task-relevant Context:** The detailed structure of the Task-relevant Context generated from our proposed representation \mathcal{R} to prompt the VLM with whole environment information for better reasoning (App. A3).
- **Detailed Experiment Results:** We provide more detailed experimental results, including the performance of STRIVE on different categories of objects on HM3D [13] and more qualitative results (App. A4).
- **Examples of VLM Reasoning:** We provide more examples of VLM reasoning results, including the reasoning process and the final decision (App. A5).
- **Details of VLM-based Verification:** We provide a detailed explanation of the VLM-based verification process, including Context-Aware Verification and Viewpoint-Optimized Re-Verification, and we also provide more examples of VLM-based verification results, including the reasoning process and the final decision (App. A6).
- **Dataset Error:** We briefly show some examples of mislabelled data in the dataset (App. A7).
- **Detailed Room Segmentation:** We provide more details about the room segmentation process (App. A8).

A2 Exploration Heuristics

In order to make the VLM better guide the navigation to complete the current task as soon as possible, we provide a general prompt to give it an overall concept of the object navigation task. We also explain the meaning of each part in the provided JSON file and explicitly require the VLM to consider the probability of the target object appearing in each candidate room and the travel cost required to explore each room when making decisions.

Listing 1: General Object Navigation Heuristic Prompt

```
PROMPT = """
You are a wheeled mobile robot operating in an indoor environment. Your goal is
to efficiently find a target object based on a human-provided instruction in a
new house. The current room you are in has been fully explored. To achieve the
goal, you must select the next room to explore from the partially explored rooms
listed in a JSON file, aiming to complete the task as quickly as possible.

### Provided Information:
1. A specific instruction describing the task.
2. A description of your current position and previous trajectories.
3. A JSON file containing details about the scene, including rooms, viewpoints,
and objects.

The JSON file contains the following information:
- **Objects**
```

```

- 'object_idx': A unique identifier for the object.
- 'position': The spatial coordinates of the object.
- 'class': The category or type of the object.
- 'confidence': The confidence level of the classification result.
- 'size': The bounding box size of the object (in meters).

- **Viewpoints**:
  - 'viewpoint_idx': A unique identifier for the viewpoint.
  - 'position': The spatial coordinates of the viewpoint.
  - 'state': The state of the viewpoint ('1' for visited, '0' for unvisited).
  - 'neighbors': A list of connected viewpoints.
  - 'has_frontier': Relevant only when the viewpoint is unvisited.
  - 'True': The viewpoint has a frontier, meaning unknown regions exist around it.
  - 'False': The area around the viewpoint has already been observed from distant viewpoints, but small objects may still be unclear.
  - 'objects': A list of objects observable from the viewpoint.

- **Rooms**:
  - 'room_idx': A unique identifier for the room.
  - 'state': The state of the room ('1' for fully explored, '0' for partially explored).
  - 'distance': The distance (in meters) the robot needs to travel to reach this room.
  - 'viewpoints': A list of viewpoints in the room.

### Task:
You must carefully analyze the JSON file, using logical reasoning and common sense, to select the next room to explore from the list of partially explored rooms. Consider the following factors:
- Evaluate how closely each room's viewpoints aligns with the overall task objective.
- Optimize the exploration path by leveraging the robot's current momentum and minimizing unnecessary backtracking or redundant movements.
- Assess the likelihood that exploring the selected room will meaningfully advance or complete the overall task.

### Output Format:
Your response should include:
- 'steps': The chain of thought leading to the decision.
- 'final_answer': The 'idx' of the next room to explore.
- 'reason': The rationale for selecting this room.

**Note:** The chosen room must be partially explored.
"""

```

A3 Json Structure

Here we provide a detailed description of the task-related context used to prompt the VLM about the current navigation process and the currently known environmental information.

As shown in Fig. 6, it consists of the following parts:

- **Target Object:** It begins by specifying the target object as "Find the <target object>".
- **Current Viewpoint and Position:** It then states the robot's current viewpoint and position as "The robot is now at Viewpoint with position [x,y,z] in Room r_i ".
- **Navigation History:** The navigation history up to the current step is provided as "The robot history trajectory is Position [x,y,z] → , Position [x,y,z] → ...".

- **Scene Representation:** The scene representation \mathcal{R} is formatted as a JSON file as the last part. This representation contains information about the layout and semantic information of the environment, which is crucial for the VLM to make informed decisions. The JSON file is structured as the format of Rooms-Viewpoints-Objects. For detailed json structure, please refer to Fig. 6.

```
Prompt_info: Find the <toilet>.
You are now at node with position [-1.017, -0.126, -0.8] in the Room 1.
The robot history trajectory is: Position [ 0.0, 0.0, -0.8] --> Position [-1.017, -0.126, -0.8]

"objects":
[
  {
    "object_idx": 0,
    "class": "door",
    "position": [-1.274, -1.083, -0.8],
    "confidence": 0.391,
    "size": [0.078, 0.853, 2.053]
  },
  {
    "object_idx": 1,
    "class": "luggage",
    "position": [-0.05, -0.921, -0.8],
    "confidence": 0.704,
    "size": [0.631, 0.331, 0.178]
  },
  {
    "object_idx": 2,
    "class": "wardrobe",
    "position": [0.972, -0.839, -0.8],
    "confidence": 0.581,
    "size": [2.909, 0.637, 2.38]
  },
  ..... ]

"rooms":
[
  {
    "room_idx": 0,
    "state": 1,
    "distance": 100000.0,
    "viewpoints": [
      {
        "viewpoint_idx": 0,
        "position": [0.0, 0.0, -0.8],
        "has_frontier": false,
        "objects": [6, 9, 0, 1, 2, 4, 5, 7, 8, 10]
      },
      {
        "viewpoint_idx": 1,
        "position": [1.625, -0.225, -0.8],
        "has_frontier": false,
        "objects": [3, 4, 5, 6, 7, 9]
      }
    ]
  },
  ..... ]
```

Figure 6: Visualization of the Json file.

Please note that when translating our representation \mathcal{R} into a JSON file, we begin by listing all the objects in the scene. This is because an object may be associated with multiple viewpoints; directly listing objects under each viewpoint would lead to redundancy and may exceed the prompt’s length limit.

A4 More Experiment Results

We provide more experimental results on the HM3D dataset. First we show the Success Rate on each target object category of the HM3D dataset in Table 4. We can see that our method achieves the best performance on most of the categories.

Table 4: **Success rate of each category on HM3D [13]**. The best and second best results are highlighted in **bold** and underline, respectively.

Method	bed	chair	plant	sofa	toilet	tv_monitor	Average
L3MVN [27]	52.9	51.6	46.4	50.1	41.5	54.2	49.5
TriHelper [43]	57.1	58.6	58.3	58.9	52.3	57.4	57.1
CogNav [6]	67.9	73.4	73.1	67.0	72.6	74.0	<u>72.5</u>
STRIVE (ours)	83.8	86.2	<u>67.6</u>	81.2	81.9	<u>73.3</u>	79.6

We also show more qualitative results of our method on the HM3D dataset in Figure 7. We visualize the trajectory of the agent and the final RGB-D image when the agent stops. The agent is able to efficiently navigate to the target object and stop at a reasonable viewpoint to observe the object.

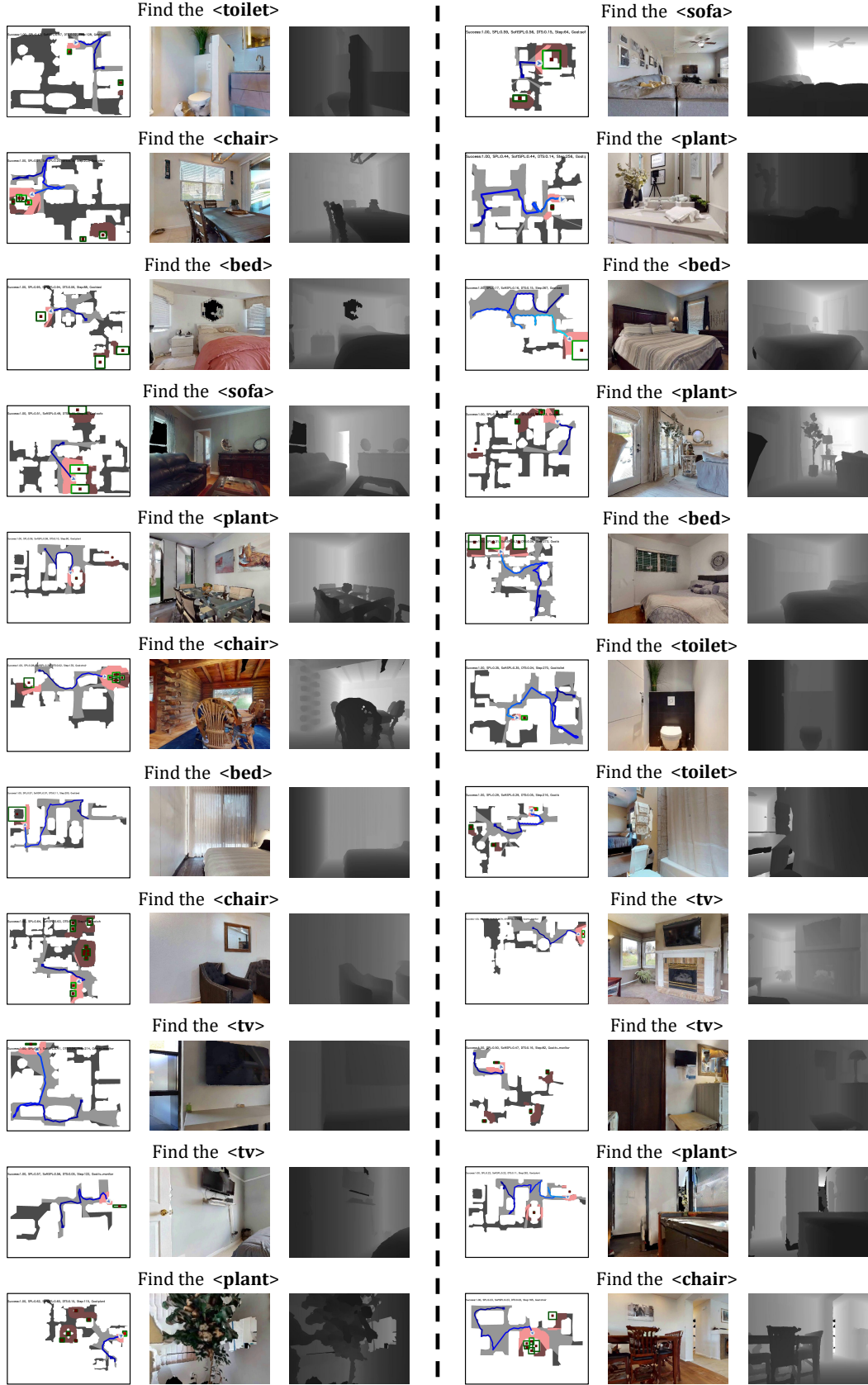


Figure 7: **Qualitative results of STRIVE on HM3D.** We show the trajectory of the agent and the final RGB-D image when the agent stops.

A5 Examples of VLM Reasoning

We provide additional examples of the VLM’s reasoning process in Figs. 8 to 10. These results demonstrate that our structured representation enables the VLM to reason effectively over both spatial layout and semantic cues, leading to more accurate room selection. Moreover, the VLM is able to balance the likelihood of finding the target object with the travel distance cost when planning room-to-room exploration.

Example 1

Chain of Thought

Thought: Analyze the task: Find the 'toilet'.
Output: Target object: toilet.

Thought: Identify the current location and fully explored rooms.
Output: Current location: Room 1 (fully explored).

Thought: Identify the partially explored rooms available for exploration and their distances.
Output: Partially explored rooms: Room 2 (distance 2.288), Room 3 (distance 3.062). Available options: [2, 3].

Thought: Evaluate the likelihood of finding the target ('toilet') in each partially explored room based on visible objects.
Output: Room 2: Visible objects (furniture, door way) do not suggest a bathroom. State is partially explored with frontiers.
Room 3: Visible objects (floor, door way) do not suggest a bathroom. State is partially explored with frontiers.

Thought: Determine the best strategy for exploration given no strong indicators of the target location in the visible areas.
Prioritize exploring closer, partially explored rooms with frontiers.
Output: Both rooms are partially explored with frontiers. Room 2 is closer than Room 3.

Thought: Select the next room based on the evaluation and strategy.
Output: Choose Room 2 as it is the closest partially explored room with frontiers.

Final Answer: Room 2
Reason: The goal is to find a 'toilet'. Neither of the partially explored rooms (Room 2 and Room 3) contain objects in their visible areas that strongly suggest the presence of a toilet or bathroom. Both rooms are partially explored and have frontiers, meaning further exploration is needed to determine if they contain the target. In the absence of specific clues, selecting the closer partially explored room (Room 2, distance 2.288) is the most efficient strategy to continue the search.

Figure 8: Example 1 of VLM reasoning.

Example 2

Chain of Thought

Thought: The task is to find a <potted_plant>.
Output: Analyze the task instruction.

Thought: Currently in Room 5, which is fully explored.
Output: Identify the current location and status.

Thought: Partially explored rooms available are 1, 2, 3, and 4.
Output: Identify the available rooms for exploration and their status.

Thought: Review the objects and viewpoints in each room: Room 1 (distance 4.612) has furniture, stuffed animal, lamp, box. Room 2 (distance 1.553) has floor, doorway, furniture, sink. Room 3 (distance 3.18) has doorway, air duct, electrical outlet, floor, wall. Room 4 (distance 2.9) has floor, furniture, blanket, window, doorway.
Output: Evaluate each partially explored room based on its content and distance.

Thought: Rooms 1, 2, and 3 don't show objects typically associated with potted plants. Room 2 (sink) seems least likely. Room 4 contains a window, which is a common location for potted plants.
Output: Assess the likelihood of finding a <potted_plant> in each room.

Thought: Despite Room 2 being closer, Room 4 is more likely to contain the target object (<potted_plant>) due to the presence of a window. Prioritize likelihood over minimum distance in this case, assuming the distance difference is not prohibitive.
Output: Select the most promising room for exploration.

Thought: Room 4's state is 0, indicating it is partially explored.
Output: Confirm the selected room is partially explored.

Final Answer: Room 4
Reason: The goal is to find a <potted_plant>. I am in Room 5, which is fully explored. I need to select a partially explored room to search. Evaluating the partially explored rooms (1, 2, 3, 4): Room 1 contains furniture, stuffed animal, lamp, and box. Room 2 contains a sink, suggesting a bathroom or utility area, less likely for plants. Room 3 contains doorway, air duct, electrical outlet, floor, and wall. Room 4 contains furniture, blanket, and crucially, windows, which are often associated with potted plants. Although Room 2 is closer, the presence of a window in Room 4 makes it a more promising location to search for a potted plant. Room 4 also has unexplored areas (frontiers). Therefore, Room 4 is the best choice to explore next.

Figure 9: Example 2 of VLM reasoning.

Example 3

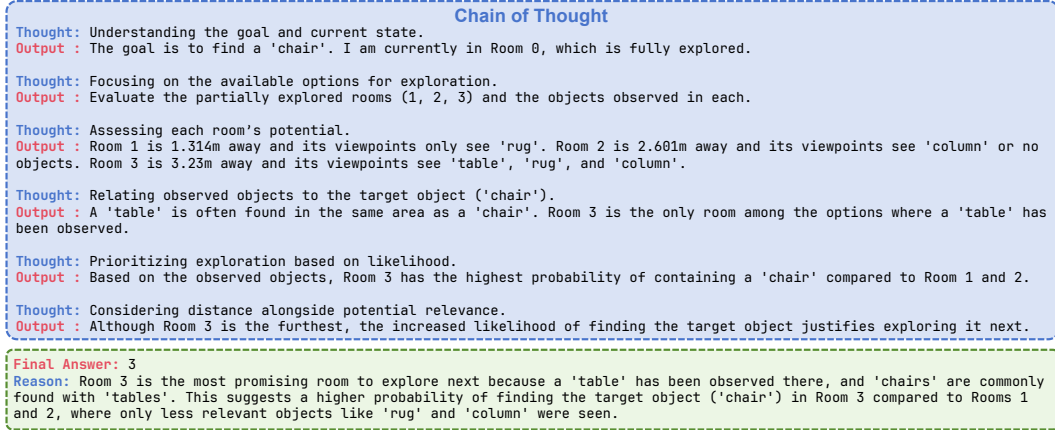


Figure 10: Example 3 of VLM reasoning.

A6 Details of VLM-based Verification

A6.1 Context-Aware Verification

For the context-aware verification, we provide more examples in Fig. 11. These examples show that the VLM can effectively utilize the surrounding context information to verify the detected target objects and avoid false positives.

A6.2 Viewpoint-Optimized Re-Verification

For the viewpoint-optimized re-verification, we first introduce how to compute a better viewpoint for the re-verification. Once we detect a target object, we plan a path from the current position to the target object and sample a series of points along this path with a step size of 0.25m (the agent's forward step length) as candidate viewpoints. For each sampled point, we compute the width (w), height (h), and visible fraction of the target object in the camera frame when the agent faces directly towards the target object. We traverse the sampled points in reverse order (from the target object back to the current position) and select the first point that satisfies all the following conditions as the re-verification viewpoint:

1. More than 95% of the target object's point cloud is visible from this viewpoint.
2. The 2D bounding box area of the target object exceeds that at the original viewpoint.
3. The estimated object width and height at this viewpoint are both greater than 80% of their original values, scaled by the square root of the bounding box area ratio.

The square root of the area ratio is adopted as a scaling factor to account for the expected increase in apparent object dimensions under improved viewing conditions, thereby ensuring a broader and clearer observation for more reliable VLM-based re-verification.

We show more examples of viewpoint-optimized re-verification in Fig. 12. These examples demonstrate that we select a better viewpoint for re-verification, which can help the VLM to verify the target object more accurately. The selected viewpoint is not only closer to the target object but also provides a clearer view of the object, allowing the VLM to make a more confident verification decision.



DINO: **Television** (Confidence 0.52)

Thought: Observing the shape and position of the object within the bounding box and in the context of the room.
Output: The bounding box contains a rectangular object with a frame, hanging on the wall.

Thought: Based on the visual characteristics and typical items found on walls.
Output: The object appears to be a piece of art or a photograph enclosed in a frame.

Thought: Matching the identified object with the options in the pre-defined list.
Output: Comparing the object to the provided list, 'picture' is the most appropriate label.

VLM: **Picture** Context-Aware Verification



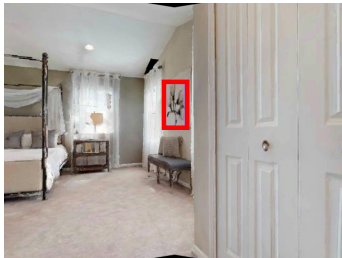
DINO: **Sofa** (Confidence 0.63)

Thought: The bounding box contains a piece of furniture covered with bedding.
Output: The bounding box contains a bed-like object.

Thought: Considering the context of a room with a window and wooden walls, the object appears to be a bed frame with a blanket or sheet.
Output: The object is likely a bed.

Thought: Comparing the provided list of labels, 'bed' is the most fitting description.
Output: The most appropriate label is 'bed'.

VLM: **Bed** Context-Aware Verification



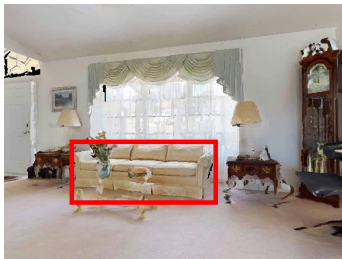
DINO: **Plant** (Confidence 0.51)

Thought: The image shows a framed object mounted on the wall.
Output: The bounding box contains a framed item hanging on the wall.

Thought: Based on the visual appearance of the object within the frame, it resembles a picture or a piece of art.
Output: The framed item appears to be a picture or artwork.

Thought: Comparing the observation with object list, 'picture' is the best fit.
Output: 'picture' is the most suitable label for a framed piece of art hanging on a wall.

VLM: **Picture** Context-Aware Verification



DINO: **Bed** (Confidence 0.66)

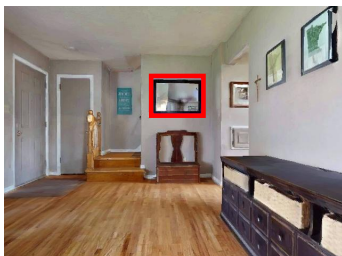
Thought: Analyzing the visual characteristics of the object within the bounding box.
Output: The image inside the bbox shows a piece of furniture designed for seating multiple people, commonly found in a living room. It has a cushioned back and seat, and armrests.

Thought: Identifying the furniture type based on its form.
Output: Based on its appearance and typical function, the object is identifiable as a sofa.

Thought: Comparing the identified object with the allowed labels.
Output: Checking the provided list of predefined objects, 'sofa' is present.

Thought: Selecting the correct label from the list.
Output: The most appropriate label for the object is 'sofa'.

VLM: **Sofa** Context-Aware Verification



DINO: **Television** (Confidence 0.84)

Thought: Observing the full image, the bbox highlights a specific item mounted on the wall.
Output: The bounding box encloses a rectangular object hanging on the wall.

Thought: The primary feature is a reflective area surrounded by a dark frame.
Output: The cropped image shows a reflective surface inside a frame.

Thought: The reflection visible in the object is of the surrounding environment, which is typical for a mirror.
Output: The reflective surface is showing a reflection of the room.

Thought: Combining the observations of its shape, mounting on the wall, frame, and reflection, the object is identifiable as a mirror.
Output: Based on appearance, location, and reflective properties, the object is a mirror.

VLM: **Mirror** Context-Aware Verification

Figure 11: Examples of context-aware verification.



Figure 12: Examples of viewpoint-optimized re-verification.

A7 Dataset Error

Here we provide many annotation errors in the HM3D [13] and MP3D [15] datasets, mainly because these two datasets are collected from real-world scenes and rely on manual annotation of their semantic information. This leads to many scenes where some objects are missed or incorrectly annotated with semantic information. In the object navigation task, this problem is particularly serious. If the target object information is missing in the scene, but the agent actually finds that missing target object, this episode will be incorrectly considered a failure. We provide some examples of annotation errors in the HM3D and MP3D datasets for reference in Fig. 13 and 14.

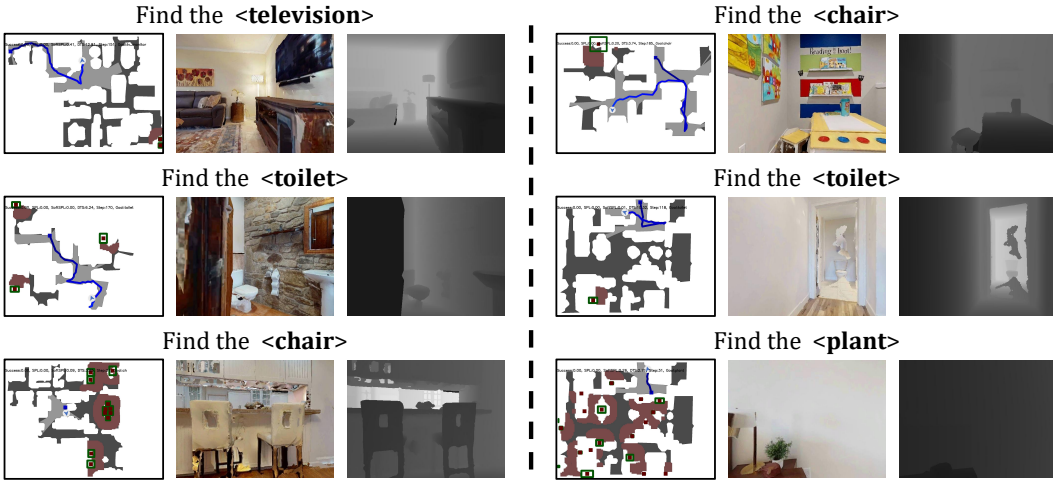


Figure 13: Example of HM3D [13] dataset error causing episode failure. The target object is not annotated in the scene, but the agent finds it. This episode will be incorrectly considered a failure.

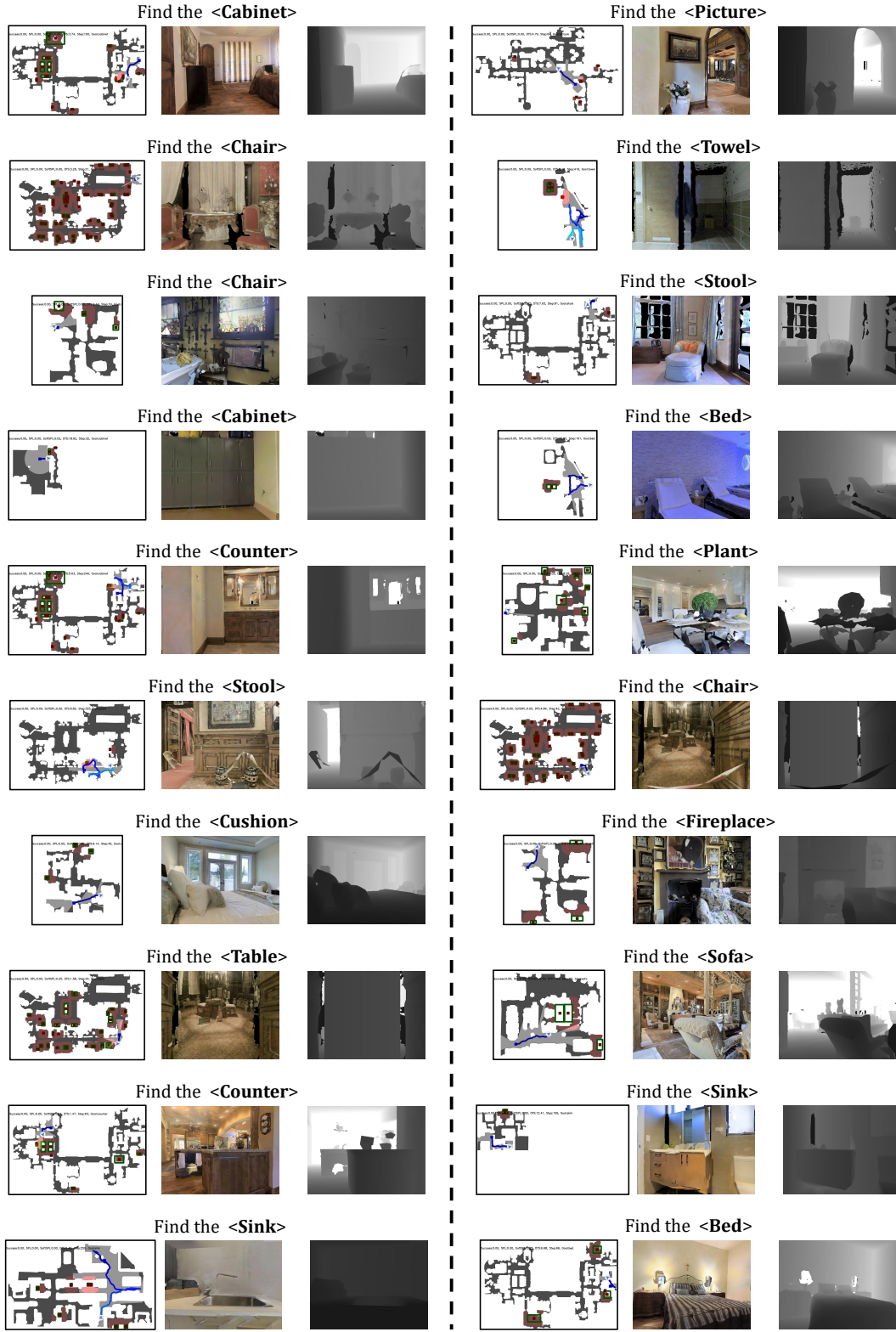


Figure 14: Example of MP3D [15] dataset error causing episode failure. The target object is not annotated in the scene, but the agent finds it. This episode will be incorrectly considered a failure.

A8 Room Segmentation

Based on the scene point cloud, we first construct a top-down-view 2D histogram and extract the wall borders, resulting in a binary mask that highlights the walls in the scene. Next, we generate a whole-scene mask by combining the detected walls and point cloud. After obtaining the whole-scene mask, we gradually dilate the background to obtain the room segmentation. To be specific, we gradually dilate the background. After each dilation step, we check all the connected components in the scene. If a component's area is smaller than a threshold, we mark it as a room and remove that area from the dilated whole-scene mask. We continue dilating the background until no new disconnected regions are found. Finally, we use the marked rooms as seeds and apply the watershed algorithm on the whole-scene mask to obtain the final room segmentation results.



Figure 15: Visualization of room segmentation process.