

HeckNav: Interpretable Heuristic Knowledge Navigation via Bayesian Probabilistic Inference

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Object goal navigation (ObjectNav) is a well defined and*
 002 *fundamental problem for robotic to navigate, that start-*
 003 *ing in a random point at unseen environment to a given*
 004 *object. Recent advanced study in ObjectNav uses Vision-*
 005 *Language Models (VLMs) or multimodal large language*
 006 *models (MLLMs) for object detection and decision-making.*
 007 *However, due to the "black-box" nature and implicit rea-*
 008 *soning of VLMs, current ObjectNav systems often rely on*
 009 *the opaque internal capabilities of the models or exhaus-*
 010 *tive prompt engineering, lacking transparency and explicit*
 011 *planning. In this paper, we propose HeckNav, a zero-shot,*
 012 *heuristic-based knowledge navigation module that enhance*
 013 *VLMs with Knowledge Graph (KG) and Bayesian prob-*
 014 *abilistic inference. To enhance planning efficiency while*
 015 *maintaining a minimal footprint, we introduce a lightweight*
 016 *KG module that employ Bayesian inference to derive navi-*
 017 *gational plans based on the robot's observed environment.*
 018 *Furthermore, our framework provides unique level of inter-*
 019 *pretability by visualizing the agent's environmental aware-*
 020 *ness and its path-finding process through an abstract graph*
 021 *based on heuristic value preferences. Experimental re-*
 022 *sults demonstrate that HeckNav achieves significant im-*
 023 *provement on the HM3D and MP3D datasets especially on*
 024 *smaller VLM.*

025 1. Introduction

026 "Knowledge is crude." - David Papineau

027 ObjectNav [4] requires an embodied agent to navigate
 028 toward a specific object category within an unseen envi-
 029 ronment. This task tests the agent's capacity for object
 030 detection, spatial awareness, sequential planning, and au-
 031 tonomous execution. The workflow of objectNav can be
 032 divided into:

- 033 • **Localization:** agent is aware of surrounding environment
 034 with semantic meaning.
- 035 • **Mapping:** agent remembers the area that has been ex-

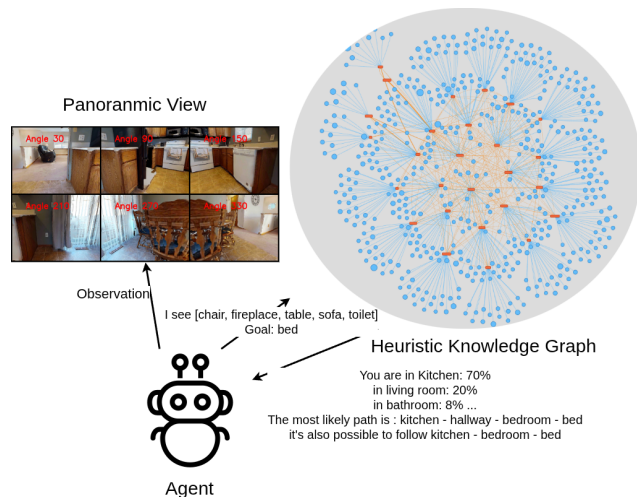


Figure 1. Overview of HeckNav module. It capture the statistical probability between regions and objects, make Bayesian inference over linguistic knowledge graph possible

explored and unexplored.

- **Planning:** agent makes plan based on the known information and area of interest for long-term.
- **Decision making:** agent use judgement to make the decision based on plan and surrounding area for the following step.
- **Goal detection:** the agent stops when the goal is reached or when it believes that the goal is reached.

In the very early stage navigation problem, the pioneer like DD-PPO [26] trains policies using reinforcement learning algorithms for decision making. Frameworks in the era doesn't come with planning process, they just merely predict next step based on the trained policy. Such method is demanding for training data, prolonged training times, and high computational resources.

Early milestones in objectNav, such as ZSON[19], use of multimodal embeddings, CLIP[23], to transfer policies trained on image navigation to open-vocabulary targets. However, ZSON relies heavily on implicit end-to-end rein-

055	forcement learning policies, which often struggle with long-	108
056	horizon exploration and lack explicit spatial memory.	
057	To address these exploration inefficiencies, subsequent	109
058	works like VLFM[33] integrated simple VLM-driven seman-	110
059	tic scoring with classical frontier-based mapping.	111
060	Based on the map, they capture images and use BLIP-2[13]	112
061	to get cosine-similarity score. The similarity score incor-	113
062	porates with the field of view to build value map. Using	114
063	these data, adapt a pre-trained PointNav policy. However,	115
064	this framework depends on BLIP to evaluate the similarity,	116
065	which would cause disaster, as BLIP is not powerful enough	117
066	to infer from an unrelated item to the goal.	118
067	Following the era of CLIP-based scoring, the rapid evolu-	119
068	tion of MLLMs such as ChatGPT and Gemini [6, 21],	120
069	or complex question-answering VLMs like LLaVa and	
070	Qwen[16, 22], a profound paradigm shift in Embodied AI	121
071	that makes vision-processing, planning, and decision mak-	
072	ing integrated into single model. Unlike simple VLMs	
073	limited to embedding alignments, MLLMs exhibit emer-	
074	gent reasoning, rich contextual understanding. On the other	
075	hand, simple VLMs with smaller parameters can be de-	
076	ployed on robots locally. However, the limited context win-	
077	ow of models makes each step as an isolated query as the	
078	navigation horizon extends. VLM requires robust external	
079	memory mechanism.	
080	With the MLLM advancement, WMNav[20] pushed the	
081	boundaries by framing the MLLM as a 'World Model'. By	
082	maintaining a Curiosity Value Map, WMNav enables the	
083	agent to simulate future states and evaluate potential out-	
084	comes before acting, thereby reducing trial-and-error. How-	
085	ever, such world-model-based approaches implicitly rely	
086	on the MLLM's internal simulation capabilities, which in-	
087	currs high computational overhead, suffers from hallucina-	
088	tion risks, and lacks a transparent, symbolic representation	
089	of common sense. In the meanwhile, the performance of	
090	WMNav decrease dramatically when switch to small VLM,	
091	makes it impossible to deploy on local robot.	
092	To bridge these gaps, our HeckNav framework intro-	
093	duces a highly optimized, lightweight knowledge graph	
094	synergized with Bayesian probabilistic inference. Instead	
095	of relying on expensive MLLM-based forward simula-	
096	tions, HeckNav leverages explicit, structured heuristic pri-	
097	ors to logically deduce target locations (e.g., inferring a mi-	
098	crowave's location based on a visible sink). This approach	
099	not only reduces computational burden but also provides an	
100	unique, interpretable visualization of the navigation reason-	
101	ing process.	
102	Extensive experiments on the standard HM3Dv1 and	
103	MP3D datasets demonstrate the superiority of HeckNav.	
104	Notably, our framework achieves an absolute improvement	
105	of 1.4% in Success Rate (SR) on HM3D, 3.1% SR on	
106	MP3D over the benchmark using Gemini-2.0-Flash and	
107	over 3.1% SR using Qwen7b/8b, while maintaining an	
	ultra-low memory footprint of $< 2.2MB$	108
	The main contributions of our work are summarized as	109
	follows:	110
	• A Novel Interpretable Module: We introduce an inter-	111
	pretable visualization mechanism that maps the agent's	112
	spatial awareness and decision-making logic for long-	113
	term.	114
	• Extreme Memory Efficiency: We introduce a highly op-	115
	timized knowledge graph without the massive computa-	116
	tional overhead.	117
	• SOTA Performance: HeckNav improves success rates	118
	standard benchmarks (HM3D and MP3D) compared with	119
	SOTA framework especially on small VLM.'	120
	2. Related Work	121
	ObjectNav Datasets. The development of Embodied AI	122
	is deeply rooted in the availability of high-fidelity simula-	123
	tion environments. A significant lineage of datasets, includ-	124
	ing Gibson [28], Matterport3D (MP3D) [3], Replica [2],	125
	HM3D [24], HM3Dv2 [30] and HM3D-OVON [34], are	126
	reconstructed from real-world physical scans. Among	127
	these, MP3D is particularly distinguished by its exhaus-	128
	tive, human-annotated semantic metadata, providing gran-	129
	ular ground-truth associations between object and room.	130
	Here, we conduct a systematic statistical representation of	131
	these human-annotated semantic distributions, transform-	132
	ing metadata into a probabilistic knowledge graph that	133
	guides the agent for long-term planning.	134
	Foundation Models for Semantic Navigation. The in-	135
	tegration of Large Language Models (LLMs) and MLLMs	136
	has revolutionized the decision-making pipeline in Ob-	137
	jectNav. Beyond the predictive world models like WM-	138
	Nav [20], several works explore the reasoning capabilities	139
	of foundation models. For instance, L3MVN [35] utilize	140
	LLMs to score frontiers by evaluating the semantic rele-	141
	vance of room-object associations. Another parallel re-	142
	search track focuses on persistent spatial representations.	143
	VLMMaps [9] and CoW [5] store pixel-level VLM embed-	144
	dings in a 2D grid, enabling open-vocabulary spatial query-	145
	ing. InstructNav [18], NavCoT [15] introduces a Chain-of-	146
	Thought prompting strategy for MLLMs to enhance step-	147
	by-step navigation reasoning. VISTA v2[10] uses diffusion	148
	models to evaluate the path quality by "imagine" the future	149
	scene. While these methods demonstrate impressive zero-	150
	shot capabilities, they often encounter a "performance de-	151
	terioration" using smaller model, making it hard to deploy	152
	on local robot. HeckNav addresses these limitations by distill-	153
	ing these complex reasoning processes into a lightweight	154
	Bayesian Knowledge Graph, effectively bridging the gap	155
	between high-level MLLM perception and efficient, low-	156
	level symbolic execution for long-term planning.	157
	Graph and Graph-Based Navigation. There are al-	158
	ready plenty of study about the 3D-structure graph [1,	159

7, 12] that extract the semantic representation into abstract graph as prior. Utilizing structured prior to encode environmental common sense has long been a pursuit in Embodied AI. Frameworks such as AKGVP [29], NavRAG[25], and HOZPlus [37] have pioneered the integration of **offline** Knowledge Graphs (KGs) to provide the agent with a global semantic context. On the other hand, there are many **online** graph frameworks like TopoNav[17], TSGM[11], UniGoal[32], DGN[14]and SG-Nav[31] that served as memory mechanism. In this paper, the graph we discuss will be severed as **offline prior**. There is a fundamental limitation of these existing KG-based frameworks is their reliance on binary or unweighted relational structures. While they successfully establish that certain objects "can" coexist (e.g., a stove is in a kitchen), they fail to explicitly model the *statistical strength* or *probabilistic tendency* of these relationships. In these graphs, all edges are treated with equal significance, neglecting the fact that some spatial correlations are far more predictive than others in real-world distributions. **HeckNav** addresses this gap by introducing a Bayesian probabilistic layer atop the knowledge graph. By deriving edge weights from the statistical analysis of human-annotated datasets like MP3D, our framework quantifies the "affinity" between entities, allowing the agent to prioritize exploration targets based on the rigorous mathematical likelihood rather than simple connectivity.

3. Method

3.1. Statistical Knowledge Grounding

To empower the agent with structured common sense, we construct a lightweight, multi-relational Knowledge Graph (KG) derived from the semantic metadata of the MP3D dataset. Unlike implicit end-to-end priors, our KG is explicitly built to encode both topological layouts and probabilistic semantic distributions.

3.1.1. Multi-relational Graph Construction

Formally, we define our Knowledge Graph as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. The node set $\mathcal{V} = \mathcal{V}_R \cup \mathcal{V}_O$ consists of room nodes \mathcal{V}_R (e.g., *bedroom*, *kitchen*) and object nodes \mathcal{V}_O (e.g., *bed*, *cabinet*). The edge set $\mathcal{E} = \mathcal{E}_{nav} \cup \mathcal{E}_{prob}$ encompasses two distinct types of relationships:

- **Navigation Edges (\mathcal{E}_{nav}):** Undirected edges connecting adjacent room nodes based on spatial annotations (e.g., portals/doors). These define the topological navigability of the environment.
- **Probabilistic Semantic Edges (\mathcal{E}_{prob}):** Directed edges connecting object nodes to room nodes. The weight of these edges represents the statistical correlation probability, bridging semantic concepts with spatial containers.

3.1.2. Probabilistic Edge Formulation

To quantify the semantic edges \mathcal{E}_{prob} , we aggregate the spatial distributions from the raw metadata. Let $C(o, r)$ denote the total occurrences of object $o \in \mathcal{V}_O$ residing within room $r \in \mathcal{V}_R$. We formulate the conditional probability $P(r|o)$ —the likelihood that an observed object o belongs to a room type r —using a normalized frequency distribution.

3.2. Heuristic Navigation via Graph Reasoning

During navigation, HeckNav dynamically queries this lightweight KG to form long-term reasoning chains. The decision-making process is divided into location belief updating and semantic path generation.

3.2.1. Zero-Shot Perception and Room Belief Update

At each navigation step t , the agent faces cognitive tasks: localizing its current semantic context. HeckNav accomplishes this by deriving two explicit probability distributions over all room categories \mathcal{V}_R , conditioned on real-time observations and the specified goal.

Current Room Distribution ($P(R_{curr})$): First, the VLM parses the current panoramic RGB observation to identify a set of salient anchor objects \mathcal{O}_{obs} . Instead of relying on geometric coordinates, the agent infers its current semantic location by aggregating the probabilistic evidence from these observed anchors. First, the model initialize the current room distribution by calculating:

$$P(R_i)_{initial} = \frac{f(R_i, \mathcal{O}_{obs})}{\sum_{j \in \mathcal{R}} f(R_j, \mathcal{O}_{obs})} \quad (1)$$

The function $f(R_i, o_{obs})$ represents the likelihood of observing o_{obs} given that the agent is in room R_i . Based on the current belief distribution over all rooms, the scores are L1-normalized across all room nodes. After movement, the belief is updated as follows:

$$P(R_i | \mathcal{O}_{obs}) = \frac{P(\mathcal{O}_{obs} | R_i)P(R_i)}{P(\mathcal{O}_{obs})} \quad (2)$$

Here, the denominator is a dynamic normalization constant. It fluctuates based on the semantic density of the current observation.

$$P(\mathcal{O}_{obs}) = \sum_{j \in \mathcal{R}} P(\mathcal{O}_{obs} | R_j)P(R_j) \quad (3)$$

Goal Room Distribution ($P(R_{goal})$): Simultaneously, the agent must estimate the destination. Using the given target object o_{target} , we query the probabilistic semantic edges \mathcal{E}_{prob} in our Knowledge Graph. The unnormalized belief for the target residing in room r is derived directly from the statistical prior: $S_{goal}(r) = P(r|o_{target})$. This is similarly normalized to yield the goal room distribution, $P(R_{goal})$.

252 By transforming both the observations and the text-
253 specified goal into unified room-level probability distribu-
254 tions, HeckNav establishes the source and destination priors
255 required for the subsequent semantic path generation.

256 3.2.2. Semantic Navigation Chain Generation

257 With the dual candidate sets \mathcal{R}_{curr}^* and \mathcal{R}_{goal}^* established,
258 HeckNav shifts from perception to cognitive planning. We
259 formulate this as a semantic path search problem over the
260 topological navigation edges \mathcal{E}_{nav} within our Knowledge
261 Graph.

262 For any presumed source room $r_s \in \mathcal{R}_{curr}^*$ and goal
263 room $r_g \in \mathcal{R}_{goal}^*$, a candidate semantic path $\Gamma =$
264 (r_s, r_1, \dots, r_g) is defined as a sequence of topologically
265 connected rooms. Let $N = |\Gamma|$ denote the number of hops
266 (i.e., edges) in this path. A naive approach to evaluate the
267 traversability of Γ would be the cumulative product of all
268 edge probabilities. However, since transition probabilities
269 strictly reside in $(0, 1]$, a pure multiplicative formulation in-
270 herently suffers from the *multi-hop attenuation problem*—
271 it penalizes longer, multi-hop trajectories, biasing the agent
272 toward strictly myopic, short-distance explorations.

273 To inherently decouple the semantic confidence from the
274 path length, we introduce a length-normalized transition
275 scoring function. By taking the N -th root of the cumula-
276 tive product, we effectively compute the **geometric mean**
277 of the transition probabilities. The transition score of a path
278 Γ is formulated as:

$$279 \quad S_{trans}(\Gamma) = \left(\prod_{i=1}^N P_{nav}(r_i|r_{i-1}) \right)^{\frac{1}{N}} \quad (4)$$

280 where $P_{nav}(r_i|r_{i-1})$ represents the navigation edge proba-
281 bility between adjacent rooms derived from the prior graph.

282 The overall confidence of a complete semantic chain
283 $\mathcal{C}_{s \rightarrow g}$ is then synergized by integrating the initial location
284 belief, the length-normalized path score, and the goal esti-
285 mation:

$$286 \quad S_{chain}(\mathcal{C}_{s \rightarrow g}) = P(r_s) \cdot S_{trans}(\Gamma(r_s \rightarrow r_g)) \cdot P(r_g) \quad (5)$$

287 The framework evaluates all possible paths between the
288 Top- K source and goal candidates, selecting the chain with
289 the highest overall score $\arg \max S_{chain}$. The agent then
290 can decompose the task and decide next step.

291 3.3. Transparent and Interpretable Reasoning

292 Unlike end-to-end black-box models or implicit world mod-
293 els, HeckNav’s decision-making process is entirely trans-
294 parent. Because both the environmental prior \mathcal{G} and the dy-
295 namic semantic chains are explicitly formulated, the agent’s
296 cognitive state can be visualized in real-time. By overlay-
297 ing the active room beliefs $P(R_{curr})$, $P(R_{goal})$, and the
298 highest-scoring semantic path onto the topological graph,

humans can intuitively trace the exact probabilistic reason-
ing that led to the agent’s navigational choices. This explicit
symbolic structure enhances the trust and debuggability of
the Embodied AI system.

4. Experiment

4.1. Data and Evaluation Metrics

Knowledge Graph Optimization and Specifications We
meticulously optimized the raw semantic categories ex-
tracted from the MP3D dataset. Architectural elements
that are ubiquitous across all indoor environments—such as
walls, ceilings, floors, and doors—are discarded for goal-
oriented semantic reasoning. Consequently, we filtered out
these non-discriminative entities into different categories
during the graph construction phase.

The resulting refined Knowledge Graph is highly com-
pact, comprising exactly 52 nodes: 26 room category nodes
(\mathcal{V}_R) and 26 salient object category nodes (\mathcal{V}_O). These
nodes are interconnected by 1,421 directed edges (\mathcal{E}_{nav} and
 \mathcal{E}_{prob}), encapsulating both the topological navigation links
and the probabilistic semantic affinities. In our implemen-
tation, the module that contains graph and retrieval func-
tion only occupies a mere $< 2.2MB$ of runtime mem-
ory and each inference step takes $< 2ms$ (module test on
 $i5 - 12600K$ with $32GB$ memory). This extremely low
memory footprint validates our claim of HeckNav being an
lightweight, plug-and-play module. The category we used
for object nodes is based on **mpcat40**.

Dataset The first version of HM3D[24] comes with 2000
validation episodes and 6 goal object categories such as
sofa, bed on 20 environments. MP3D[3] has 2195 valida-
tion episodes and 21 goal object categories on 11 scenes.

Metrics Like other frameworks, we adopt Success Rate
(SR) and Success Rate Weighted by Inverse Path Length
(SPL). SR is a robustness test and SPL is efficiency test

4.2. Experimental Setup

MLLM Perception Engine To execute the zero-shot se-
mantic perception and anchor extraction, we intentionally
deploy **Gemini 2.0-Flash** as our primary Vision-Language
Model (VLM) backend as we explicitly avoid performance
improvement by merely utilizing heavier, more powerful
models such as Gemini 1.5-Pro or Gemini 3.0-Flash.

Implementation Details To evaluate the versatility and ef-
fectiveness of our proposed Bayesian Knowledge Graph,
we implement **HeckNav** as a lightweight, plug-and-play
cognitive enhancer built upon the state-of-the-art **WMNav**
architecture.

Crucially, we do **not** alter or bypass WMNav’s origi-
nal computational pipeline, **nor** change any config from
the WMNav. The core sensory modules and the com-
plex curiosity-driven world-model simulations are fully re-

Method	VLM	Reason/Plan	MP3D		HM3D	
			SR (%) ↑	SPL (%) ↑	SR (%) ↑	SPL (%) ↑
CoW [5]	CLIP	FM	-	-	9.2	4.9
VLFM [33]	BLIP	FM + DDPPPO	36.4	17.5	52.5	30.4
ESC [38]	GLIP	DeBERTa v3	28.7	14.2	39.2	22.3
L3MVN [35]	Mask R-CNN	RoBERTa-large	-	-	50.4	23.1
VoroNav [27]	BLIP	GPT-3.5	-	-	42.0	26.0
HOZ++ [37]	CNN+GCN	GPT-R	37.0	15.2	-	-
SG-Nav [31]	LLaVA-1.6 (7B)	GPT-4	40.2	16.0	54.0	24.9
WMNav [20]	Gemini-1.5-Pro	Gemini-1.5-Pro	45.4	17.2	58.1	31.2
WMNav [20]	Gemini-2.0-Flash	Gemini-2.0-Flash	43.3	15.3	54.1	28.5
HeckNav	Gemini-2.0-Flash	Gemini-2.0-Flash	44.7	15.4	57.2	29.9

Table 1. Comparison of performance on the MP3D and HM3D datasets. We separate the models used for detection and planning.

349 tained to ensure a strictly fair baseline comparison. Instead,
350 our HeckNav module is integrated orthogonally: while
351 WMNav simulates implicit future states, our pre-computed
352 topological and probabilistic graph acts as an explicit cog-
353 nitive guide. By injecting our *Semantic Navigation Chains*
354 into the agent’s reasoning loop, HeckNav provides statisti-
355 cal priors that assist and constrain the heavy MLLM compu-
356 tations, effectively preventing the agent from hallucinating
357 highly improbable spatial layouts.

358 4.3. Comparison with SOTA Methods

359 As HeckNav’s KG extract from MP3D, we will compare the
360 performance of HeckNav method on MP3D and HM3Dv1
361 dataset for generalizability in Table 1. Because Gemini-1.5-
362 pro api is not access at the time of experiment, we pick a
363 less powerful model to compare with WMNav. **FM** means
364 **Frontier Map**, DDPPPO means **Decentralized Distributed**
365 **Proximal Policy Optimization**

366 The result shows that even with model degradation from
367 1.5-pro to 2.0-flash, HeckNav has a good improvement over
368 WMNav, and still outperform most SOTA models. This re-
369 sults shows that explicit reasoning module would benefit the
370 VLM’s performance.

371 In addition to ultra-lightweight and fast inference, the
372 module, based on Bayesian inference, can track the proba-
373 bility of possible paths and generate a colored graph rea-
374 soning process for explainability.

375 4.4. Robustness Study

376 4.4.1. Performance on Different Models and Datasets

377 As the ultimate goal is to deploy the agent on local,
378 resource-constrained devices, we evaluate HeckNav against
379 the WMNav baseline using Qwen-VL series in Table 2,
380 across the HM3Dv1 and HM3Dv2 datasets. Traditional
381 navigation metrics fail to analysis the model hallucina-

tion. Therefore, in this comparison, we introduce three
additional perception-centric metrics: **Hallucination Rate**
(Hal.), **Precision (Prec.)**, and **F1-score**.

- **Hallucination Rate (Hal.):** Every VLMs have potential
of visual hallucinations, often claiming to see a target ob-
ject that is not present. As a result, the success rate won’t
achieve 100% as the robot already stop and claim a false
positive stop.
- **Precision (Prec.):** Addition to Hallucination Rate, this
metric measures that when robot claim it reaches a goal,
what’s the actual rate it has reach the goal.
- **F1-score:** By providing the harmonic mean of precision
and recall, the F1-score offers a holistic evaluation of
the agent’s semantic perception capabilities during explo-
ration.

By incorporating these metrics alongside SR and SPL,
we prove that HeckNav not only improves physical naviga-
tion efficiency but enhances the semantic and cognitive re-
liability of open-source local models. We also find that use
a more powerful VLM model will not only lead to higher
success rate but also higher hallucination rate.¹

4.4.2. HM3Dv2 Comparison

Due to curiosity about the performance of HeckNav on lat-
est Qwen3-VL-8B-Instruct, we compare it with other SOTA
models that has HM3Dv2 dataset tested in table 3. We
find that even with local deployed Qwen3 8B model, the
performance of HeckNav is still competitive to those GPT
or Gemini based API-calling models. It outperforms those
GPT-4 based framewoks on the success rate, and the SPL
also achieves second place in the table. What’s more, the
entire framework takes about 23GB of VRAM including

¹Due to computational constraints, these experiments are not exhaus-
tive. However, both HeckNav and the baseline are tested on the exact same
500 episodes to ensure fairness.

Model	Method	HM3Dv1 [24]					HM3Dv2 [30]				
		SR↑	SPL↑	Hal.↓	Prec.↑	F1↑	SR↑	SPL↑	Hal.↓	Prec.↑	F1↑
Qwen-2.5(7B)	WMNav	39.2	18.4	13.4	74.5	51.4	51.6	21.1	8.0	86.6	64.7
	HeckNav	45.2	20.5	13.8	76.6	56.9	53.6	22.9	6.0	89.9	67.2
Qwen-3(8B)	WMNav	55.1	28.0	17.2	76.2	64.0	66.0	28.1	11.1	85.6	74.5
	HeckNav	56.7	29.7	16.0	78.0	65.7	69.3	29.9	10.5	86.9	77.1

Table 2. Performance evaluation using Qwen on HM3D. All metrics are reported as percentages (%). By using explicit reasoning with the HeckNav graph, the agent consistently achieves higher SR and semantical F1-scores are also notably outperform across different datasets.

413 deploy Qwen3-VL-8B, which makes it possible to transfer
414 to local device later on.²

415 4.5. Failure Case Analysis

416 Despite the substantial performance gains, we conducted a
417 thorough investigation into the failure cases of **HeckNav** to
418 identify remaining challenges in semantic navigation. Our
419 analysis reveals three primary modes of failure:

- 420 1. **Dataset Artifacts and Initialization Constraints:** In
421 certain MP3D episodes, the agent spawns exterior to the
422 building structure, devoid of any meaningful indoor seman-
423 tic reference objects. This initialization fundamentally
424 contradicts the core premise of the indoor Object-
425 Nav task. Without valid visual anchors ($\mathcal{O}_{obs} = \emptyset$), our
426 Bayesian graph struggles to formulate an initial seman-
427 tic belief $P(R_{curr})$, inadvertently degrading the overall
428 success rate for these outlier episodes.
- 429 2. **Local Trapping in Expansive Regions (Memory
430 Mechanism Deficits):** When navigating massive, open-
431 plan spaces (e.g., churches or grand halls), the agent may
432 exhibit localized roaming behaviors, repeatedly visiting
433 the same sub-regions before eventually breaking out of
434 the loop. Although the agent often ultimately reaches the
435 target, this redundant exploration severely penalizes the
436 SPL metric.
- 437 3. **Out-of-Distribution (OOD) Semantic Contexts:** Cer-
438 tain episodes require navigating toward highly uncom-

²23GB of VRAM is tested on RTX 6000 pro with 22%
gpu_memory_utilization plus running simulation. It's not the exhaustive
experiments that are conducted by running 50 instance at the same time.

Method	VLM (Largest)	SR (%) ↑	SPL (%) ↑
L3MVN [35]	Mask R-CNN	36.3	15.7
SG-Nav [31]	GPT-4	54.0	24.9
InstructNav [18]	GPT-4V	58.0	20.9
SGM [36]	GPT-4	60.2	30.8
DORAEMON [8]	Gemini-1.5-Pro	66.5	20.6
HeckNav (Ours)	Qwen3-VL (8B)	64.6	28.5

Table 3. Comparison with SOTA methods on the HM3Dv2.

439 mon target objects (e.g., *gym equipment*) starting from
440 atypical locations (e.g., *a lobby*). These rare seman-
441 tic pairs possess near-zero co-occurrence probabilities in
442 our statistical prior. Consequently, the agent may be-
443 come trapped in the starting region, unable to establish a
444 confident semantic chain. Notably, these OOD scenarios
445 are also nearly impossible for baseline frameworks rely-
446 ing purely on zero-shot LLM/MLLM reasoning, high-
447 lighting a shared bottleneck in current Embodied AI.

- 448 4. **Complexity of Environment Detection:** The idea of
449 knowledge graph is to separate rooms by the mostly cor-
450 related items detected. However boundaries between
451 rooms are not clear (e.g. livingroom and kitchen may
452 be connected without walls), or the object behind doors or
453 windows will be detected to undermine the result of lo-
454 calization.

455 This analysis suggests that future iterations could heav-
456 ily benefit from an explicit **spatial memory mechanism**,
457 where the agent dynamically updates its topological graph
458 by masking previously visited nodes and penalizing the
459 traversal weights of explored regions. Such an episodic
460 memory would not only prevent redundant looping to dras-
461 tically improve SPL, but also force the agent to explore
462 novel frontiers when trapped in OOD semantic contexts.

463 5. Conclusion

464 In this paper, we introduce **HeckNav**, a zero-shot,
465 neuro-symbolic navigation module that integrates explicit
466 Bayesian probabilistic inference with MLLMs. By extract-
467 ing priors into an ultra-lightweight Knowledge Graph mod-
468 ule with , HeckNav effectively offloads complex spatial
469 reasoning from the MLLM to generate deterministic, length-
470 normalized semantic chains. Operating as a plug-and-play
471 module over the state-of-the-art WMNav, our framework
472 achieves up to a 3.1% absolute improvement in Success
473 Rate using Gemini 2.0-Flash, and consistently outperform
474 baseline (1.6% to 6.0% SR, 1.8% SPL) using local models
475 using Qwen. Furthermore, our explicit graph formulation
476 enables real-time visualization of path-finding decisions,
477 leveraging the interpretability gap in current black-box Em-
478 bodied AI systems.

479

References

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

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

- [1] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5664–5673, 2019. 2
- [2] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R. Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2
- [3] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. 2, 4
- [4] Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *In Neural Information Processing Systems (NeurIPS)*, 2020. 1
- [5] Samir Yitzhak Gadre, Mitchell Wortsman, Gabriel Ilharco, Ludwig Schmidt, and Shuran Song. Cows on pasture: Bases and benchmarks for language-driven zero-shot object navigation. *CVPR*, 2023. 2, 5
- [6] Gemini Team, Google. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. 2
- [7] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5021–5028. IEEE, 2024. 3
- [8] Tianjun Gu, Linfeng Li, Xuhong Wang, Chenghua Gong, Jingyu Gong, Zhizhong Zhang, Yuan Xie, Lizhuang Ma, and Xin Tan. Doraemon: Decentralized ontology-aware reliable agent with enhanced memory oriented navigation, 2025. 6
- [9] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023. 2
- [10] Yanjia Huang, Xianshun Jiang, Xiangbo Gao, Mingyang Wu, and Zhengzhong Tu. Vistav2: World imagination for indoor vision-and-language navigation, 2025. 2
- [11] Nuri Kim, Obin Kwon, Hwiyeon Yoo, Yunho Choi, Jeongho Park, and Songhawi Oh. Topological Semantic Graph Memory for Image Goal Navigation. In *CoRL*, 2022. 3
- [12] Ue-Hwan Kim, Jin-Man Park, Taek-jin Song, and Jong-Hwan Kim. 3-d scene graph: A sparse and semantic representation of physical environments for intelligent agents. *IEEE Transactions on Cybernetics*, 50(12): 4921–4933, 2020. 3
- [13] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning*. JMLR.org, 2023. 2
- [14] Shiyao Li, Ziyang Meng, Jiansong Pei, Jiahao Chen, Bingcheng Dong, Guangsheng Li, Shenglan Liu, and Feilong Wang. Knowledge-driven visual target navigation: Dual graph navigation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7547–7554, 2025. 3
- [15] Bingqian Lin, Yunshuang Nie, Ziming Wei, Jiaqi Chen, Shikui Ma, Jianhua Han, Hang Xu, Xiaojun Chang, and Xiaodan Liang. Navcot: Boosting llm-based vision-and-language navigation via learning disentangled reasoning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 2
- [16] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 2
- [17] Peiran Liu, Qiang Zhang, Daojie Peng, Lingfeng Zhang, Yihao Qin, Hang Zhou, Jun Ma, Renjing Xu, and Yiding Ji. Toponav: Topological graphs as a key enabler for advanced object navigation, 2025. 3
- [18] Yuxing Long, Wenzhe Cai, Hongcheng Wang, Guanqi Zhan, and Hao Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment, 2024. 2, 6
- [19] Arjun Majumdar, Gunjan Aggarwal, Bhavika Devnani, Judy Hoffman, and Dhruv Batra. Zson: Zero-shot object-goal navigation using multimodal goal embeddings. In *Neural Information Processing Systems (NeurIPS)*, 2022. 1
- [20] Dujun Nie, Xianda Guo, Yiqun Duan, Ruijun Zhang, and Long Chen. Wmnav: Integrating vision-language models into world models for object goal navigation. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2392–2399, 2025. 2, 5
- [21] OpenAI. Introducing GPT-5.2. <https://openai.com/index/introducing-gpt-5-2/>, 2025. 2
- [22] Qwen Team. Qwen3.5: Towards native multimodal agents, 2026. 2
- [23] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 1
- [24] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 2, 4, 6
- [25] Zihan Wang, Yaohui Zhu, Gim Hee Lee, and Yachun Fan. NavRAG: Generating user demand instructions for embodied navigation through retrieval-augmented LLM. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 8430–8440, Vienna, Austria, 2025. Association for Computational Linguistics. 3
- [26] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. 535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592

- 593 Dd-ppo: Learning near-perfect pointgoal navigators from 2.5
594 billion frames, 2020. 1
- 595 [27] Pengying Wu, Yao Mu, Bingxian Wu, Yi Hou, Ji Ma, Shang-
596 hang Zhang, and Chang Liu. Voronav: Voronoi-based zero-
597 shot object navigation with large language model. *arXiv*
598 *preprint arXiv:2401.02695*, 2024. 5
- 599 [28] Fei Xia, Amir R. Zamir, Zhi-Yang He, Alexander Sax, Jiten-
600 dra Malik, and Silvio Savarese. Gibson Env: real-world per-
601 ception for embodied agents. In *Computer Vision and Pat-*
602 *tern Recognition (CVPR), 2018 IEEE Conference on*. IEEE,
603 2018. 2
- 604 [29] Nuo Xu, Wen Wang, Rong Yang, Mengjie Qin, Zheyuan
605 Lin, Wei Song, Chunlong Zhang, Jason Gu, and Chao Li.
606 Aligning knowledge graph with visual perception for object-
607 goal navigation. In *2024 IEEE International Conference on*
608 *Robotics and Automation (ICRA)*, pages 5214–5220. IEEE,
609 2024. 3
- 610 [30] Karmesh Yadav, Ram Ramrakhya, Santhosh Kumar Ramakr-
611 ishnan, Theo Gervet, John Turner, Aaron Gokaslan, Noah
612 Maestre, Angel Xuan Chang, Dhruv Batra, Manolis Savva,
613 Alexander William Clegg, and Devendra Singh Chaplot.
614 Habitat-matterport 3d semantics dataset. In *2023 IEEE/CVF*
615 *Conference on Computer Vision and Pattern Recognition*
616 *(CVPR)*, pages 4927–4936, 2023. 2, 6
- 617 [31] Hang Yin, Xiuwei Xu, Zhenyu Wu, Jie Zhou, and Ji-
618 wen Lu. Sg-nav: Online 3d scene graph prompting
619 for llm-based zero-shot object navigation. *arXiv preprint*
620 *arXiv:2410.08189*, 2024. 3, 5, 6
- 621 [32] Hang Yin, Xiuwei Xu, Linqing Zhao, Ziwei Wang, Jie Zhou,
622 and Jiwen Lu. Unigoal: Towards universal zero-shot goal-
623 oriented navigation. *arXiv preprint arXiv:2503.10630*, 2025.
624 3
- 625 [33] Naoki Yokoyama, Sehoon Ha, Dhruv Batra, Jiuguang Wang,
626 and Bernadette Bucher. Vlfm: Vision-language frontier
627 maps for zero-shot semantic navigation. In *International*
628 *Conference on Robotics and Automation (ICRA)*, 2024. 2,
629 5
- 630 [34] Naoki Yokoyama, Ram Ramrakhya, Abhishek Das, Dhruv
631 Batra, and Sehoon Ha. Hm3d-ovon: A dataset and
632 benchmark for open-vocabulary object goal navigation. In
633 *IEEE/RSJ International Conference on Intelligent Robots*
634 *and Systems (IROS)*, 2024. 2
- 635 [35] Bangguo Yu, Hamidreza Kasaei, and Ming Cao. L3mvn:
636 Leveraging large language models for visual target naviga-
637 tion. In *2023 IEEE/RSJ International Conference on Intel-*
638 *ligent Robots and Systems (IROS)*, page 3554–3560. IEEE,
639 2023. 2, 5, 6
- 640 [36] Sixian Zhang, Xinyao Yu, Xinhang Song, Xiaohan Wang,
641 and Shuqiang Jiang. Imagine before go: Self-supervised
642 generative map for object goal navigation. In *2024*
643 *IEEE/CVF Conference on Computer Vision and Pattern*
644 *Recognition (CVPR)*, pages 16414–16425, 2024. 6
- 645 [37] Sixian Zhang, Xinhang Song, Xinyao Yu, Yubing Bai, Xin-
646 long Guo, Weijie Li, and Shuqiang Jiang. Hoz++: Versa-
647 tile hierarchical object-to-zone graph for object navigation.
648 *IEEE Transactions on Pattern Analysis and Machine Intelli-*
649 *gence*, 47(7):5958–5975, 2025. 3, 5
- [38] Kaiwen Zhou, Kaizhi Zheng, Connor Pryor, Yilin Shen, 650
Hongxia Jin, Lise Getoor, and Xin Eric Wang. Esc: Ex- 651
ploration with soft commonsense constraints for zero-shot 652
object navigation. In *Proceedings of the 40th International* 653
Conference on Machine Learning (ICML). JMLR.org, 2023. 654
5 655