Mem2Ego: Empowering Vision-Language Models with Global-to-Ego Memory for Long-Horizon Embodied Navigation

Anonymous CVPR submission

Paper ID 13

Abstract

Recent advancements in Large Language Models (LLMs) 001 002 and Vision-Language Models (VLMs) have made them powerful tools in embodied navigation, enabling agents to 003 leverage commonsense and spatial reasoning for efficient 004 exploration in unfamiliar environments. Existing LLM-005 based approaches convert global memory, such as semantic 006 or topological maps, into language descriptions to guide 007 navigation. While this improves efficiency and reduces re-008 dundant exploration, the loss of geometric information in 009 010 language-based representations hinders spatial reasoning, especially in intricate environments. To address this, VLM-011 based approaches directly process ego-centric visual inputs 012 to select optimal directions for exploration. However, re-013 lying solely on a first-person perspective makes navigation 014 a partially observed decision-making problem, leading to 015 016 suboptimal decisions in complex environments. In this paper, we present a novel vision-language model (VLM)-based 017 018 navigation framework that addresses these challenges by adaptively retrieving task-relevant cues from a global mem-019 020 ory module and integrating them with the agent's egocentric observations. By dynamically aligning global con-021 022 textual information with local perception, our approach 023 enhances spatial reasoning and decision-making in longhorizon tasks. Experimental results demonstrate that the 024 proposed method surpasses previous state-of-the-art ap-025 proaches in object navigation tasks, providing a more ef-026 027 fective and scalable solution for embodied navigation.

028 1. Introduction

Embodied navigation is a crucial component of embodied artificial intelligence, with widespread applications in
diverse scenarios such as domestic environments, office
settings, logistics and delivery, and factory inspections
[4, 22, 34]. Its significance stems from its ability to enable
agents to autonomously navigate and perform tasks within
physical environments [13, 17].

Embodied navigation poses two key challenges. First, 036 unlike autonomous driving, which typically occurs in struc-037 tured outdoor environments, embodied navigation requires 038 operating in diverse indoor and industrial settings such as 039 factories, shopping malls, and offices. These spaces feature 040 intricate layouts and obstacles, demanding advanced per-041 ception and planning [4, 16, 20, 26]. Second, it necessitates 042 a high degree of autonomy, as agents must adapt to unfamil-043 iar environments without relying on predefined maps. They 044 must interpret human instructions and dynamically interact 045 with their surroundings to navigate effectively. This work 046 focuses on Object Goal Navigation (ObjectNav), a task in 047 which agents must locate specified objects within complex 048 spaces [2, 15]. 049

In recent years, the rapid advancement of large language 050 models (LLMs) has opened new possibilities for embodied 051 navigation [29, 32]. These models enable robots to leverage 052 commonsense reasoning, improving their understanding of 053 natural language commands and enhancing the integration 054 of perceptual data. This allows for navigation decisions 055 that better align with human intentions [10, 28]. Further-056 more, recent ObjectNav research underscores the impor-057 tance of historical information in improving environmen-058 tal understanding, decision-making, and grounding naviga-059 tion instructions [5, 27]. This has led to the incorporation 060 of memory systems into LLMs, such as episodic memory 061 for past experiences and scene graph memory for structur-062 ing environmental data. However, because these memory 063 systems often represent memories using natural language, 064 which lacks geometric information, the spatial reasoning 065 capacity of LLMs is compromised. 066

Alongside these advancements, there is growing interest 067 in using images as a primary source of guidance by integrat-068 ing foundation models with low-level planners [3, 18]. This 069 approach takes advantage of the advanced visual and lan-070 guage understanding of foundation models, offering an ef-071 fective alternative to traditional map-based methods, which 072 often rely on costly and disruption-prone depth sensing and 073 localization. However, these methods predominantly rely 074 on a first-person perspective without incorporating global 075

133

134

135

136

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

memory into the decision-making process. As a result, theytend to lead to redundant exploration and reduced efficiencyin complex environments.

079 In this paper, we propose a novel Vision-Language 080 Model (VLM)-based navigation framework that addresses these challenges by adaptively retrieving task-relevant cues 081 from a global memory module, which are then mapped to 082 the agent's ego-view visual observations. By integrating 083 084 global contextual information with local perceptual inputs, our framework enables more informed action decisions, 085 086 enhancing the agent's situational awareness and decisionmaking capabilities. The approach significantly enhances 087 the agent's ability to navigate complex, long-horizon tasks 088 by dynamically aligning global context with egocentric rea-089 soning, offering a more effective and scalable solution for 090 091 embodied navigation. The experimental results demonstrate that our proposed navigation pipeline outperforms state-of-092 the-art baselines. Through an ablation study, we verified 093 094 the essential nature of each component of our method. Using our proposed data collection approach, the supervised 095 096 fine-tuned Llama3.2-11B model exhibited superior performance compared to both the vanilla Llama3.2-11B model 097 and GPT-4o. 098

099 2. Related Work

Existing studies that leverage VLMs and LLMs for naviga-tion can be categorized into the following directions.

102 2.1. LLM-based Navigation

These approaches often construct a global memory map
based on image observations and use natural language to
describe candidate points for navigation, with action decisions driven by large language models (LLMs).

Several methods fall within this category, including LFG 107 108 [23], VoroNav [28], ESC [35], and openFMNav [9]. LFG 109 uses frontier-based exploration and large language models 110 to score potential subgoals and guide navigation based on 111 the robot's observations and exploration progress. VoroNav introduces Reduced Voronoi Graphs (RVGs) to optimize 112 the robot's exploration by identifying intersections that pro-113 vide the best observational opportunities, while the LLM 114 predicts the next best waypoint. ESC uses commonsense 115 116 knowledge and frontier-based exploration to navigate toward objects in the environment, while openFMNav ad-117 118 dresses challenges related to human instructions that imply target objects and zero-shot generalization. These methods 119 120 employ LLMs to dynamically update a semantic map as the 121 robot explores, enhancing memory and reducing redundant 122 exploration.

While these methods offer the advantage of maintaining a global map and using high-level reasoning, they also
face limitations. The language-based reasoning used for
decision-making sacrifices high-dimensional semantic in-

formation, such as spatial and geometric details, which can
constrain performance in complex environments. Further-
more, translating raw ego-view observations into abstract
linguistic descriptions may weaken the model's capacity for
precise spatial reasoning.127
128
129

2.2. Value Map-based Navigation

In this class of methods, a global value function is computed based on ego-view observations, and actions are chosen based on the generated value map instead of using VLMs for decision-making.

Notable approaches in this category include VLFM [31] 137 and InstructNav [14]. VLFM uses a pre-trained vision-138 language model to generate a language-grounded value 139 map, guiding the agent to explore optimal frontiers. In-140 structNav extends the idea of goal-directed navigation by 141 introducing a Dynamic Chain of Navigation that breaks 142 down tasks into sequences of actions and landmarks. These 143 methods partially address memory forgetting by integrating 144 global value maps, but they still face challenges. The value 145 map is still constructed based on local observations, and 146 decision-making driven by vision-language models (VLMs) 147 often lacks a comprehensive global perspective. As a re-148 sult, these approaches frequently lead to suboptimal solu-149 tions constrained by local decision-making. 150

2.3. VLM-based Navigation

These approaches directly leverage first-person perspective images as the input of vision-language models (VLMs) to generate action decisions. By using the spatial reasoning capabilities of VLMs, these methods enable the model to interpret complex environmental features from the robot's current viewpoint, facilitating more informed and contextaware navigation decisions.

CoNVOI [21] and **PIVOT** [18] exemplify approaches that process first-person images with VLMs to facilitate real-time navigation and decision-making. While effective in leveraging immediate visual inputs, these methods lack mechanisms for incorporating historical observations, often resulting in redundant exploration. This limitation poses challenges in long-horizon tasks, where maintaining contextual awareness of past actions is critical for efficient navigation. **VLMNav** [6] addresses some of these limitations by integrating both RGB-D images and the robot's pose information to construct a navigability mask that identifies reachable regions. The model incrementally builds a voxelbased map and refines its action proposals by prioritizing unexplored areas.

NoMaD [24] unifies goal-directed navigation and explo-173ration by using the robot's current image and the goal's im-174age as input. The model includes a transformer backbone175for processing visual data and a diffusion model for pre-176dicting action sequences. A binary mask is applied to the177

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

265

266

input to focus on either exploration (excluding the goal) or
goal-reaching (including the goal). Despite its innovative
design, NoMaD remains constrained by the absence of a
global memory, relying solely on the most recent three observations. This limitation restricts its capacity for sustained
long-term exploration.

Recent methods have sought to integrate VLMs more 184 effectively for embodied navigation. OpenIN [25] fo-185 cuses on navigation tasks where the robot must locate spe-186 187 cific objects that have been moved, introducing a Carrier-Relationship Scene Graph (CRSG) to track objects and their 188 locations. The system uses VLMs to process multimodal in-189 structions and commonsense knowledge to guide navigation 190 decisions. 191

192 **Uni-NaVid** [33] takes a significant step toward unifying different navigation tasks in a single model. It processes 193 both video streams and natural language instructions as in-194 put, creating a framework that can generalize across a range 195 196 of navigation tasks. By training on diverse data, including video question answering and captioning tasks, Uni-NaVid 197 improves its performance in real-world scenarios and en-198 ables asynchronous execution for efficiency. 199

These methods move toward integrating both global and
local information more effectively, enabling the robot to
navigate complex environments with a better understanding
of spatial context. However, challenges remain in optimizing the trade-off between VLMs' generalization capabilities
and the need for precise, real-time navigation.

3. Method

3.1. Problem Formulation

In this work, we focus on the object navigation (ObjNav) 208 209 task, where an agent begins at a random location within an unseen environment and is tasked with finding and navigat-210 ing to a target object, denoted by q. The agent has no access 211 to a pre-built map and must rely entirely on its sensory in-212 213 puts for navigation. At each time step t, the agent captures 214 an egocentric RGB-D image, denoted by o^t , from its on-215 board RGB-D camera. Additionally, the agent has access to its current location and orientation, which are represented 216 by the extrinsic matrix M_{ext} of the camera. Using these 217 inputs, the agent must compute and execute a low-level ac-218 tion, a^t , that efficiently guides it toward the target object. 219

The workflow of our proposed method is illustrated in Figure 1. The VLM-based navigation relies on the integration of a memory module that encompasses three distinct types of memories. The construction and maintenance of this memory module, as well as the VLM-based navigation process, will be discussed in detail in the subsequent sections.

3.2. Memory Construction

The memory module is composed of three distinct types of 228 memories, each serving a different purpose: 229

- Frontier Map: Denoted as M_f , frontier map has been 230 proven to be effective for environment exploration in 231 object navigation tasks, as demonstrated in Shah et al. 232 [23], Zhou et al. [35]. We adopt an approach similar to 233 that used in ESC [35] to construct the frontier map. Using 234 the agent's position and camera parameters, RGB-D im-235 ages are transformed into 3D space, where each 2D pixel 236 is mapped to a 3D voxel in the global coordinate system. 237 Voxels located near the floor, with no obstacles along the 238 height dimension, are classified as free space. A fron-239 tier in this map is defined as the boundary between free 240 and unexplored areas. This frontier map is maintained 241 throughout the navigation task. 242
- Landmark Semantic Memory: Denoted as M_l , this memory stores descriptions of the landmarks that the agent has seen in the past. Each entry includes the global coordinates of the landmark and a description of the nearby semantic information, such as objects or decoration texture. For example: "[13.2, 5.4]: Located on the floor near a sink. There is a bath tub nearby.". The description of each landmark is generated by the VLMs, as explained in Section 3.3.
- Visitation Memory: Denoted as M_v , this memory keeps track of the landmarks that the agent has already visited. By maintaining a record of visited locations, the visitation memory serves as a crucial mechanism to prevent redundant exploration and improve overall exploration efficiency.

3.3. Mem2Ego Navigation

At each time step t, given the image-based observation o^t and the three types of memories— M_f^t , M_l^t , and M_v^t —introduced in section 3.2, the proposed memory-toegocentric (Mem2Ego) navigation process can be formulated as follows:

$$a^{t} = f_{\theta}(o^{t}, M_{f}^{t}, M_{l}^{t}, M_{v}^{t}, g)$$
 (1) 264

Further details are provided in the following sections.

3.3.1. Panoramic Observation Generation

After the environment is initialized or the agent reaches a 267 new location, the agent captures four egocentric RGB-D 268 images by rotating its viewpoint 90 degrees at each step. 269 These images are then stacked to construct a 360-degree 270 panoramic observation o_{pano}^t (see Equation 2), offering a 271 comprehensive representation of the surrounding environ-272 ment. Compared to navigation methods relying on a sin-273 gle egocentric view, this panoramic approach enhances the 274

309

310

326



Figure 1. Workflow of our proposed method. Our method maintains three types of memories and project cues from them onto the egocentric images for goal location prediction. Further details are provided in Section 3

agent's spatial awareness and scene understanding. A simi-275 276 lar strategy has been employed in Long et al. [14].

277
$$o_{\text{pano}}^t = \text{Concatenate}([o_0^t, o_{\pi/2}^t, o_{\pi}^t, o_{3\pi/2}^t])$$
 (2)

3.3.2. Frontier and Visitation Memory Projection 278

279 Based on the agent's position and the newly captured depth 280 images, the navigation map and corresponding frontiers are updated following the method outlined in Section 3.2. Can-281 didate locations, denoted as $[C_1, ..., C_N]$ in Equation 3, are 282 generated by combining frontier clustering and grid-based 283 sampling. The centroid of each frontier segment is com-284 puted by clustering all points within the segment. However, 285 using the centroid directly as a candidate may result in un-286 287 reachable goal positions. To mitigate this, we identify the 288 nearest grid point on the floor area to the centroid, ensuring that the candidate is accessible to the agent. Similarly, 289 visited locations, $[\mathbf{V}_1, ..., \mathbf{V}_M]$, are extracted from the visi-290 291 tation memory M_v^t , as shown in Equation 4.

2
$$[\mathbf{C}_1, ..., \mathbf{C}_N] = \text{CandidatesGeneration}(M_f^t)$$

$$[\mathbf{V}_1, ..., \mathbf{V}_M] = \text{VisitationExtraction}(M_v^t)$$
 (4)

295 Once determined, the global coordinates of these candidates and visitations are projected onto the egocentric im-296 age plane as pixel locations $[\mathbf{c}_1, ..., \mathbf{c}_N]$ and $[\mathbf{v}_1, ..., \mathbf{v}_M]$, as 297 shown in Equation 5, where K and M_{ext} represent the cam-298 era intrinsics and extrinsics, respectively. 299

$$\mathbf{c}_{i} = \operatorname{Projection}(\mathbf{C}_{i}), \quad \mathbf{v}_{i} = \operatorname{Projection}(\mathbf{V}_{i})$$
where $\mathbf{c}_{i} = (x_{i}, y_{i}), \quad \mathbf{C}_{i} = (X_{i}, Y_{i}, Z_{i}), \text{ similar for } \mathbf{v}_{i} \text{ and } \mathbf{V}_{i}$

$$[x'_{i}, y'_{i}, w_{i}]^{T} = K \cdot M_{\text{ext}} \cdot [X_{i}, Y_{i}, Z_{i}, 1]^{T}$$

$$(x_{i}, y_{i}) = (\frac{x'_{i}}{w_{i}}, \frac{y'_{i}}{w_{i}})$$
(5)

An annotation function is then applied to map these locations onto the panoramic observation o_{pano}^t , resulting in an 302 annotated observation o_{anno}^t , as outlined in Equation 6. In 303 the annotated image, candidate locations are depicted as 304 green circles, each labeled with a unique identifier corre-305 sponding to its position in the image. Similarly, visited lo-306 cations are marked as blue circles, but only if they are visi-307 ble within the current view. 308

$$o_{\text{anno}}^{t} = \text{AnnotateImage}(o_{\text{pano}}^{t}, [\mathbf{c}_{1}, ..., \mathbf{c}_{N}], [\mathbf{v}_{1}, ..., \mathbf{v}_{M}])$$
(6)

3.3.3. Landmark Memory Retrieval

The panoramic image, augmented with frontier candidates, 311 highlights potential navigation targets within the agent's im-312 mediate field of view. However, it is common that no suit-313 able targets are visible, and more promising options may 314 exist among the landmarks the agent has previously en-315 countered but not yet explored. These previously encoun-316 tered landmarks are stored in the dynamic landmark seman-317 tic memory M_{I}^{t} . To manage the rapid expansion of this 318 memory during navigation, we utilize large language mod-319 els (LLMs) to retrieve the top-k landmarks most relevant to 320 the target object. This retrieval process generates an addi-321 tional observation from memory, o_{mem}^t , which is then incor-322 porated into the decision-making process. The prompt used 323 for memory retrieval is detailed in Appendix 5.1. 324

$$o_{\text{mem}}^t = \text{MemoryRetrieval}_{\text{LLMs}}(M_l^t, k)$$
 (7) 325

3.3.4. Memory Augmented Decision Making

At this stage, the panoramic image with annotations, o_{anno}^t , 327 along with the top-k landmarks retrieved from memory, 328 o_{mem}^t , is used to query VLMs to select the next target lo-329 cation to visit (described in Equation 8). The VLMs are 330 tasked with identifying the marker on the image most likely 331 to lead to the target object, while avoiding markers that are 332

300

(3)

too close to previously visited locations. To ensure consis-333 334 tency in the output format, the top-k landmarks are num-335 bered, and their descriptions are considered only if no suit-336 able marker is identified directly from the panoramic im-337 age. A Chain-of-Thought (CoT) prompting strategy is employed to guide the VLMs in generating a structured rea-338 soning process before producing a single numerical output 339 corresponding to the selected marker. The full prompt used 340 341 for decision-making is provided in Appendix 5.1.

$$a^{t} = f_{\text{VLMs}}\left(\text{prompt}(g), o_{\text{anno}}^{t}, o_{\text{mem}}^{t}\right)$$
(8)

343 3.3.5. Action Execution

344 The marker selected in step 3.3.4 is transformed to the 345 global coordinate system to determine the global coordinates of the target location. Shortest path follower provided 346 347 by habitat simulator is then executed to navigate agent to the target location while avoiding obstacles. Object detec-348 tion is performed each time the agent moves or adjusts its 349 350 viewing angle. The task is deemed successful if the target object is detected within the agent's field of view and 351 the agent successfully navigates to the target object's view-352 points provided by the Habitat dataset. If the target object 353 354 is not detected, the process continues until the agent either reaches the designated viewpoints or exceeds the maximum 355 356 allowed number of exploration steps.

357 3.3.6. Memory Update

358 While only one landmark from the current view is selected as the next-step navigation target, other landmarks may still 359 be valuable for future exploration. The landmark seman-360 361 tic memory is updated before target position navigation described in Section 3.3.5. VLMs are prompted to describe 362 the surrounding environment near each marker annotated on 363 the panoramic image. The output from the VLMs includes 364 a list of marker IDs paired with corresponding descriptions. 365 The marker IDs are then converted to global coordinates 366 367 and, together with their descriptions, saved to the landmark semantic memory for use in future exploration processes. 368 369 The prompt used for landmark description is provided in Appendix 5.1. 370

Meanwhile, the navigation map is updated along the navigation process, using the RGB-D images captured along the way. Additionally, the agent's most recent location is added to the visitation memory to facilitate future exploration.

376 3.4. Data Collection and Model Finetuning

To enhance the capabilities of open-sourced VLMs and
narrow their performance gap with GPT-4o, we design a
pipeline to collect training data for supervised finetuning
(SFT). The data collection pipeline is illustrated in Figure 2.
To improve data diversity and validate the generalization
ability of the model, we gather 40 new categories of objects

from the HSSD dataset, rather than using the original 6 cat-383 egories provided. First, new target objects are sampled from 384 the HSSD scenes. For each frame of data, ground-truth tra-385 jectories from the current position to these targets are calcu-386 lated based on the A^* algorithm and subsequently smoothed 387 using Bézier curves. Egocentric images and the correspond-388 ing ground-truth target pixel (x, y) (defined as the endpoint 389 of the ground-truth trajectory shown in the image) for each 390 image are saved. To construct the multiple marker anno-391 tated image that VLMs encounter in the marker selection 392 task, we generate a few candidate landmarks for each image 393 by sampling from the edge of the floor area. Both ground-394 truth and sampled candidate landmarks are annotated on the 395 egocentric image in the same way as in the Section 3. 396

We collect two types of data for VLM fine-tuning: 397 marker description and target marker selection with ra-398 tionale. To generate marker description data, we use GPT-399 40 to describe the surrounding environment of each marker 400 on the image. For example, "Marker Number: 1 Descrip-401 tion: Positioned near a dining chair...; Marker Number: 2 402 ...". Each target marker selection data entry includes both 403 a rationale and the ID of the selected marker. To ensure 404 a robust rationale, we utilizes egocentric images annotated 405 with the ground-truth trajectory and employ a dual-phase 406 prompting strategy: first, GPT-40 is prompted to describe 407 all the objects along the ground-truth trajectory, then to 408 predict the location of the target marker based on its rela-409 tionship to these objects. Importantly, the rationale gen-410 erated by GPT-40 must not reference the ground-truth tra-411 jectory itself; the trajectory is only used to guide the gen-412 eration of the rationale. The generated rationale is then 413 automatically validated using GPT-40, assessing both the 414 accuracy of detected objects and the correctness of the ra-415 tionale. This dual-phase prompting strategy has proven to 416 be more reliable than a single-phase prompting approach. 417 The prompts used for rationale generation are provided in 418 Appendix 5.2. The validated rationale is then concate-419 nated with the ground-truth marker ID to enforce a CoT-420 like thinking process. An example of the resulting response 421 is "Think: The candle is most likely located on the shelf on 422 the right side ... Action: 2". Note that the resulting marker 423 selection data used for model fine-tuning relies on images 424 annotated with numerical markers, rather than those anno-425 tated with the ground-truth trajectory. In total, we generated 426 30,352 visual question answer (VQA) pairs of data from 427 104 scenes and 5678 object navigation tasks. This data was 428 used to fine-tune a Llama3.2-11B-Vision model [7] follow-429 ing the configuration recommended by official Llama repos-430 itory. The model was fine-tuned for 3 epochs with a learning 431 rate of 1e-5 and an effective batch size of 128. 432



473 474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

Image Annotatio 2 idate Mark Marker 1 (0.15, 0.11) Marker 3 (0.74, 0.13) Ego Image with Marker Floor Edge Marker Description Data Marker Selection Data පු පු t relative to each Data Generatio Dual Phase tub is most likelv lo ß corner of a kitcher \$ Rationale Verification & Data Processing Data Processing cal markers present in Thought: The bathtub is most likely.

Marker 2 (0.84, 0.32)

Target Marke

ound-truth Trajectory

Figure 2. Pipeline of SFT data collection. The ground-truth trajectory and floor edge are used to extract target marker and candidate markers, respectively. Marker description and selection data is generated for model fine-tuning.

4. Experiments 433

4.1. Experimental Setup 434

We evaluated our method on the navigation tasks using 435 the Habitat 3.0 [19] simulation platform. We adopt simi-436 lar setup as the Habitat ObjectNav 2022 challenge [30] for 437 all the experiments. The action space of the agent consists 438 of: STOP, MOVE_FORWARD, TURN_LEFT, TURN_RIGHT, 439 440 with a forward movement distance of 0.25 meters and a turning angle of 30 degrees per step. For low-level move-441 ment control, we utilized Habitat's built-in shortest-path 442 follower. The maximum number of steps allowed per task 443 is set to 500 by default. Due to limitations in the image 444 445 quality within the Habitat environment and the suboptimal performance of state-of-the-art perception modules, such as 446 GroundingDINO [12], we opted for Habitat's built-in se-447 mantic ground truth with object size conditions as the per-448 ception module. In this context, we can assume that the 449 perception module is sufficiently effective. The LLMs and 450 451 VLMs used in this study was GPT-40 and Llama3.2-11B.

4.2. Datasets 452

Our method is evaluated on the following two object navi-453 gation datasets: 454

455 • Habitat Synthetic Scenes Dataset (HSSD) [8]: We use the HSSD validation dataset to evaluate our method. 456 HSSD consists of 41 scenes and six object goal cat-457 egories: chair, couch, potted plant, bed, 458 toilet, and tv. To ensure task diversity, we select only 459 one episode per scene-object pair. After filtering out er-460 461 roneous episodes-such as cases where the agent's initial

position was in mid-air-the final number of evaluated episodes is 213.

• **HSSD-Hard:** Since some HSSD episodes are relatively easy, with the agent finding the target object in just a few steps, we created a more challenging dataset, HSSD-Hard, by selecting HSSD episodes with longer search distance. We calculated the geodesic distance from the agent's starting point to the target object for each episode and selected the top 50% of episodes with the longest searching distances to form the HSSD-Hard dataset. The total number of episodes in HSSD-Hard is 102.

4.3. Baselines

We compare our method against the following state-of-theart (SOTA) baselines that represent different strategies to address the object navigation problem:

- **PIVOT** [18]: This approach casts the navigation task as an iterative visual question answering problem by annotating the image with numerical markers that represent the navigation subgoals. The method is adapted for the HSSD object navigation tasks. Without a frontier map, visitation memory, and landmark semantic memory, our proposed navigation pipeline degenerates to PIVOT.
- LFG [23]: This method employs frontier-based exploration and LLMs to score potential subgoals and guide the navigation.
- VLFM [31]: The approach utilizes VLMs to generate a language-grounded value map, from which the location with the highest value is selected as the next subgoal for navigation.
- InstructNav [14]: InstructNav introduces a Dynamic Chain of Navigation, breaking down navigation tasks into sequences of actions and landmarks. It employs four value maps, each with different semantic representations, to assist in selecting the appropriate landmark.
- VLMNav [6]: This approach relies on a voxel map built from RGB-D images and the agent's pose to narrow down action space.

To ensure a fair comparison, all experimental setups are conducted under the same conditions described earlier.

4.4. Metrics

We employ the following metrics to evaluate the performance of all the methods:

- Success Rate (SR): A task was deemed successful when the distance between the agent and any viewpoint of the target object was less than 0.2 meters.
- Success Weighted by Path Length (SPL) [1]: This met-507 ric evaluates how efficient the agent's path is compared to 508 the optimal path. SPL is calculated as: 509

$$SPL = \frac{1}{N} \sum_{i=1}^{N} S_i \frac{l_i}{max(p_i, l_i)}$$
 (9) 510

Table 1. Main results. Our proposed method is compared to the baselines on HSSD and HSSD-Hard datasets. SR: Success Rate. SPL: Success Weighted by Path Length. All experiments are conducted using gpt-40.

	HSSD		HSSD-Hard	
	SR ↑	SPL ↑	SR ↑	SPL \uparrow
LFG	0.6244	0.3371	0.6176	0.3454
VLMNav	0.6526	0.3620	0.5294	0.1973
InstructNav	0.7605	0.3722	0.6372	0.4187
VLFM	0.7652	0.5574	0.6078	0.4270
PIVOT	0.7840	0.5658	0.6372	0.4744
Ours	0.8685	0.5788	0.7647	0.4790

511 where l_i is the length of the optimal path for episode *i*. p_i 512 is the length of path taken by the agent. S_i is the binary

is the length of path taken by the agent. S_i is the binary indicator of success in episode *i*.

514 4.5. Main Results

We evaluate our method and all baselines on both HSSD 515 and HSSD-Hard datasets using the same setup described in 516 Section 4.1, with results summarized in Table 1. Perfor-517 518 mance is evaluated based on Success Rate (SR) and Suc-519 cess weighted by Path Length (SPL). While SR indicates the overall ability to find the target object, SPL measures 520 the efficiency of the navigation process. Notably, these two 521 metrics are not correlated, as a method can achieve a high 522 SR by sacrificing navigation efficiency. As shown in Ta-523 ble 1, on the HSSD dataset, our proposed method achieves 524 an SR of 0.8685 and an SPL of 0.5788, both of which are 525 526 higher than all the baseline methods.

Compared to HSSD, tasks in the HSSD-Hard dataset are 527 more challenging due to the relatively longer search dis-528 tance, requiring additional steps to locate target objects. 529 530 As shown in Figure 1, the performance of all methods de-531 creases on the HSSD-Hard dataset, though the impact varies by model. Notably, our method demonstrate an even greater 532 advantages in these more difficult scenarios, achieving an 533 SR that is 12.75% higher than the second-best baseline 534 (PIVOT). Additonally, our method outperforms others in 535 SPL as well, further highlighting its efficiency. These re-536 sults underscore the effectiveness and robustness of our ap-537 proach in tackling challenging navigation tasks. 538

Figure 3 illustrates the memory-augmented decisionmaking process from a real HSSD episode. VLMs, such
as GPT-4o, analyze all memory cues on the image before
reasoning about the most likely location of the target object
and selecting the next marker to explore. The full prompt
and responses for this case are provided in Appendix 5.3.

545 Most failed cases with our method result from reaching
546 the maximum allowable number of steps. This can occur
547 due to the VLM selecting a suboptimal position or the task



Figure 3. Demonstration of memory-augmented decision making process from a real HSSD episode. The full prompts and GPT-40 response are provided in Appendix 5.3

being inherently challenging. Additionally, we occasionally observe that even state-of-the-art VLMs like GPT-40 can exhibit visual hallucinations [11], where they select a marker ID that does not appear in the image or prompt. We have provided two examples in the Appendix 5.4.

Table 2. Results with different VLM models. The fine-tuned Llama3.2-11B model (SFT Llama3.2-11B) trained on our collected data outperforms GPT-40, achieving the best overall performance.

	HSSD		HSSD-Hard	
	SR ↑	SPL \uparrow	SR ↑	SPL \uparrow
GPT-40	0.8685	0.5788	0.7647	0.4790
Vanilla Llama	0.7511	0.5582	0.7352	0.4626
SFT Llama	0.8732	0.5995	0.7843	0.5274

4.6. VLM Model Supervised Fine-tuning (SFT)

To evaluate the impact of the VLM used, we assess the per-
formance of various VLMs within our proposed method.554As shown in Table 2, the vanilla Llama3.2 model performs
significantly worse than GPT-40. Given the substantial dif-
ferences in model size and training data, it is not surpris-
ing that smaller open-source VLMs like Llama3.2-11B un-554

553

617

632

derperform compared to state-of-the-art proprietary models 560 such as GPT-40. Failure analysis reveals that Llama3.2-11B 561 562 is more prone to visual hallucinations and struggles with instruction following, particularly in marker selection and de-563 564 scription tasks. This may stem from a lack of relevant training data for Llama3.2-11B, limiting its ability to generalize 565 effectively in these scenarios. 566

To improve the performance of VLMs, we collected over 567 568 30,000 VQA samples by generating data from 40 new object categories within the HSSD dataset. Rationales gener-569 570 ated using a dual-phase prompting strategy are used to guide the marker selection process. We then fine-tuned Llama3.2-571 572 11B following the approach detailed in Section 3.4. As shown in Table 2, the performance of Llama3.2 is improved 573 significantly after supervised fine-tuning (SFT), surpassing 574 575 GPT-40 in both SR and SPL metrics on the HSSD and HSSD-Hard datasets. These results are particularly promis-576 ing, considering that Llama3.2-11B is substantially smaller 577 than GPT-40 (11B vs. an estimated 175B) and more cost-578 579 effective to train and deploy. This highlights the effec-580 tiveness of our data collection strategy and fine-tuning approach. This performance improvement can be attributed 581 to enhanced instruction adherence and a more effective rea-582 soning process grounded in the given environment. 583



Figure 4. Results for different maximum steps. Experiments are conducted on the HSSD dataset. While SFT Llama3.2-11B consistently achieves the best performance, its advantage over the other two models is most pronounced at 300 maximum steps.

586 587 588

584 Numerous factors influence the performance of the navigation pipeline, particularly the maximum number of al-585 lowed steps. Consequently, we assessed the impact of this factor on performance. As shown in Figure 4, with the maximum number of steps increasing, both the Success Rate (SR) and Success weighted by Path Length (SPL) improve 589 for each model. When the maximum number of steps is 590 relatively low, such as 200 or fewer, all methods exhibit 591 592 suboptimal performance, with only minor differences between models. This is likely because the limited step count 593 prevents the full exploitation of each model's capabilities. 594 Conversely, when the maximum number of steps is large, 595 such as 500 or more, all models reach their respective per-596 formance limits, making the differences between them less 597 598 distinct. It is noteworthy that the performance gap between

different models becomes more apparent when the maxi-599 mum number of steps is set to 300 or 400. 600

4.7. Ablation Study

To assess the contribution of each component in our pro-602 posed method, we conducted an ablation study by evaluat-603 ing its performance without landmark semantic memory or 604 visitation memory. As shown in Table 3, removing either of 605 these memory modules leads to a decline in performance. 606 The frontier map could not be excluded from this analy-607 sis, as it is essential for marker generation. The absence of 608 visitation memory results in redundant exploration in some 609 cases, thus reducing both the success rate and navigation 610 efficiency. Meanwhile, without landmark semantic mem-611 ory, the agent is unable to select a navigation goal globally 612 when no suitable marker is present in its current view, which 613 harms the performannce as well. These findings highlight 614 the crucial role of both memory modules in fully leveraging 615 the potential of our proposed method. 616

Table 3. Ablation study. Removing any of the main components of our design leads to degraded performance on HSSD dataset.

	HSSD		HSSD-Hard	
	SR ↑	SPL \uparrow	SR ↑	$SPL\uparrow$
Ours	0.8685	0.5788	0.7647	0.4790
w/o Visitation	0.8450	0.5761	0.7450	0.4961
w/o Landmark Semantic	0.8356	0.5669	0.7352	0.4795

5. Conclusion

This study proposes an efficient fusion strategy that inte-618 grates task-relevant global memory information with first-619 person perspective information, thereby overcoming the 620 suboptimal solution problem associated with existing mul-621 timodal navigation frameworks due to local observabil-622 ity. Moreover, this method can simultaneously activate and 623 utilize the complex spatial understanding, reasoning, and 624 commonsense reasoning capabilities of VLMs, thus signif-625 icantly enhancing the ability and efficiency of navigation 626 decisions in complex spatial scenarios. Theoretically, en-627 hanced spatial cognitive abilities can reduce the required 628 travel distance and number of actions, thereby increasing 629 the task completion success rate and overall efficiency of 630 the navigation scheme. 631

References

[1] Peter Anderson, Angel Chang, Devendra Singh Chaplot, 633 Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana 634 Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, 635

637

662

663

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

et al. On evaluation of embodied navigation agents. *arXiv* preprint arXiv:1807.06757, 2018. 6

- [2] Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva,
 Alexander Toshev, and Erik Wijmans. Objectnav revisited:
 On evaluation of embodied agents navigating to objects. *arXiv preprint arXiv:2006.13171*, 2020. 1
- [3] Wenzhe Cai, Siyuan Huang, Guangran Cheng, Yuxing Long,
 Peng Gao, Changyin Sun, and Hao Dong. Bridging zeroshot object navigation and foundation models through pixelguided navigation skill. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 5228–5234.
 IEEE, 2024. 1
- [4] Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and
 Cheston Tan. A survey of embodied ai: From simulators to
 research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022. 1
- [5] Rui Fukushima, Kei Ota, Asako Kanezaki, Yoko Sasaki,
 and Yusuke Yoshiyasu. Object memory transformer for object goal navigation. In 2022 International conference on robotics and automation (ICRA), pages 11288–11294. IEEE,
 2022. 1
- [6] Dylan Goetting, Himanshu Gaurav Singh, and Antonio Loquercio. End-to-end navigation with vision language models: Transforming spatial reasoning into question-answering. *arXiv preprint arXiv:2411.05755*, 2024. 2, 6
 - [7] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, and et al. The llama 3 herd of models, 2024. 5
- [8] Mukul Khanna, Yongsen Mao, Hanxiao Jiang, Sanjay 664 665 Haresh, Brennan Shacklett, Dhruv Batra, Alexander Clegg, 666 Eric Undersander, Angel X Chang, and Manolis Savva. Habitat synthetic scenes dataset (hssd-200): An analysis of 667 668 3d scene scale and realism tradeoffs for objectgoal naviga-669 tion. In Proceedings of the IEEE/CVF Conference on Com-670 puter Vision and Pattern Recognition, pages 16384-16393, 671 2024. 6
- [9] Yuxuan Kuang, Hai Lin, and Meng Jiang. Openfmnav: Towards open-set zero-shot object navigation via
 vision-language foundation models. *arXiv preprint arXiv:2402.10670*, 2024. 2
- [10] Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu
 Zhou, Sanjana Srivastava, Cem Gokmen, Tony Lee, Li Erran
 Li, Ruohan Zhang, et al. Embodied agent interface: Benchmarking llms for embodied decision making. *arXiv preprint arXiv:2410.07166*, 2024. 1
- [11] Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng.
 A survey on hallucination in large vision-language models. *arXiv preprint arXiv:2402.00253*, 2024. 7
- [12] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao
 Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang,
 Hang Su, et al. Grounding dino: Marrying dino with
 grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55.
 Springer, 2025. 6
- [13] Yang Liu, Weixing Chen, Yongjie Bai, Xiaodan Liang,Guanbin Li, Wen Gao, and Liang Lin. Aligning cyber space

with physical world: A comprehensive survey on embodied693ai. arXiv preprint arXiv:2407.06886, 2024. 1694

- [14] Yuxing Long, Wenzhe Cai, Hongcheng Wang, Guanqi Zhan, and Hao Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. arXiv preprint arXiv:2406.04882, 2024. 2, 4, 6
- [15] Arjun Majumdar, Gunjan Aggarwal, Bhavika Devnani, Judy Hoffman, and Dhruv Batra. Zson: Zero-shot object-goal navigation using multimodal goal embeddings. *Advances in Neural Information Processing Systems*, 35:32340–32352, 2022. 1
- [16] Christoforos Mavrogiannis, Patrícia Alves-Oliveira, Wil Thomason, and Ross A Knepper. Social momentum: Design and evaluation of a framework for socially competent robot navigation. ACM Transactions on Human-Robot Interaction (THRI), 11(2):1–37, 2022. 1
- [17] Ronja Möller, Antonino Furnari, Sebastiano Battiato, Aki Härmä, and Giovanni Maria Farinella. A survey on humanaware robot navigation. *Robotics and Autonomous Systems*, 145:103837, 2021. 1
- [18] Soroush Nasiriany, Fei Xia, Wenhao Yu, Ted Xiao, Jacky Liang, Ishita Dasgupta, Annie Xie, Danny Driess, Ayzaan Wahid, Zhuo Xu, et al. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *arXiv preprint arXiv:2402.07872*, 2024. 1, 2, 6
- [19] Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey, Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. Habitat 3.0: A co-habitat for humans, avatars and robots. arXiv preprint arXiv:2310.13724, 2023. 6
- [20] Pericle Salvini, Diego Paez-Granados, and Aude Billard. Safety concerns emerging from robots navigating in crowded pedestrian areas. *International Journal of Social Robotics*, 14(2):441–462, 2022. 1
- [21] Adarsh Jagan Sathyamoorthy, Kasun Weerakoon, Mohamed Elnoor, Anuj Zore, Brian Ichter, Fei Xia, Jie Tan, Wenhao Yu, and Dinesh Manocha. Convoi: Context-aware navigation using vision language models in outdoor and indoor environments. arXiv preprint arXiv:2403.15637, 2024. 2
- [22] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9339–9347, 2019. 1
- [23] Dhruv Shah, Michael Robert Equi, Błażej Osiński, Fei Xia, Brian Ichter, and Sergey Levine. Navigation with large language models: Semantic guesswork as a heuristic for planning. In *Conference on Robot Learning*, pages 2683–2699. PMLR, 2023. 2, 3, 6
- [24] Ajay Sridhar, Dhruv Shah, Catherine Glossop, and Sergey Levine. Nomad: Goal masked diffusion policies for navigation and exploration. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 63–70.
 [10] TEEE, 2024. 2
 [24] Ajay Sridhar, Dhruv Shah, Catherine Glossop, and Sergey 743
 [24] Ajay Sridhar, Dhruv Shah, Catherine Glossop, and Sergey 743
 [24] Ajay Sridhar, Dhruv Shah, Catherine Glossop, and Sergey 744
 [25] Teiter Analysis
 [26] Teiter Analysis
 [27] Teiter Analysis
 [26] Teiter Analysis
 [27] Teiter Analysis
 [26] Teiter Analysis
 [27] Teiter Analysis
 [28] Teiter Analysis
 [29] Teiter Analysis
 [29] Teiter Analysis
 [20] Teiter Analysi
- [25] Yujie Tang, Meiling Wang, Yinan Deng, Zibo Zheng, Jingchuan Deng, and Yufeng Yue. Openin: Open-vocabulary
 749

instance-oriented navigation in dynamic domestic environ-750 ments. arXiv preprint arXiv:2501.04279, 2025. 3

- [26] Kristinn Thórisson and Helgi Helgasson. Cognitive archi-752 753 tectures and autonomy: A comparative review. Journal of 754 Artificial General Intelligence, 3(2):1–30, 2012. 1
- 755 [27] Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, and 756 Shuqiang Jiang. Gridmm: Grid memory map for vision-757 and-language navigation. In Proceedings of the IEEE/CVF 758 International Conference on Computer Vision, pages 15625-759 15636, 2023. 1
- [28] Pengying Wu, Yao Mu, Bingxian Wu, Yi Hou, Ji Ma, Shang-760 761 hang Zhang, and Chang Liu. Voronav: Voronoi-based zero-762 shot object navigation with large language model. arXiv 763 preprint arXiv:2401.02695, 2024. 1, 2
- [29] Yuchen Wu, Pengcheng Zhang, Meiying Gu, Jin Zheng, and 764 765 Xiao Bai. Embodied navigation with multi-modal informa-766 tion: A survey from tasks to methodology. Information Fu-767 sion, page 102532, 2024. 1
- [30] Karmesh Yadav, Santhosh Kumar Ramakrishnan, John 768 769 Turner, Aaron Gokaslan, Oleksandr Maksymets, Rishabh 770 Jain, Ram Ramrakhya, Angel X Chang, Alexander Clegg, 771 Manolis Savva, Eric Undersander, Devendra Singh Chap-772 lot, and Dhruv Batra. Habitat challenge 2022. https: 773 //aihabitat.org/challenge/2022/, 2022. 6
- 774 [31] Naoki Yokoyama, Sehoon Ha, Dhruv Batra, Jiuguang Wang, 775 and Bernadette Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In 2024 IEEE In-776 ternational Conference on Robotics and Automation (ICRA), 777 pages 42-48. IEEE, 2024. 2, 6 778
- 779 [32] Fanlong Zeng, Wensheng Gan, Yongheng Wang, Ning Liu, 780 and Philip S Yu. Large language models for robotics: A 781 survey. arXiv preprint arXiv:2311.07226, 2023. 1
- [33] Jiazhao Zhang, Kunyu Wang, Shaoan Wang, Minghan Li, 782 783 Haoran Liu, Songlin Wei, Zhongyuan Wang, Zhizheng Zhang, and He Wang. Uni-navid: A video-based vision-784 785 language-action model for unifying embodied navigation 786 tasks. arXiv preprint arXiv:2412.06224, 2024. 3
- 787 [34] Tianyao Zhang, Xiaoguang Hu, Jin Xiao, and Guofeng 788 Zhang. A survey of visual navigation: From geometry to 789 embodied ai. Engineering Applications of Artificial Intelli-790 gence, 114:105036, 2022. 1
- 791 [35] Kaiwen Zhou, Kaizhi Zheng, Connor Pryor, Yilin Shen, 792 Hongxia Jin, Lise Getoor, and Xin Eric Wang. Esc: Ex-793 ploration with soft commonsense constraints for zero-shot 794 object navigation. In International Conference on Machine Learning, pages 42829-42842. PMLR, 2023. 2, 3 795

Mem2Ego: Empowering Vision-Language Models with Global-to-Ego Memory for Long-Horizon Embodied Navigation

Supplementary Material

5.1. Prompts for inference

Prompt 1

You are an automated system with the capability to analyze the provided image. Based on the numerical markers present in the image, please describe the surrounding environment relative to each marker's position. Ensure that descriptions of different markers are distinct to maintain the uniqueness of each marker. The marker number should not appear in the description. Please adhere to the following format: Marker Number: [insert the number of the first marker here] Description: [provide a description corresponding to the first marker here] Marker Number: [insert the number of the second marker here] Description: [provide a description corresponding to the second marker here]

Marker Number: [insert the number of the last marker here] Description: [provide a description corresponding to the last marker here]

Prompt 2

Based on the provided descriptions for each number, please select at most three number whose corresponding descriptions are most likely to help identify the {goal_object}. {Number: Description; Number: Description; ... Number: Description} If the total number is less than 3, please use -1 to occupy the empty position. Please adhere to the following format for the output: Number List: [first number, second number, third number]

799

Prompt 3

You are a robot and after 360 degrees observation, you can see the given panorama image. The panorama image combines 4 images from different angles. Your task is to find the {goal_object}. Based on the numerical markers in the image, select one of these numbers to move next. If you're not confident in moving to the marker to find the {goal_object}, you can choose one of the numerical markers located outside of this image. The descriptions of these markers are as follows: {the top 3 numbers with descriptions, corresponding to the response of Prompt 2}. If you're still not confident in moving to the marker to find the {goal_object}, your action should be 'None'. The blue circle marker on the floor indicates the previously explored position. Ιt is better to choose a numeric marker that is not close to the blue circle marker. Please note all closed doors cannot be opened. Please follow the format like this, Thought: [put your step-by-step thinking process here] Action: [put a single marker id or None here]

797

803

5.2. Prompts for SFT data generation

Prompt for dual-phase rational generation

You are given an image with a red movement trajectory on it. Please first identify the objects near the red line in the given image. If there is no red trajectory in the image, please directly return "None". Second, knowning that {goal_object} could be found after following the red trajectory, you need to predict the location of goal_object or the region where {goal_object} could be most likely located. This can be achieved by reasonably imagining the unseen areas after the red trajectory based on the room layout. **Do not mention the red trajectory/line or "the image" in your output!** Please structure your output in the following way: OBJECTS_RED_LINE: LOCATION_PREDICTION_AND_REASONING:

Prompt for rationale filtering

You are given an image with a movement trajectory marked in a red line. Please first verify if all of the objects in a given list are present near the red line in the given image. If there is no red line in the image or any of the objects not present, please ignore the rest and directly return "NONE". Second, verify if the reasonings of why {goal_object} may be put close to the objects in the list. A good reasoning should be logical and perfectly reflect common sense knowledge. A good reasoning gives convincing reasons while a bad reasoning gives vague or untruthful reasons. If the reasonings are good, output "GOOD REASONINGS", otherwise, output "BAD REASONINGS". Example of a good reasoning: "The book is most likely located on the shelves in the background. The shelves are a

common place for storing books, and they are visible in the room at the end of the path.". Example of a bad reasoning: "The mirror is most likely located on the dark wall to the right of the doorway. This is inferred from the visible portion of the mirror reflecting the room, indicating its position on the dark wall.". Please structure your output in the following way: OBJECTS_PRESENCE_CHECK: REASONING_CHECK: Object list: object_list Reasonings: reasonings

5.3. Case Study

Prompt and response for case study

Prompt: You are a robot and after 360 degrees observation, you can see the given panorama image. The panorama image combines 4 images from different angles. Your task is to find the tv. Based on the numerical markers in the image, select one of these numbers to move next. If you're not confident in moving to the marker to find the tv, you can choose one of the numerical markers located outside of this image. The descriptions of these markers are as follows: 97: The environment near this marker is a well-furnished bedroom featuring a bed with grey and white bedding. There is a grey headboard attached to the wall, from which two contemporary pendant lights hang. Next to the bedside is a white nightstand with drawers. In the background, large windows or glass doors bring light into the room. To the far right, through the doorway, part of an adjacent room is visible where a small trolley or shelf holding white objects can be seen. 98: The marker is located near the doorway of a

807

CVPR

#13

room. To the left of the marker, there is a green potted plant. Through the doorway, the room is furnished with a bed positioned near the wall with large windows that extend almost to the floor, allowing for an expansive view of the outside. Additionally, there is a chair or possibly a sofa situated to the left inside the room. The floor is covered with a dark blue surface adding a touch of color to the environment. If you're still not confident in moving to the marker to find the tv, your action should be 'None'. The blue circle marker on the floor indicates the previously explored position. It is better to choose a numeric marker that is not close to the blue circle marker. Please note all closed doors cannot be opened. Please follow the format like this, Thought: [put your step-by-step thinking process here] Action: [put a single marker id or None here]

Response: Thought: I need to find the TV in this panorama image. The available markers inside the image are 0, 1, 2, 3, and 4. The markers outside of the image, specifically 97 and 98, describe a bedroom and another room with a bed. Since a TV is more likely to be in a living room or common area rather than a bedroom, and given that we're currently in what appears to be a common area or living room (as indicated by the dining area and kitchen), it makes more sense to explore the nearby markers in the current area. Marker 0: Leads to a kitchen area, not likely to have a TV. Marker 1: Near a closed door with a blue circle, so less likely. Marker 2: In the direction of a long room which might be a living room and could have a TV. Marker 3: Looks like



Figure 5. Visual Hallucination Example 1



Figure 6. Visual Hallucination Example 2

another part of the same long room, potentially a living area. Marker 4: Seems to be in the dining area, also less likely to have a TV. Action: 2

5.4. Visual Hallucinations

Visual Hallucination Example 1

Thought: 1. The top-left, bottom-left, and bottom-right sections of the panorama show closed doors and walls. 2. The top-right section shows an open doorway leading to another room that might be more likely to have a toilet. 3. The numeric marker "2" is located near the open doorway in the top-right section. Action: 2 809

811

CVPR

#13

CVPR

#13

Visual Hallucination Example 2

Thought: The given panorama shows parts of a room and possibly adjacent areas. The upper left image shows a workspace with a computer, the upper right shows a hallway and another room, the lower left has a part of a bed, and the lower right includes windows with curtains. There is no indication of a couch in these views. The only promising lead is the room visible through the doorway marked with a green number "1" in the upper right image, which might show a larger area including a couch. Action: 1