
Think, Remember, Navigate: Zero-Shot Object-Goal Navigation with VLM-Powered Reasoning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While Vision-Language Models (VLMs) are set to transform robotic navigation,
2 existing methods often underutilize their reasoning capabilities. To unlock the
3 full potential of VLMs in robotics, we shift their role from passive observers to
4 active strategists in the navigation process. Our framework outsources high-level
5 planning to a VLM, which leverages its contextual understanding to guide a frontier-
6 based exploration agent. This intelligent guidance is achieved through a trio of
7 techniques: structured chain-of-thought prompting that elicits logical, step-by-step
8 reasoning; dynamic inclusion of the agent’s recent action history to prevent getting
9 stuck in loops; and a novel capability that enables the VLM to interpret top-down
10 obstacle maps alongside first-person views, thereby enhancing spatial awareness.
11 When tested on challenging benchmarks like HM3D, Gibson, and MP3D, this
12 method produces exceptionally direct and logical trajectories, marking a substantial
13 improvement in navigation efficiency over existing approaches and charting a path
14 toward more capable embodied agents.

15 1 Introduction

16 A primary challenge in robotics is enabling autonomous navigation in unknown environments, a skill
17 crucial for tasks like search and rescue or industrial inspection. Object Goal Navigation (ObjectNav),
18 where an agent must locate a specific object in a new setting, is a particularly demanding version
19 of this problem (3; 18). It requires a combination of advanced spatial awareness and semantic
20 comprehension. Traditional approaches, which often depend on geometric maps and predetermined
21 planning strategies, have difficulty generalizing to novel environments and thus fall short of achieving
22 genuinely intelligent exploration.

23 The emergence of powerful Vision-Language Models (VLMs) has opened a new frontier, giving
24 robots the potential to interpret their surroundings with a human-like understanding of context.
25 Despite this promise, VLMs have typically been integrated into navigation systems in a limited
26 capacity. Many current methods assign the VLM a passive function, such as scene description or
27 query answering, rather than making it the primary strategist. This underuse is a result of superficial
28 knowledge integration, inflexible prompting techniques, and a lack of memory, all of which lead to
29 inefficient navigation and a restricted role for the VLM.

30 To overcome these shortcomings, we propose a novel framework that delegates high-level planning
31 to a VLM, using its inherent contextual knowledge to direct navigation. Our approach reimagines the
32 VLM’s function, elevating it from a simple reactive component to the main navigator. By combining
33 the VLM’s emergent planning skills with frontier-based exploration, dynamic prompting, and multi-
34 view fusion, our framework achieves robust and efficient navigation without relying on conventional
35 planners. Our main contributions include:

- 36 • **Chain-of-Thought (CoT) for Navigation:** We employ CoT reasoning within the VLM,
37 which allows it to produce more logical and context-aware instructions by methodically
38 thinking through each step of the navigation task.
- 39 • **Dynamic Prompts with Action History:** Our system uses advanced prompts that include
40 the agent’s recent actions. This helps to avoid common issues like getting stuck in loops or
41 indecisive movements, leading to more reliable exploration.
- 42 • **Top-Down Map Interpretation:** We enhance the VLM’s reasoning by enabling it to analyze
43 top-down obstacle maps in conjunction with its first-person view, giving it a better sense of
44 the overall space for long-term planning.

45 2 Related Work

46 The challenge of Object Goal Navigation (ObjectNav) has been approached from several angles,
47 beginning with end-to-end learning frameworks that often struggled to generalize beyond their
48 training data (14; 15). Modular pipelines emerged as an alternative, deconstructing the problem into
49 perception, mapping, and planning stages (3). While more robust, these systems could be brittle and
50 prone to cascading errors between components (11). This has motivated a recent surge in zero-shot
51 methodologies that harness the world knowledge of large pre-trained models, thereby avoiding the
52 need for extensive task-specific fine-tuning (18).

53 The integration of Vision-Language Models (VLMs) first involved using them as powerful feature
54 extractors, with methods like CLIP on Wheels (CoWs) (7) leveraging embeddings to link visual
55 scenes with a target object’s name. The role of language models soon became more active, with
56 systems like ESC applying common-sense reasoning to guide exploration (29). More recently, the
57 frontier of ObjectNav has shifted toward offloading high-level strategy entirely to large models.
58 This paradigm includes diverse techniques such as translating maps into text for an LLM to score
59 exploration paths (L3MVN (26)), employing a VLM to directly evaluate the semantic promise
60 of frontiers (VLFM (25)), or even using a VLM to imagine and select optimal future viewpoints
61 (ImagineNav (28)).

62 Our work advances this paradigm by enhancing the VLM’s cognitive role in navigation. We introduce
63 a framework that empowers the VLM with a more profound understanding of the task through a
64 synergistic combination of three key techniques: structured Chain-of-Thought (CoT) prompting to
65 elicit more logical, step-by-step analysis (22); the incorporation of a memory of recent actions to
66 prevent stagnation; and a novel method for providing multimodal spatial context by enabling the
67 VLM to interpret top-down obstacle maps in conjunction with its egocentric view. This holistic
68 approach results in a more effective zero-shot navigator capable of generating more coherent and
69 efficient trajectories.

70 3 Methodology

71 Our method’s navigation is organized into three phases: initialization, exploration, and goal navigation.
72 Our primary contribution is in the exploration phase, where we use a VLM to bring contextual
73 intelligence to a standard frontier-based exploration framework. As shown in Figure 1, the VLM
74 processes the agent’s first-person view, a top-down obstacle map, and a specially designed prompt. It
75 then serves as a high-level guide, suggesting the next move based on its comprehensive understanding
76 of the scene and the goal. Algorithm 1 outlines the comprehensive procedure for the exploration
77 stage, detailing how the VLM guidance, Value Map updates, and frontier prioritization are integrated
78 to navigate the environment.

79 3.1 Waypoint Generation via Frontier Exploration

80 The agent is equipped with a 2D action space, including moving forward, turning, and stopping.
81 We use the VER algorithm (23) for low-level motion control to guide the agent toward specific
82 2D waypoints. These waypoints are identified through a frontier-based exploration method (19),
83 which analyzes depth data to create an obstacle map and find the boundaries between explored and
84 unexplored areas. The midpoints of these frontiers are chosen as potential targets for exploration,
85 promoting a thorough search of the environment.

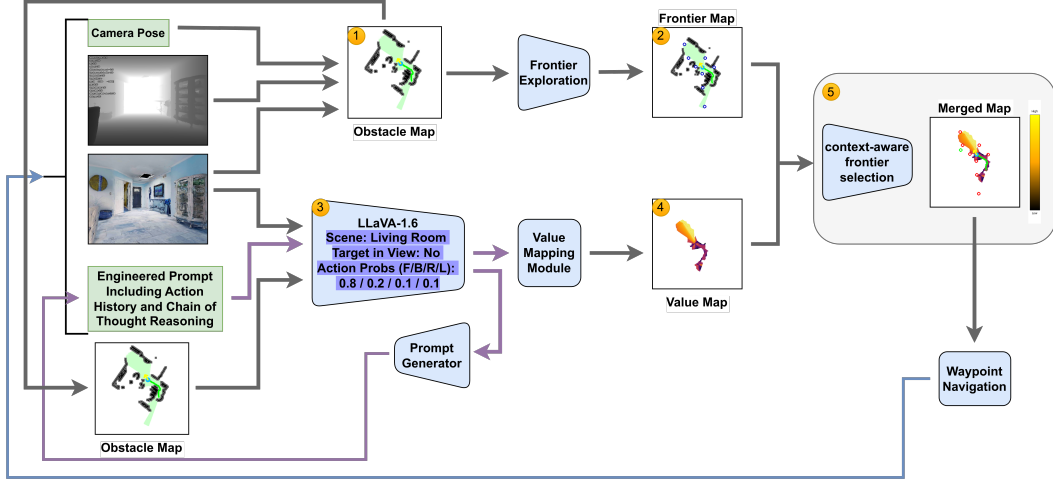


Figure 1: Our System Pipeline. (1) Sensor data is used to create an Obstacle Map. (2) Geometric frontiers are detected on this map. (3) The LLaVA-1.6 VLM analyzes the agent’s egocentric view, the map, and a dynamic prompt that includes action history. (4) The VLM produces semantic scores, which are then used to build a Value Map that indicates the relevance of different areas. (5) The Frontier and Value Maps are combined to prioritize waypoints, directing the agent toward the most promising regions.

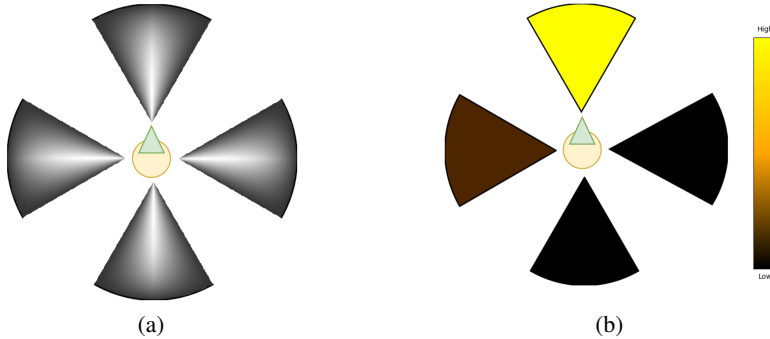


Figure 2: Conceptual illustration of the value map components. (a) Robot’s field of view (FOV) and the associated viewing uncertainty cone. (b) Example action space visualization with VLM-assigned scores: [Forward: 0.9, Backward: 0, Right: 0, Left: 0.1].

86 3.2 Value Map Creation using VLM Guidance

87 Since geometric frontiers do not have semantic meaning, we use the VLM to assign probability scores
 88 to actions like moving forward or turning left, based on the current view, action history, and our
 89 prompt. These scores are used to build a *value map* that reflects how promising different areas are.
 90 The scores are projected onto the local map and are adjusted by a *viewing uncertainty* factor, which
 91 accounts for the fact that the VLM’s judgments may be less accurate for distant or peripheral regions,
 92 as conceptually depicted in Figure 2. The confidence level c for a point at a distance d and angle θ is
 93 calculated as:

$$c(d, \theta) = e^{-\lambda d} \cdot \cos^2 \left(\frac{\theta}{\theta_{\text{fov}}/2} \cdot \frac{\pi}{2} \right) \quad (1)$$

94 Here, θ_{fov} is the camera’s field of view, and λ is a parameter for distance decay. When the agent
 95 observes areas that overlap with previously seen regions, the semantic value $v_{i,j}^{\text{new}}$ and confidence
 96 $c_{i,j}^{\text{new}}$ for each pixel (i, j) are updated through a confidence-weighted average. This ensures that more

97 reliable observations have a stronger influence:

$$v_{i,j}^{\text{new}} = \frac{c_{i,j}^{\text{curr}} v_{i,j}^{\text{curr}} + c_{i,j}^{\text{prev}} v_{i,j}^{\text{prev}}}{c_{i,j}^{\text{curr}} + c_{i,j}^{\text{prev}}} \quad \text{and} \quad c_{i,j}^{\text{new}} = \frac{(c_{i,j}^{\text{curr}})^2 + (c_{i,j}^{\text{prev}})^2}{c_{i,j}^{\text{curr}} + c_{i,j}^{\text{prev}}} \quad (2)$$

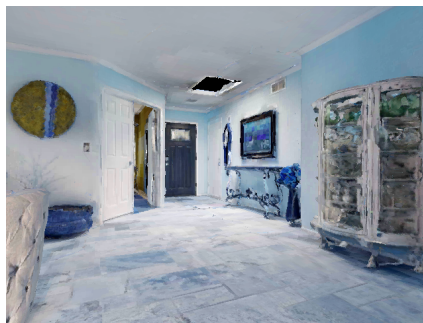
98 The resulting value map is then combined with the geometric frontier map, enabling the system to
99 prioritize exploration of the most semantically relevant areas.

100 3.3 Top-Down View Map Parsing by VLM

101 Drawing inspiration from human spatial reasoning, our framework equips the VLM with both a first-
102 person egocentric view and a top-down obstacle map (Figure 3). Although VLMs are not typically
103 trained on map data, we provide clear instructions in the prompt, such as: "*The second image is a*
104 *top-down obstacle map...*". This merging of first-person and overhead perspectives enhances the
105 VLM’s spatial awareness, leading to better-informed navigation choices, a finding supported by our
106 ablation studies.



(a) Top-Down Obstacle Map



(b) Egocentric View

Figure 3: Dual visual inputs for the VLM: (a) the top-down map shows the spatial layout with obstacles (in gray) and the agent’s heading (arrow), while (b) the egocentric view provides a first-person perspective. This combination improves the VLM’s spatial understanding.

107 Our choice of VLM was guided by a balance of reasoning capability and computational efficiency.
108 After evaluating several open-source models, we selected LLaVA-1.6 (7B) (12) as it provided strong
109 performance on multi-step reasoning tasks while maintaining an inference speed suitable for our
110 navigation loop. During the final goal navigation phase, our system employs a robust hybrid strategy
111 to mitigate potential failures in object detection. If an object is detected with high confidence by
112 standard detectors (e.g., YOLOv7, Grounding-DINO), the system verifies the finding with the VLM
113 and uses a lightweight segmentation model (Mobile-SAM) to delineate the object’s boundaries,
114 ensuring a more reliable final approach to the target.

115 4 Prompt Engineering

116 Successful interaction with the VLM depends heavily on the design of the prompt. We use several
117 techniques to obtain structured and logical responses that are useful for navigation. CoT prompting
118 improves a model’s reasoning by guiding it through a series of intermediate steps to solve a complex
119 problem (22). Instead of just asking for the best move, our prompt (see listing 1) leads the VLM
120 through a logical sequence: (1) *Determine the room type.* (2) *Evaluate if the target object is likely to*
121 *be in that room.* (3) *Confirm if the target is currently visible.* (4) *Suggest the most sensible action*
122 *and assign scores.* This methodical approach, which was refined through extensive testing, ensures
123 that the VLM’s understanding of the scene is consistent with its recommended actions, reducing the
124 chances of logical errors.

125 A common issue in navigation is when an agent gets stuck in a loop or becomes indecisive, especially
126 in visually similar environments like long hallways. To address this, we keep a record of the agent’s

127 last 10 actions. This history is included in the VLM prompt with instructions such as: "*Maintain*
 128 *forward progress... and avoid repetitive actions.*" If the agent still gets stuck, a backup plan is
 129 activated, which temporarily repeats the last successful action to break the cycle. Our ablation study
 130 in Section 5.1 confirms that this memory-based strategy is essential for reliable long-term navigation.

131 5 Experimental Evaluation

132 We tested our framework in the Habitat simulator (9), using the HM3D, Gibson, and MP3D ObjectNav
 133 benchmarks (16; 24; 4). The main metrics for evaluation were Success weighted by Path Length
 134 (SPL) and Success Rate (SR) (1). We compared our approach with other top zero-shot methods.
 135 As Table 1 shows, our approach delivers exceptional efficiency, achieving the highest SPL on the
 136 demanding HM3D and MP3D datasets, and strong performance on Gibson. This indicates that
 137 our method’s deep integration of VLM reasoning leads to more direct and intentional paths. This
 138 high efficiency is a direct result of our CoT and action history techniques, which prevent aimless
 139 exploration and keep the agent on a more logical course.

Table 1: Performance Comparison on ObjectNav Datasets. Our approach utilizes the LLaVA-1.6 (7B) model. SR and SPL are reported in percent. Top values are in bold.

Method	HM3D		MP3D		Gibson	
	SR	SPL	SR	SPL	SR	SPL
CoW (20)	-	6.3	3.7	7.4	-	-
ZSON (21)	25.5	12.6	15.3	4.8	-	-
ESC (29)	39.2	22.3	28.7	14.2	-	-
VLFM (25)	52.5	30.4	36.4	17.5	84.0	52.2
L3MVN (26)	50.4	23.1	-	-	76.1	37.7
GAMap (27)	53.1	26.0	-	-	85.7	55.5
Our Approach	54.3	31.1	36.0	17.7	80.2	53.0

140 5.1 Ablation Studies

141 To determine the impact of each part of our system, we conducted a series of ablation studies on
 142 a smaller set of benchmark episodes (Table 2). The findings show that every component adds
 143 value. Removing the top-down map resulted in a small drop in performance, while taking out
 144 the Chain-of-Thought reasoning caused a more significant decrease. Notably, the largest decline
 145 in performance happened when the **action history module** was removed, with the Success Rate
 146 on HM3D plummeting from 54.3% to 44.0%. This severe drop underscores that providing the
 147 agent with a short-term memory of its recent path is a necessity for robust, long-horizon navigation.
 148 Without this temporal context, the agent is prone to re-evaluating the same frontiers and becoming
 149 trapped in unproductive, oscillatory loops—a common failure mode observed in visually ambiguous
 150 environments.

Table 2: Results of the Ablation Study on a 50-episode subset. Each row indicates performance when the corresponding component is removed. The removal of action history has the most negative effect.

Ablation Condition	HM3D		MP3D		Gibson	
	SR (%)	SPL (%)	SR (%)	SPL (%)	SR (%)	SPL (%)
Full Framework	54.3	31.1	36.0	17.7	80.2	53.0
Without Chain-of-Thought	51.2	29.0	33.9	16.5	75.6	49.4
Without Action History	44.0	23.7	29.2	13.5	65.0	40.4
Without Obstacle Map	53.6	29.6	35.5	16.9	79.2	50.5

151 5.1.1 Effectiveness of Chain-of-Thought Prompting

152 To validate our structured reasoning approach, we tested prompts with varying levels of CoT com-
 153 plexity, based on the hypothesis that guiding the VLM through a logical sequence of questions would

154 yield more coherent and effective navigation decisions. We evaluated four configurations ranging in
 155 complexity: a baseline **No CoT** prompt that only requested action scores without any reasoning; a
 156 **Basic CoT** prompt asking for the single best action to find the target; an **Intermediate CoT** prompt
 157 that added a query for scene identification (e.g., "what room are you in?") before asking for an action;
 158 and our **Full NaviGen CoT**, a complete, multi-step prompt that requires the VLM to first identify the
 159 scene, then assess the likelihood of finding the target object there, and finally recommend an action
 160 based on this analysis.

161 The results, shown in Table 3, confirm a direct positive correlation between the depth of the CoT
 162 prompt and navigation performance. For instance, on HM3D, progressing from a non-CoT prompt to
 163 our full framework improved the SR from 51.2% to 54.3% and boosted the efficiency (SPL) from
 164 29.0 to 31.1. This shows that compelling the VLM to "think step-by-step" is a functional requirement
 165 for translating visual input into successful, goal-oriented action.

Table 3: Ablation Study on Chain-of-Thought Prompt Complexity. Performance evaluated on a 50-episode subset.

Prompt Variation	HM3D		MP3D		Gibson	
	SR	SPL	SR	SPL	SR	SPL
No CoT (scores only)	51.2	29.0	32.5	15.1	75.4	49.5
Basic CoT ("Best action?")	52.0	29.5	33.1	15.9	77.0	50.3
Intermediate CoT (Scene ID)	53.3	30.2	34.8	16.8	78.9	51.7
Full CoT (multi-step, 1)	54.3	31.1	36.0	17.7	80.2	53.0

166 5.2 Qualitative Analysis and Limitations

167 To investigate the reasons for our framework’s efficiency, we conducted a qualitative study of the
 168 VLM’s decision-making. Figure 4 offers a direct comparison between an agent using our full,
 169 multi-step Chain-of-Thought (CoT) prompt and a basic agent with a simpler, non-CoT prompt. This
 170 comparison clearly demonstrates how structured thinking leads to better navigation.

171 The agent without CoT acts impulsively and with little foresight. Its reasoning is basic, connecting
 172 a visual stimulus directly to an action without considering the larger situation (e.g., "the TV is on
 173 the right... turn right"). This results in aimless wandering and an inability to complete the task on
 174 time. In contrast, the agent with our full CoT system follows a more logical and cohesive plan. It
 175 methodically analyzes the problem by first identifying the room type ("Bathroom"), then using its
 176 general knowledge to determine the probability of finding the target there ("Can a TV be Found
 177 Here?: No"), and finally making a sensible decision to search a more likely location ("Recommended
 178 Action: Go forward"). This step-by-step approach avoids pointless detours and creates direct and
 179 efficient paths. This structured thinking prevents logical inconsistencies in planning. Furthermore,
 180 the inclusion of the action history module provides crucial temporal context, allowing the agent to
 181 recognize and break out of unproductive loops, as visually demonstrated in Figure 5. This ability to
 182 avoid stagnation is critical for successful long-horizon navigation.

183 Our analysis also revealed limitations in the standard evaluation protocol, as we observed episodes
 184 marked as failures even when the agent found a valid object of the target category, simply because
 185 it was not the specific instance required by the ground truth. A detailed manual analysis of failed
 186 episodes provided further insights, categorizing them into five primary modes: annotation incom-
 187 pleteness (25%), premature episode termination due to loops (24%), in-view target oversight (19%),
 188 cross-level navigation deficits (18%), and semantic misclassification (14%). These findings highlight
 189 key areas for improvement in both dataset fidelity and the agent’s long-horizon exploration and
 190 perception capabilities.

191 6 Conclusion and Future Work

192 We have presented a new framework for zero-shot object navigation that successfully combines
 193 the semantic reasoning of a VLM with methodical frontier-based exploration. Our main contribu-
 194 tions—dynamic prompt engineering with CoT reasoning and action history—result in more logical



Figure 4: Qualitative analysis of navigation with and without our full CoT framework. The agent’s view and the VLM’s reasoning are shown at various timesteps. The **top row (No CoT)** displays an agent with basic reasoning that wanders aimlessly and fails to locate the target. The **bottom row (Full CoT)** illustrates how a structured, step-by-step reasoning process (such as identifying a bathroom, realizing a TV is not there, and choosing to leave) results in a more intelligent exploration strategy and a direct, successful path. This comparison underscores the crucial role of CoT in achieving more effective and intelligent navigation.

195 and reliable navigation, achieving state-of-the-art efficiency on widely used benchmarks. Despite
 196 these encouraging results, our approach has some limitations that point to areas for future research.
 197 The high computational demand of large VLMs is a barrier to real-time use, which calls for advance-
 198 ments in model compression and more efficient architectures. Additionally, our prompt engineering
 199 is currently a manual process; automating prompt optimization could improve performance. Future
 200 work could also improve the VLM’s spatial reasoning by adding more detailed semantic information
 201 to the top-down map or by training VLMs specifically on map interpretation.

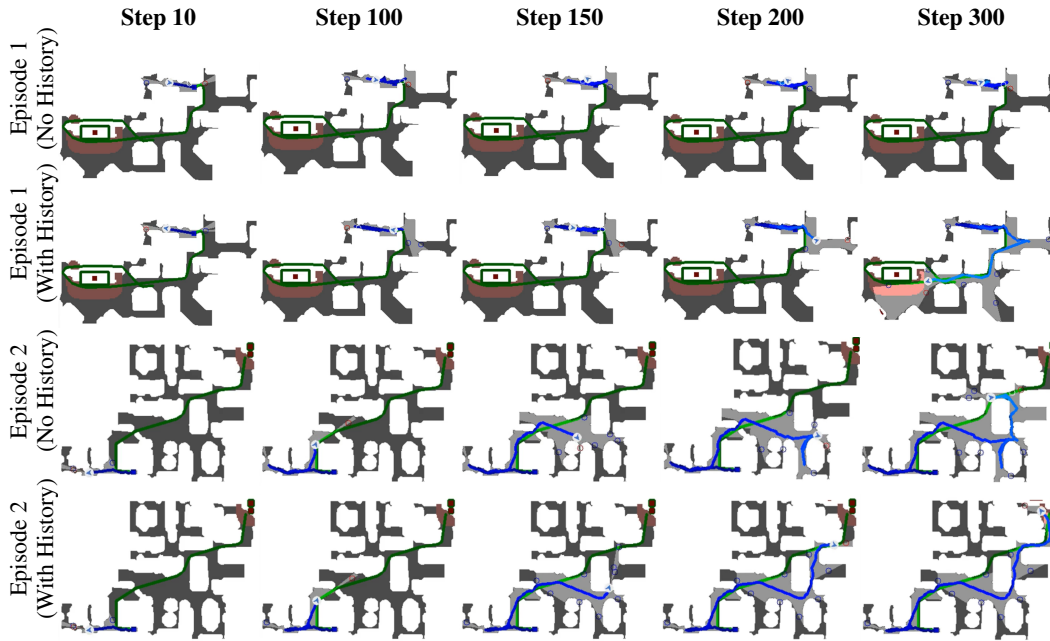


Figure 5: Qualitative comparison across two episodes, demonstrating the critical role of the action history module in resolving common navigation failures. **Episode 1 (Rows 1-2):** The top row shows an agent without action history getting trapped in a corridor. It falls into a persistent decision loop, oscillating between two frontier points and failing to make progress. The second row illustrates how the full framework, using its action history, recognizes the repetitive pattern, breaks the cycle, and successfully finds the target. **Episode 2 (Rows 3-4):** The third row depicts another failure where the agent without temporal memory gets stuck in indecisive movements, ultimately running out of time. In the final row, the action history module prevents this stagnation, enabling the agent to explore efficiently and reach the goal. Each column shows the agent’s view at a specific simulation step.

References

- 202
- 203 [1] P. Anderson, et al. (2018). Vision-and-language navigation: Interpreting visually-grounded navigation
204 instructions in real environments. In *CVPR*.
- 205 [2] W. Cai, et al. (2024). Bridging zero-shot object navigation and foundation models through pixel-guided
206 navigation skill. In *ICRA*.
- 207 [3] D. S. Chaplot, et al. (2020). Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*.
- 208 [4] A. Chang, et al. (2017). Matterport3d: Learning from rgb-d data in indoor environments. In *3DV*.
- 209 [5] V. S. Dorbala, J. F. Mullen, and D. Manocha. (2023). Can an embodied agent find your “cat-shaped mug”?
210 llm-based zero-shot object navigation. *IEEE Robotics and Automation Letters*.
- 211 [6] V. S. Dorbala, et al. (2022). Clip-nav: Using clip for zero-shot vision-and-language navigation. *arXiv
212 preprint arXiv:2211.16649*.
- 213 [7] S. Y. Gadre, et al. (2022). CLIP on wheels: Zero-shot object navigation as object localization and
214 exploration. *arXiv preprint arXiv:2203.10421*.
- 215 [8] N. Gireesh, et al. (2022). Object goal navigation using data regularized q-learning. In *CASE*.
- 216 [9] K. Yadav et al. (2023). Habitat Challenge 2023. <https://aihabitat.org/challenge/2023/>.
- 217 [10] A. Khandelwal, et al. (2022). Simple but Effective: CLIP Embeddings for Embodied AI. In *CVPR*.
- 218 [11] G. Kumar, et al. (2021). Gcexp: Goal-conditioned exploration for object goal navigation. In *RO-MAN*.
- 219 [12] H. Liu, C. Li, Y. Li, and Y. J. Lee. (2024). Improved baselines with visual instruction tuning. *arXiv preprint
220 arXiv:2310.03744*.

- 221 [13] A. Majumdar, et al. (2022). Zson: Zero-shot object-goal navigation using multimodal goal embeddings. In
222 *NeurIPS*.
- 223 [14] A. Mousavian, et al. (2019). Visual representations for semantic target driven navigation. In *ICRA*.
- 224 [15] X. Ye and Y. Yang, "Efficient robotic object search via hiem: Hierarchical policy learning with intrinsic-
225 extrinsic modeling," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4425–4432, 2021.
- 226 [16] S. K. Ramakrishnan et al. (2021). Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments
227 for embodied ai. *arXiv preprint arXiv:2109.08238*.
- 228 [17] B. Sun, et al. (2025). FrontierNet: Learning Visual Cues to Explore. *arXiv preprint arXiv:2501.04597*.
- 229 [18] J. Sun, J. Wu, Z. Ji, and Y.-K. Lai. (2024). A survey of object goal navigation. *IEEE Transactions on*
230 *Automation Science and Engineering*.
- 231 [19] A. Topiwala, P. Inani, and A. Kathpal. (2018). Frontier Based Exploration for Autonomous Robot. *arXiv*
232 *preprint arXiv:1806.03581*.
- 233 [20] S. Y. Gadre, K. Jiang, M. Wortsman, G. Ilharco, G. Gkioxari, D. Fried, L. Schmidt, A. Farhadi, and
234 M. Rastegari, "CLIP on wheels: Zero-shot object navigation as object localization and exploration," in
235 *Conference on Robot Learning (CoRL)*, 2022, pp. 70–80.
- 236 [21] A. Majumdar, Z. Al-Halah, and K. Grauman, "ZSON: A fast and efficient zero-shot open-set object
237 navigator," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2023, pp. 22744–22753.
- 238 [22] J. Wei, et al. (2023). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In
239 *NeurIPS*.
- 240 [23] E. Wijmans, I. Essa, and D. Batra. (2022). Ver: Scaling on-policy rl leads to the emergence of navigation
241 in embodied rearrangement. In *NeurIPS*.
- 242 [24] F. Xia, et al. (2018). Gibson Env: Real-World Perception for Embodied Agents. In *CVPR*.
- 243 [25] N. Yokoyama, et al. (2024). Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In
244 *ICRA*.
- 245 [26] B. Yu, H. Kasaei, and M. Cao. (2023). L3mvm: Leveraging large language models for visual target
246 navigation. In *IROS*.
- 247 [27] S. Yuan, et al. (2024). GAMap: Zero-Shot Object Goal Navigation with Multi-Scale Geometric-Affordance
248 Guidance. *arXiv preprint arXiv:2410.23978*.
- 249 [28] X. Zhao, et al. (2024). ImagineNav: Prompting Vision-Language Models as Embodied Navigator through
250 Scene Imagination. *arXiv preprint arXiv:2410.09874*.
- 251 [29] K. Zhou, et al. (2023). Esc: Exploration with soft commonsense constraints for zero-shot object navigation.
252 In *ICML*.

253 **A Appendix**

254 This appendix provides the pseudocode for the exploration phase of our approach and the full,
 255 optimized VLM prompt that resulted from our iterative refinement process.

Algorithm 1 Object Navigation Exploration Phase

```

1: Initialize map  $\mathcal{M}_0$ , Value map  $\mathcal{V}_0(m) \leftarrow 0.5 \forall m$ , Confidence  $C_0(m) \leftarrow 0 \forall m$ , pose  $p_0$ , action
   history  $\mathcal{H}_a \leftarrow \emptyset$ 
2: function Explore( $c_{\text{target}}, p_{\text{curr}}, \mathcal{M}_{\text{curr}}, \mathcal{V}_{\text{curr}}, C_{\text{curr}}, \mathcal{H}_{\text{curr}}$ )
3:  $p_t \leftarrow p_{\text{curr}}; \mathcal{M}_{t-1} \leftarrow \mathcal{M}_{\text{curr}}; \mathcal{V}_{t-1} \leftarrow \mathcal{V}_{\text{curr}}; C_{t-1} \leftarrow C_{\text{curr}}; \mathcal{H}_a \leftarrow \mathcal{H}_{\text{curr}}$ 
4: for  $t = 0$  to MaxExplorationSteps do
5:   Get observation  $O_t = (I_t, D_t)$  at pose  $p_t$ 
6:    $\mathcal{M}_t \leftarrow \text{UpdateMap}(\mathcal{M}_{t-1}, D_t, p_t)$ 
7:    $\mathcal{F}_g \leftarrow \text{FindGeometricFrontiers}(\mathcal{M}_t, p_t)$ 
8:    $\pi_{\text{gen},t} \leftarrow \text{GeneratePrompt}(\pi_{\text{template}}, \mathcal{H}_a, c_{\text{target}})$ 
9:    $(S_{\text{act},t}, \beta_{\text{obj},t}, R_{\text{vlm},t}) \leftarrow \text{VLMQuery}(I_t, \mathcal{M}_t, \mathcal{H}_a, \pi_{\text{gen},t})$ 
10:   $\mathcal{V}_t \leftarrow \mathcal{V}_{t-1}; C_t \leftarrow C_{t-1}$ 
11:  for each map cell  $m$  in current FOV of  $O_t$  do
12:     $v_m^{\text{vlm}} \leftarrow \text{ProjectVLMscore}(S_{\text{act},t}, m, p_t, I_t)$ 
13:     $c_m^{\text{view}} \leftarrow \text{ViewingConfidence}(m, p_t, \text{FOV}_{\text{params}})$  {Eq. 1}
14:     $\mathcal{V}_t(m) \leftarrow \frac{c_m^{\text{view}} v_m^{\text{vlm}} + C_{t-1}(m) \mathcal{V}_{t-1}(m)}{c_m^{\text{view}} + C_{t-1}(m) + \epsilon}$  {Eq. 2}
15:     $C_t(m) \leftarrow \frac{(c_m^{\text{view}})^2 + (C_{t-1}(m))^2}{c_m^{\text{view}} + C_{t-1}(m) + \epsilon}$  {Eq. 2}
16:  end for
17:   $\mathcal{F}_p \leftarrow \emptyset$ 
18:  for each frontier  $f \in \mathcal{F}_g$  do
19:     $v_f \leftarrow \text{QueryValueMap}(\mathcal{V}_t, C_t, f)$ 
20:    Add  $(f, v_f)$  to  $\mathcal{F}_p$ 
21:  end for
22:  Sort  $\mathcal{F}_p$  by  $v_f$  descending
23:   $w^* \leftarrow \text{SelectExplorationWaypoint}(\mathcal{F}_p)$ 
24:   $(p_{t+1}, a_t) \leftarrow \text{MoveTo}(w^*)$ 
25:   $\mathcal{H}_a \leftarrow \text{UpdateActionHistory}(\mathcal{H}_a, a_t)$ 
26:   $o_{\text{sensor}} \leftarrow \text{ObjectDetect}(I_t, c_{\text{target}})$ 
27:  if  $\beta_{\text{obj},t}$  and  $o_{\text{sensor}}$  then
28:    NavigateToTarget(...) {Transition to goal phase}
29:  break
30:  end if
31: end for
32: end function

```

```

256 You are a robot navigating an indoor environment in search of a [
257   TARGET_OBJECT].
258 Once you find it, move near the [TARGET_OBJECT] and stop.
259 The first image is your current observation and the second image is a
260 top downview obstacle map of the environment. The grey areas are
261 obstacles and The robots direction is visible with an arrow.
262 You must think step by step and ensure that all parts of your response
263 are consistent.
264
265 Here are your recent actions: [ACTION_HISTORY]
266 Maintain forward progress and avoid getting stuck in a loop.
267
268 Here are the tasks:
269 1. Identify what part of the house we are about to enter (choose from:
270   [bedroom, living room, kitchen, corridor, bathroom]).
271 2. Assess whether a [TARGET_OBJECT] can realistically be found in this
272   area, based on common sense and the current observation.
273 3. Is there a [TARGET_OBJECT] in the current scene?
274

```

```

275 4. Determine the most logical next action for the robot (choose from:
276     [go forward, go backward, turn right, turn left]).
277 - The chosen action must prioritize exploring areas likely to contain
278   a [TARGET_OBJECT].
279 - Avoid suggesting actions that contradict previous observations (e.g
280   ., don't explore a bathroom if couches aren't found there).
281 - If you are in a corridor, continue your path and Try to exit the
282   corridor and describe where it leads.
283 5. Provide a probability score for each possible action. Each
284   probability score should be a number between 0 and 1.
285
286 When providing your response, use this structure:
287 1. Part of the House: [Your answer]
288 - Reasoning: [Explain why you think this is the correct part of the
289   house based on the observation and map.]
290 2. Can a [TARGET_OBJECT] Be Found Here?: [Yes/No]
291 - Reasoning: [Explain why or why not.]
292 3. Have You Found the [TARGET_OBJECT]?: [Yes/No]
293 4. Recommended Action: [Your action]
294 - Reasoning: [Explain why this action is the most logical based on
295   steps 1, 2, and the action history.]
296 5. Probability Scores for Each Action:
297 - Go forward: [Score]
298 - Go backward: [Score]
299 - Turn right: [Score]
300 - Turn left: [Score]

```

Listing 1: The full, optimized CoT prompt for our approach.

302 **NeurIPS Paper Checklist**

303 **1. Claims**

304 Question: Do the main claims made in the abstract and introduction accurately reflect the
305 paper’s contributions and scope?

306 Answer: [Yes]

307 Justification: The abstract and introduction claim that our framework improves navigation
308 efficiency (SPL) by deeply integrating a VLM using CoT and action history. This is directly
309 supported by the main results in Table 1 and the ablation studies in Table 2.

310 **2. Limitations**

311 Question: Does the paper discuss the limitations of the work performed by the authors?

312 Answer: [Yes]

313 Justification: The conclusion section explicitly discusses limitations, including the computa-
314 tional cost of VLMs and the reliance on manual prompt engineering.

315 **3. Theory assumptions and proofs**

316 Question: For each theoretical result, does the paper provide the full set of assumptions and
317 a complete (and correct) proof?

318 Answer: [NA]

319 Justification: This paper is empirical and does not present theoretical results or proofs.

320 **4. Experimental result reproducibility**

321 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
322 perimental results of the paper to the extent that it affects the main claims and/or conclusions
323 of the paper (regardless of whether the code and data are provided or not)?

324 Answer: [Yes]

325 Justification: The paper specifies the VLM used (LLaVA-1.6 7B), the datasets (HM3D,
326 MP3D, Gibson), the simulation environment (Habitat), and the core components of the
327 method, including the full prompt in the appendix. This provides a clear path for reproduc-
328 tion.

329 **5. Open access to data and code**

330 Question: Does the paper provide open access to the data and code, with sufficient instruc-
331 tions to faithfully reproduce the main experimental results, as described in supplemental
332 material?

333 Answer: [No]

334 Justification: The paper does not include a link to the source code at this time. However, it
335 relies entirely on publicly available datasets and models, and the methodology is described
336 in sufficient detail to allow for independent reimplementation.

337 **6. Experimental setting/details**

338 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
339 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
340 results?

341 Answer: [Yes]

342 Justification: The paper specifies that it uses the standard validation splits for the HM3D,
343 MP3D, and Gibson datasets as defined by the Habitat Challenge. As a zero-shot method,
344 no training is performed. Key algorithmic details like the decay parameter λ and the action
345 history length are provided.

346 **7. Experiment statistical significance**

347 Question: Does the paper report error bars suitably and correctly defined or other appropriate
348 information about the statistical significance of the experiments?

349 Answer: [No]

350 Justification: Error bars are not reported. The evaluation is performed over large, standard
351 benchmark datasets containing hundreds to thousands of episodes, and the aggregate metrics
352 (SR, SPL) are standard for this domain, providing a robust measure of performance.

353 8. Experiments compute resources

354 Question: For each experiment, does the paper provide sufficient information on the com-
355 puter resources (type of compute workers, memory, time of execution) needed to reproduce
356 the experiments?

357 Answer: [Yes]

358 Justification: The full paper specifies the VLM model (LLaVA-1.6 7B), and the hardware
359 used (NVIDIA RTX A6000), which provides sufficient information to understand the
360 computational requirements.

361 9. Code of ethics

362 Question: Does the research conducted in the paper conform, in every respect, with the
363 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

364 Answer: [Yes]

365 Justification: The research uses publicly available datasets and models for the task of robotic
366 navigation and does not involve human subjects or sensitive data.

367 10. Broader impacts

368 Question: Does the paper discuss both potential positive societal impacts and negative
369 societal impacts of the work performed?

370 Answer: [Yes]

371 Justification: A broader impact statement is included in the conclusion, discussing potential
372 positive applications in assistive robotics and acknowledging potential dual-use concerns.

373 11. Safeguards

374 Question: Does the paper describe safeguards that have been put in place for responsible
375 release of data or models that have a high risk for misuse (e.g., pretrained language models,
376 image generators, or scraped datasets)?

377 Answer: [NA]

378 Justification: The paper does not release new models or datasets. It uses an existing, publicly
379 available VLM. Therefore, safeguards for release are not applicable.

380 12. Licenses for existing assets

381 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
382 the paper, properly credited and are the license and terms of use explicitly mentioned and
383 properly respected?

384 Answer: [Yes]

385 Justification: The paper properly cites the creators of the datasets (HM3D, MP3D, Gibson),
386 the simulation environment (Habitat), and the VLM (LLaVA) used. These are standard
387 academic assets with permissive licenses.

388 13. New assets

389 Question: Are new assets introduced in the paper well documented and is the documentation
390 provided alongside the assets?

391 Answer: [NA]

392 Justification: The paper does not introduce any new assets like datasets or models.

393 14. Crowdsourcing and research with human subjects

394 Question: For crowdsourcing experiments and research with human subjects, does the paper
395 include the full text of instructions given to participants and screenshots, if applicable, as
396 well as details about compensation (if any)?

397 Answer: [NA]

398 Justification: This research does not involve crowdsourcing or human subjects.

399 **15. Institutional review board (IRB) approvals or equivalent for research with human**
400 **subjects**

401 Question: Does the paper describe potential risks incurred by study participants, whether
402 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
403 approvals (or an equivalent approval/review based on the requirements of your country or
404 institution) were obtained?

405 Answer: [NA]

406 Justification: This research does not involve human subjects.

407 **16. Declaration of LLM usage**

408 Question: Does the paper describe the usage of LLMs if it is an important, original, or
409 non-standard component of the core methods in this research?

410 Answer: [Yes]

411 Justification: The Vision-Language Model (a type of LLM) is the central and core method-
412 ological component of this research. Its usage is described in detail in Section 3 and
413 4.