

ADANAV: ADAPTIVE REASONING WITH UNCERTAINTY FOR VISION-LANGUAGE NAVIGATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Vision-Language Navigation (VLN) requires agents to follow natural language instructions by grounding them in sequential visual observations over long horizons. Explicit reasoning could enhance temporal consistency and perception–action alignment, but reasoning at fixed steps often leads to suboptimal performance and unnecessary computation. To address this, we propose **AdaNav**, an uncertainty-based adaptive reasoning framework for VLN. At its core is the **Uncertainty-Adaptive Reasoning Block (UAR)**, a lightweight plugin that dynamically triggers reasoning. We introduce *Action Entropy* as a policy prior for UAR and progressively refine it through a *Heuristics-to-RL* training method, enabling agents to learn difficulty-aware reasoning policies under the strict data limitations of embodied tasks. Results show that with only $6K$ training samples, AdaNav achieves substantial gains over closed-source models trained on *million-scale* data, improving success rate by 20% on R2R val-unseen, 11.7% on RxR-CE, and 11.4% in real-world scenes.

1 INTRODUCTION

As a fundamental capability for embodied agents, Vision-Language Navigation (VLN) requires agents to interpret natural language instructions and continuously ground them in sequential visual observations to execute long-horizon navigation trajectories (Gu et al., 2022; Park & Kim, 2023). Existing VLM-based methods either rely on augmenting navigation with auxiliary modalities (Krantz et al., 2021; Xu et al., 2023; Yin et al., 2024), such as depth, odometry, or topological maps to strengthen spatial understanding, or instead scale up training on VLN data *without* auxiliary inputs to improve quality and generalization (Zheng et al., 2024; Wei et al., 2025; Yu et al., 2025). However, despite these advances, current methods still hindered by two major challenges of VLN: (1) Consistent temporal grounding: continuously capturing progress along the trajectory, tracking completed parts, and deciding the next action; (2) Robust perception–action mapping: grounding language in the current spatial context, recognizing landmarks, localizing itself, and selecting appropriate navigation actions.

To address these challenges, explicit reasoning has been introduced to VLN (Zhou et al., 2024b; Wang et al., 2024; Lin et al., 2025a; Chen et al., 2024a), enabling agents to better align language, perception, and action over long-horizon navigation trajectories. However, current straightforward reasoning at each step not only incurs prohibitive computational overhead, but also results in overthinking (Sui et al., 2025; Wu et al., 2025; Shen et al., 2025) that degrades navigation quality (Figure 4 and Table 6 show higher quality with fewer reasoning steps). Ideally, VLN agents should exhibit adaptive reasoning, i.e., deciding *when and how* to reason. However, achieving such adaptivity and mitigating the overconfidence issue of LLMs (Sun et al., 2025; Groot & Valdenegro-Toro, 2024; Yoo, 2024) typically require large-scale supervised fine-tuning (SFT) with task-specific data (Wen et al., 2024; Lin et al., 2025b). However, embodied interaction data is costly to collect and far from web-scale. Under such limited data conditions, it remains difficult for models to learn when and how to adaptively invoke reasoning.

To avoid the data limitation, we propose **uncertainty-based adaptive reasoning for navigation (AdaNav)**, as shown in Figure 1. By defining *Action Entropy* as an indicator for uncertainty, AdaNav utilizes this as an objective and interpretable heuristic prior to decide when and how to reason, and then refine this prior gradually through reinforcement learning (RL) to optimize the reasoning trigger

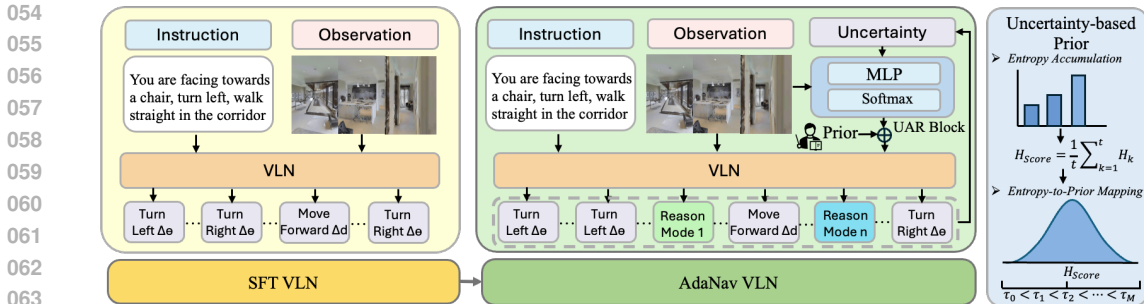


Figure 1: AdaNav augments a base VLN model by integrating the Uncertainty-Adaptive Reasoning Block (UAR Block). UAR Block leverages model uncertainty to autonomously trigger reasoning modes and timing, enhancing consistent temporal grounding and perception–action mapping while significantly improving efficiency and mitigating overthinking.

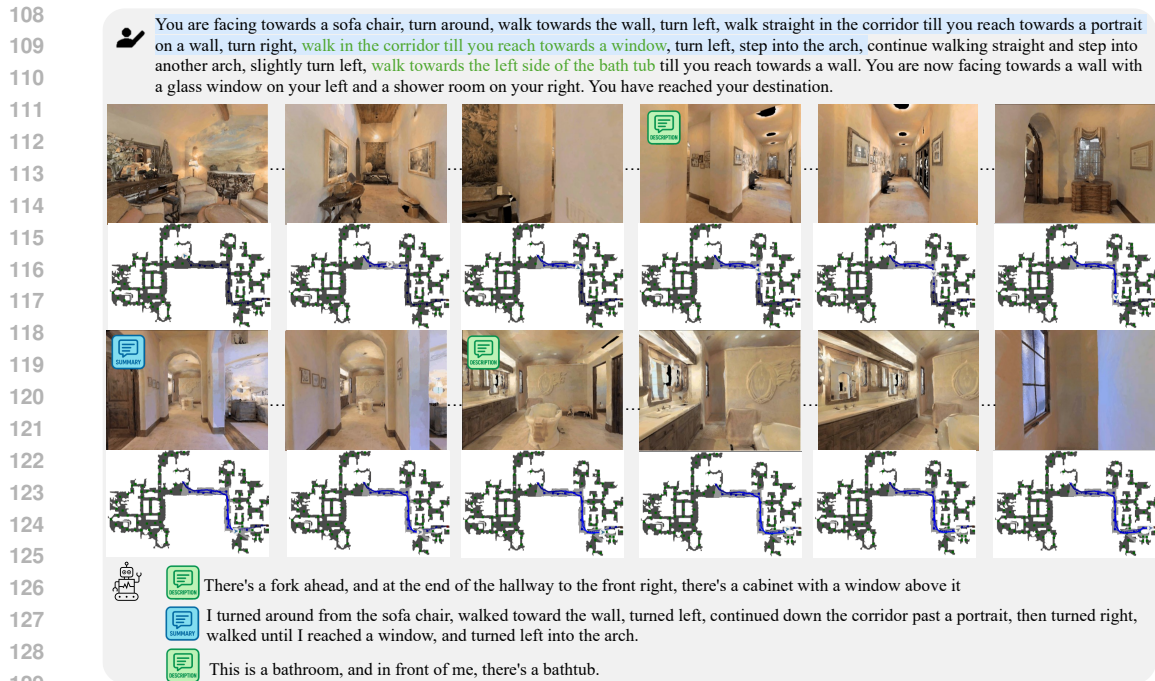
policy. By combining the efficiency of heuristic guidance with the optimality of RL, AdaNav do not involve costly labeled reasoning triggering data, but enable the agent to automatically invoke reasoning when necessary to maintain temporal grounding and robust perception–action mapping in the long-horizon navigation. See Figure 2 as an example.

To realize AdaNav, we introduce a *Uncertainty-Adaptive Reasoning Block (UAR Block)* and the *Heuristic-to-RL* training mechanism. UAR block, as a plugin for available VLN models, collects historical, embodied-interaction-dependent uncertainty signals and generates vectorized control signals to dynamically trigger VLN for appropriate reasoning modes. Leveraging the interpretable signals from the UAR Block, the Heuristic-to-RL training first explores the action space under these heuristic priors (e.g., triggering reasoning when uncertainty exceeds a threshold) to guide decision-making at critical moments. As training progresses, the influence of these priors is gradually annealed, allowing RL to autonomously refine the UAR reasoning-trigger policy for optimal reward.

To demonstrate the effectiveness of AdaNav, we integrate it with state-of-the-art open-source VLN backbones and evaluate on classic benchmarks. Remarkably, **with only 6K training samples, AdaNav significantly surpasses closed-source models trained on million-scale data.** Specifically, our method achieves an average 20% improvement in success rate on R2R val-unseen (Krantz et al., 2020), and even without training on the larger and more challenging RxR-CE (Ku et al., 2020), AdaNav yields a 11.7% gain, demonstrating the cross-dataset generality. Additionally, AdaNav exhibits strong robustness in Sim-to-Real deployment, achieving approximately a 11.4% success rate improvement over 150 instructions across four **real-world indoor scenes**. As training proceeds, AdaNav reduces the average number of reasoning steps per trajectory to only **2.5** (over trajectories with an average length of **55 steps**), while the success rate increases 7% compared to reasoning at fixed steps. Notably, 71% of reasoning steps are concentrated on hard trajectories. These results indicate that training makes reasoning more difficulty-aware and mode-adaptive.

2 RELATED WORK

VLN with Auxiliaries. VLN requires agents to follow free-form linguistic instructions and visual cues to reach a target location. Early studies relied on pre-defined waypoints for discrete navigation (Qi et al., 2020b; Thomason et al., 2020) in the Habitat-Matterport3D simulator (Chang et al., 2017), while more recent works (Qi et al., 2020a; An et al., 2021; Hong et al., 2020; Tan et al., 2019; Wang et al., 2019) use continuous environments, namely *VLN-CE*, like Habitat (Krantz et al., 2020), enabling low-level actions (e.g., move forward, rotate) for more realistic navigation. With the rise of Transformers, many works introduced pre-trained methods with auxiliary modalities, e.g., depth, odometry, or topological maps, for VLN (Ma et al., 2019; Wang et al., 2019). DUET (Chen et al., 2022) and ETPNav (An et al., 2024) build topological maps for global navigation understanding, while GridMM (Wang et al., 2023) introduced a dynamic egocentric grid memory. Although these methods improve spatial awareness, they inevitably limit generalization and introduce computational overhead and noise (Zhang et al., 2024b). Modern works increasingly target video-only



131 Figure 2: A visualization example of AdaNav. It autonomously invokes reasoning, e.g., summariza-
 132 tion and description when necessary to maintain consistent temporal grounding and robust percep-
 133 tion–action mapping.

135 general solutions for VLN without auxiliaries (Zhang et al., 2024b; Cheng et al., 2024; Zhang et al.,
 136 2024a). *VLN-CE with only videos captured by the monocular camera is also the target of this paper.*

138 **Vision-Language Models for Navigation.** With the rapid development of Vision-Language Mod-
 139 els (VLMs), RT-2 (Zitkovich et al., 2023) demonstrates the potential of transferring web-scale
 140 knowledge from VLMs to generalizable robotic manipulation. Recent work has focused on scal-
 141 ing VLN training data and fine-tuning large VLMs. For example, Navid (Zhang et al., 2024b) used
 142 550k navigation samples to fine-tune Vicuna for navigation; NaVILA (Cheng et al., 2024) expanded
 143 to 3–5M samples combining real and simulated navigation data plus general VQA supervision;
 144 Uni-NaVid (Zhang et al., 2024a) further incorporated 3.6M multi-task trajectories from Habitat-
 145 Matterport3D (Chang et al., 2017; Krantz et al., 2020) and real video QA data (Azuma et al., 2022;
 146 Chen et al., 2024b; Li et al., 2024) for cross-task generalization. Despite these advances, VLM-based
 147 VLN agents still fall short in task quality, struggling with consistent temporal grounding and robust
 148 perception–action mapping, particularly in long-horizon trajectories and complex environments.

150 **Explicit Reasoning for Navigation.** To mitigate these challenges, recent works introduce ex-
 151 plicit reasoning via off-the-shelf LLMs, where pre-defined programming rules constrain when and
 152 how reasoning modes—description, summarization, or error correction—are applied. For exam-
 153 ple, LLM-Planner (Song et al., 2023) parses instructions into sub-goals; NavGPT (Zhou et al.,
 154 2024b) generates step-wise textual scene descriptions and historical trajectories; NavGPT-2 (Zhou
 155 et al., 2024a) further integrates visual grounding; MiC (Qiao et al., 2023) organizes reasoning into a
 156 “summarization–planning–correction” loop; DiscussNav (Long et al., 2024b), MCGPT (Zhan et al.,
 157 2024), and InstructNav (Long et al., 2024a) leverage expert collaboration or memory graphs for
 158 error correction and historical summarization.

159 While these rule-driven frameworks offer interpretability, they inherently restrict flexibility in open-
 160 ended environments, hinder efficiency, and may lead to overthinking (Fang et al., 2025; Dai et al.,
 161 2025). In contrast, our method will employ a learnable mechanism that enables agents to au-
 tonomously decide when and how to reason.

3 METHOD

3.1 PROBLEM FORMULATION OF ADANAV

The central problem investigated in this work is how to enable an embodied agent to adaptively decide *when* and *how* to invoke reasoning during VLN. Unlike conventional approaches that either disable reasoning or enforce rule-based reasoning at fixed steps, our goal is to learn an autonomous reasoning policy that dynamically determines the timing and mode of reasoning, optimizing both efficiency and navigation performance.

Vision-Language Navigation. We consider a standard VLN setting where an agent is placed in a 3D environment \mathcal{E} with state space \mathcal{S} and action space $\mathcal{A} = \{\text{turn_left}(\Delta\theta), \text{turn_right}(\Delta\theta), \text{move_forward}(\Delta d), \text{stop}\}$, where $\Delta\theta$ and Δd denote the angle and distance, respectively. Given a natural language instruction I and sequential visual observations $\{o_1, o_2, \dots\}$, the agent executes a trajectory $\tau = \{(s_t, a_t)\}_{t=1}^H$ toward a goal s^* specified implicitly by I , aiming to maximize task success:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} [\mathbf{1}(s_H = s^*)]. \quad (1)$$

Adaptive Reasoning Navigation. To improve VLN performance in long-horizon and complex environments, we allow explicit reasoning at step t with a *mode* variable $m_t \in \{\emptyset\} \cup \mathcal{M}$ and reasoning content r_t . Here, $m_t = \emptyset$ denotes *no reasoning* (so $r_t = \emptyset$), while $m_t \in \mathcal{M}$ denotes invoking a reasoning mode from a predefined set. In this work, we consider three reasoning modes: *description*, *summary*, and *error correction* (see Figure 6 and Appendix D for instances). The agent’s policy then consists of two coupled processes: 1) a navigation policy $\pi_{\text{nav}}(a_t | h_t^{\text{nav}}, I, r_{\leq t})$, and 2) a reasoning policy $\pi_{\text{rea}}(m_t, r_t | h_t^{\text{rea}}, I)$ that jointly decides *when* to reason (via $m_t = \emptyset$ vs. $m_t \neq \emptyset$) and *which mode* to use (via $m_t \in \mathcal{M}$).

The overall joint policy is

$$\pi^*(a_t, m_t, r_t | h_t, I) = \pi_{\text{nav}}(a_t | h_t^{\text{nav}}, I, r_{\leq t}) \cdot \pi_{\text{rea}}(m_t, r_t | h_t^{\text{rea}}, I) \quad (2)$$

where $h_t = (h_t^{\text{nav}}, h_t^{\text{rea}})$ denotes the full history, with h_t^{nav} and h_t^{rea} representing the navigation-related and reasoning-related information, respectively. For clarity, we factorize the reasoning policy as:

$$\pi_{\text{rea}}(m_t, r_t | h_t^{\text{rea}}, I) = \underbrace{\pi_{\text{txt}}(r_t | m_t, h_t^{\text{rea}}, I)}_{\text{reasoning content}} \cdot \underbrace{\pi_{\text{sel}}(m_t | h_t^{\text{rea}}, I)}_{\text{when/which mode}} \quad (3)$$

With the constraint $r_t = \emptyset$ if $m_t = \emptyset$. Here, π_{txt} shares the same network as π_{nav} .

By integrating navigation and reasoning, the overall learning objective is to jointly optimize both policies, aiming to maximize task performance while maintaining computational efficiency.

$$\pi^* = \arg \max_{(\pi_{\text{rea}}, \pi_{\text{nav}})} \mathbb{E}_{\tau \sim (\pi_{\text{nav}}, \pi_{\text{rea}})} [R_{\text{task}}(\tau)] \quad (4)$$

where $R_{\text{task}}(\tau)$ jointly accounts for navigation success (e.g., progress or success indicator) and the latency penalty induced by reasoning calls.

3.2 METHODOLOGY OF ADANAV

Motivation. Adaptive reasoning requires the agent to selectively decide *when* reasoning is beneficial and *which mode* to invoke. However, native VLMs are neither sensitive nor objective in perceiving task difficulty, often resulting in overconfidence. In LLM research, similar issues (e.g., in mathematical reasoning) have been alleviated by incorporating high-quality reasoning traces with supervised fine-tuning (Zhang et al., 2025; Tian et al., 2025; Guo et al., 2025). In contrast, for embodied agents, collecting such high-quality interaction traces is prohibitively expensive. This motivates the development of alternative approaches that enable agents to acquire adaptive reasoning capabilities without relying on large-scale reasoning supervision.

To this end, we propose Adaptive Reasoning Navigation (AdaNav), which leverages interpretable uncertainty signals to dynamically trigger reasoning only when necessary. By combining a learnable reasoning policy with a navigation policy, AdaNav enables efficient, difficulty-aware, and mode-adaptive reasoning, achieving high performance in long-horizon and complex VLN tasks.

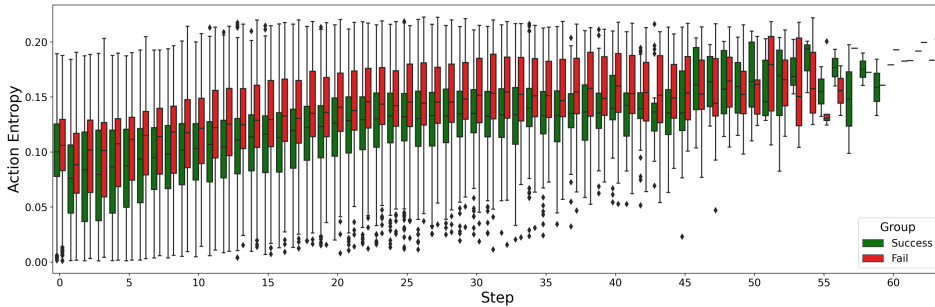


Figure 3: Action entropy per step for successful (green) vs. failed (red) trajectories. Failed trajectories show higher, especially in later steps, indicating that policy uncertainty correlates with navigation errors.

3.2.1 UNCERTAINTY-ADAPTIVE REASONING BLOCK

Recent works (Kazemnejad et al., 2024; Wang et al., 2025a; Fu et al., 2025) in language reasoning have shown that high-entropy tokens exert a disproportionately large influence on single-step text generation. Inspired by this, we explore whether a similar principle can serve as a signal for identifying “forking steps” in navigation. Specifically, we define action entropy H as an uncertainty measure:

$$H = -\frac{1}{N} \sum_{t=1}^N \sum_{v=1}^V p_t[v] \log p_t[v], \quad (5)$$

where N is the number of tokens generated, V is the size of the vocabulary, and $p_t[v]$ is the probability of the v -th vocabulary token at time step t .

To validate the effectiveness of action entropy, we conduct a diagnostic study on navigation trajectories. As shown in Figure 3, episodes with high and sustained entropy are strongly correlated with failures, while successful trajectories maintain consistently low entropy (Different Means). Furthermore, instantaneous entropy alone is insufficient: short-lived spikes do not necessarily imply failure, and many successful trajectories exhibit temporary fluctuations without requiring reasoning (Comparable Extremes).

Conversely, *combining historical action entropy trends with current action entropy states provides a more reliable signal $H_{\leq t}$* : successful trajectories show relatively lower entropy over time, while failure-prone ones accumulate persistently high entropy.

Method Design. Motivated by these findings, we design a lightweight Uncertainty-Adaptive Reasoning Block that fuses $H_{\leq t}$ with the current observation o_t , forming the reason-related information $h_t^{\text{rea}} = (H_{\leq t}, o_t)$. These signals are combined into a compact control vector:

$$p_{\text{mode}}^t = W_1 H_{\leq t} + W_2 o_t + W_3 I + b, \quad (6)$$

which directly parametrizes the reasoning mode logits. From this, the mode selection policy (cf. Equation 3) is given by:

$$\pi_{\text{sel}}(m_t | h_t^{\text{rea}}, I) = \text{Softmax}(p_{\text{mode}}^t). \quad (7)$$

3.2.2 HEURISTIC-TO-RL TRAINING

Benefiting from the interpretable signals of the UAR Block, we do not require large-scale reasoning annotations. Instead, we propose a Heuristic-to-RL training mechanism, shown in Algorithm 1, that bootstraps policy learning with uncertainty-based heuristics. These priors provide a stable cold-start exploration, enabling the agent to discover useful reasoning patterns without exhaustive supervision. As training progresses, the heuristic influence is gradually annealed, allowing reinforcement learning to refine the reasoning trigger policy. This approach integrates the efficiency of heuristic guidance with the optimality of RL, leading to adaptive long-horizon reasoning strategies that generalize to novel environments.

Uncertainty-based Prior. In the cold-start phase, the RL policy has not yet learned meaningful mode selection. We therefore initialize training with an uncertainty-based prior. Intuitively, a higher

entropy indicates a higher uncertainty, which requires stronger reasoning. We compute the scalar entropy score as the mean of past entropies, $H_{\text{score}} = \frac{1}{t} \sum_{k=1}^t H_k$, and map it into a soft prior distribution over $|\mathcal{M}| + 1$ reasoning modes (including the “no reasoning” option):

$$p_{\text{prior}} = \frac{\exp(-|H_{\text{score}} - \tau_m|/\sigma)}{\sum_{i=0}^{|\mathcal{M}|} \exp(-|H_{\text{score}} - \tau_i|/\sigma)}, \quad m = 0, \dots, |\mathcal{M}| \quad (8)$$

where $\{\tau_0, \tau_1, \dots, \tau_{|\mathcal{M}|}\}$ are mode-specific entropy thresholds ($\tau_k = \tau_0 + k\delta$), and σ controls the smoothness of the prior.

Heuristic-to-RL Transition. To gradually shift control from heuristic priors to learned RL policies, we fuse the prior distribution with the model prediction as:

$$p_{\text{final}}^t = \lambda_t \cdot p_{\text{prior}} + (1 - \lambda_t) \cdot p_{\text{model}}, \quad (9)$$

where λ_t is annealed from 1 to 0 over training steps, allowing the RL policy p_{model} to progressively take over from the uncertainty-based heuristic prior p_{prior} . Accordingly, Equation 7 can be expressed as:

$$\pi_{\text{sel}}(m_t | h_t^{\text{rea}}, I) = \text{Softmax}(p_{\text{final}}^t). \quad (10)$$

Reward Design. We first define the *reasoning cost* as a normalized penalty based on the relative reasoning length:

$$c_{\text{rea}}(t) = \text{clip}\left(\frac{L_t - L_{\text{shortest_success}}}{L_{\text{window}}}, 0, 1\right) \quad (11)$$

where L_i is the reasoning length for the current step, $L_{\text{shortest_success}}$ is the minimal generation length among success samples within the group, and L_{window} is a constant penalty window.

For the navigation objective, we adopt the common extrinsic reward based on distance reduction, where the immediate reward is defined as $r(s_t, a_t) = D_{\text{target}}(s_t) - D_{\text{target}}(s_{t+1})$; $t < T$, with $D_{\text{target}}(s_t)$ denoting the geodesic distance from the current state s_t to the target location s_{target} .

Finally, by combining extrinsic reward and reasoning cost, the overall task reward formulated in 4 is defined as the discounted cumulative return:

$$R_{\text{task}}(\tau) = \sum_{t=1}^T \beta^{t-1} (r(s_t, a_t) - c_{\text{rea}}(t)) \quad (12)$$

where $\beta \in (0, 1]$ is the discount factor controlling the weight of future rewards. This formulation encourages the agent to navigate efficiently toward the goal while avoiding unnecessary reasoning overhead.

Overall, this Heuristic-to-RL scheme combines the efficiency of uncertainty-based priors with the optimality of RL, allowing the agent to progressively acquire adaptive reasoning strategies.

4 EXPERIMENTS

We conduct experiments to answer the following questions: (1) **Performance Gain:** How much does our proposed AdaNav improve over existing models on VLN-CE benchmarks and general spatial scene understanding tasks? (2) **Reasoning Coordination:** What scheduling strategy has the UAR Block learned, and does it affect navigation efficiency? (3) **Real-World Effectiveness:** How effective is AdaNav when deployed in real-world scenarios?

4.1 PERFORMANCE GAIN

Implementation details. 1. Base models. AdaNav is designed to be general and can be integrated into existing VLN models with minimal modifications. To demonstrate its strong generalization ability, we adopt two SOTA open-source VLN models, NAVID (Zhang et al., 2024b) and

Algorithm 1 Heuristic-to-RL

- 1: Initialize navigation policy π_{nav} , reasoning selector π_{sel} , annealing schedule λ_t
- 2: **for** each training episode τ **do**
- 3: **for** each step $t = 1 \dots T$ **do**
- 4: Observe state o_t and entropy $H_{\leq t}$
- 5: Compute control vector p_{model}
- 6: Estimate Uncertainty prior p_{prior}
- 7: Fuse prior and model (Eq. 9)
- 8: Sample reasoning mode $m_t \sim p_{\text{final}}$
- 9: **if** $m_t \neq \emptyset$ **then**
- 10: Generate reasoning r_t
- 11: **end if**
- 12: Select action $a_t \sim \pi_{\text{nav}}(a_t | o_t, r_{\leq t})$
- 13: Execute a_t , observe next state s_{t+1}
- 14: Compute extrinsic reward $r(s_t, a_t)$ and reasoning cost $c_{\text{rea}}(t)$
- 15: **end for**
- 16: Compute task reward (Eq. 12)
- 17: Update policy (Eq. 4)
- 18: **end for**

NAVILA (Cheng et al., 2024), as our base models. **2. Training setup.** Training is conducted on 4 NVIDIA RTX A100 GPUs. We construct the training set by randomly sampling 3,000 episodes from the training splits of both R2R (Krantz et al., 2020) and RxR (Ku et al., 2020). For rollout collection during training, each episode is rolled out 5 turns, and the learning rate is set to 1×10^{-6} . **3. Benchmarks.** To assess both navigation and spatial scene understanding, we evaluate on the val-unseen splits of R2R and RxR for navigation, and on the ScanQA validation set for scene understanding. Detailed settings are provided in Appendix B.

VLN-CE Benchmarks. We compare AdaNav with **recent million-scale closed-source models**, including NAVID-4D (Liu et al., 2025), UNI-NAVID (Zhang et al., 2024a), and MON-ODREAM (Wang et al., 2025b). As shown in Table 1, although closed-source models generally outperform open-source ones, AdaNav achieves substantial gains with only 6K training episodes, improving NAVID and NAVILA by an average of 20% on R2R and 14.6% on RxR, respectively, and surpassing all closed-source baselines.

Cross-dataset Evaluation. As shown in Table 2, we test cross-dataset generalization by training AdaNav solely on 3K R2R samples and evaluating zero-shot on RxR Val-Unseen. AdaNav substantially improves base models, surpassing closed-source systems and demonstrating strong transferability.

Table 1: Comparison with the state-of-the-art method on Val-Unseen split of R2R-CE and RxR-CE.

Method	Observation				R2R-CE Val-Unseen				RxR-CE Val-Unseen				Training Data
	S.RGB	Pano.	Depth	Odo.	NE↓	OS↑	SR↑	SPL↑	NE↓	OS↑	SR↑	nDTW↑	
Open-Source													
Seq2Seq	✓		✓		7.77	37.0	25.0	22.0	12.10	13.9	11.9	30.8	-
CMA	✓		✓		7.37	40.0	32.0	30.0	-	-	-	-	-
RGB-Seq2Seq	✓				10.10	8.0	0.0	0.0	-	-	-	-	-
RGB-CMA	✓				9.55	10.0	5.0	4.0	-	-	-	-	-
LAW	✓		✓	✓	6.83	44.0	35.0	31.0	10.90	8.0	8.0	38.0	150K
AO-Planner		✓	✓		5.55	59.0	47.0	33.0	7.06	43.3	30.5	50.1	40K (Distill)
NaVid	✓				5.47	49.0	37.0	35.0	6.79	46.2	40.5	52.2	550K
NaVILA	✓				5.22	62.5	54.0	49.0	6.77	49.3	44.0	58.8	~3000K
Close-Source													
NaVid-4D	✓	✓			5.99	55.7	43.8	37.1	-	-	-	-	1840K
Uni-NaVid	✓				5.58	53.5	47.0	42.7	6.24	48.7	40.9	-	3600K
MonoDream	✓				5.45	61.5	55.8	49.1	6.38	55.8	49.4	-	1420K
AdaNav													
NaVid-AdaNav	✓				5.39	57.89	47.7	42.34	6.38	58.1	47.01	56.8	+6K
NaVILA-AdaNav	✓				5.01	66.62	60.19	50.0	6.21	60.51	49.8	62.2	+6K

Table 2: Cross-dataset performance on the RxR-CE [30] ValUnseen split, without training on RxR-CE training set. * indicates our reproduction following the original papers.

Method	RxR-CE Val-Unseen			
	NE↓	OS↑	SR↑	SPL↑
Open-Source				
LAW	10.87	21.0	8.0	8.0
CM2	8.98	25.3	14.4	9.2
Seq2Seq	11.8	5.02	3.51	3.43
CMA	11.7	10.7	4.41	2.47
NaVid*	8.57	32.21	21.3	20.01
NaVILA*	8.96	43.35	32.5	26.82
Close-Source				
Uni-NaVid	8.08	40.9	29.5	28.1
MonoDream	8.57	35.9	25.1	21.6
AdaNav				
NaVid-AdaNav	8.21	39.21	28.95	27.73
NaVILA-AdaNav	8.25	48.65	38.82	31.21

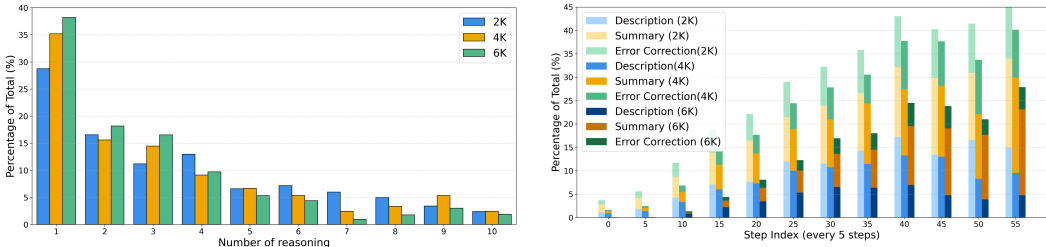
Table 3: Evaluation of spatial scene understanding performance on the ScanQA Validation split. * and † denote the use of 8 frames and 64 frames, respectively.

Method	ScanQA Validation				
	Bleu-4↑	Rouge↑	Cider↑	Meteor↑	EM↑
3D Large Multi-modal Models					
3D-LLM	7.2	32.3	59.2	12.2	20.4
LL3DA	13.5	37.3	76.8	15.9	-
Chat-3Dv2	14.0	-	87.6	-	-
Scene-LLM	12.0	40.0	80.0	16.6	27.2
LEO	13.2	49.2	101.4	20.0	24.5
2D Vision-Language-ActionModels					
Uni-NaVid	-	45.74	94.72	19.24	28.01
NavILLM	12.0	38.4	75.9	15.4	23.0
NaVILA*	14.8	46.4	95.1	18.7	27.0
NaVILA†	16.9	49.3	102.7	20.1	28.6
NaVILA-AdaNav*	15.33	47.75	97.34	19.12	27.27
NaVILA-AdaNav†	16.65	50.6	102.81	21.25	29.57

Spatial Scene Understanding Benchmarks. As a general navigation agent, robust spatial scene understanding (e.g., object localization, referring, and spatial reasoning) is crucial. To verify whether AdaNav fine-tuning affects such capability, we evaluate on the ScanQA validation benchmark (Azuma et al., 2022), a widely used dataset for 3D question answering, as shown in Table 3.

Table 4: Comparing in three diverse real-world environments scenes.

Method	Meeting Room				Home				Office			
	Simple		Complex		Simple		Complex		Simple		Complex	
	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑
Navid	2.0	67.5	2.8	50.2	1.55	65.5	1.88	55.5	2.5	61.0	3.0	52.5
Navid-AdaNav	1.6	78.5	2.2	60.5	1.3	78	1.5	75.5	2	70	2.5	66.5
Navila	1.8	74	2.0	65	1.0	82.5	1.4	82.5	2.1	76.8	2.2	70
Navila-AdaNav	1.0	82.5	1.6	73.5	0.85	95	1.1	88	1.5	87.6	2.1	75.5



(a) Distribution of reasoning steps per trajectory. As training data increases, the frequency of reasoning gradually decreases, with the model learning to invoke reasoning at more critical moments, thereby balancing efficiency and effectiveness.

(b) Number of reasoning invocations at each step, broken down by reasoning modes, across different training scales. As data increases, reasoning becomes more concentrated on high-uncertainty steps, with stronger preference for *summary* mode at later stages.

Figure 4: Analysis of reasoning behaviors in AdaNav across training scales.

Results show that after Heuristic-to-RL training, AdaNav not only preserves its general scene understanding ability without using ScanQA training data, but also achieves slight improvements, indicating enhanced robustness and transferability.

Real-World Evaluation To demonstrate the effectiveness of AdaNav in real-time settings, we conduct experiments in real-world environments using 25 simple or complex instructions. Each instruction requires the agent to complete 5–10 sequential landmark-following sub-tasks (e.g., “After passing through the ticket gate, walk straight to the sofa, turn right, take the elevator, continue walking straight, pass through the door until reaching a supermarket, and finally stop at the counter”). Each instruction is executed three times across three types of environments: *Meeting Room*, *Home*, *Office*, and *Outdoor*, following the protocol in prior works (Cheng et al., 2024; Zhang et al., 2024b). The results are summarized in Table 4 and Table 15 in Appendix G.

4.2 ANALYSIS OF UAR BLOCK

To better understand how the UAR Block adapts over training, we conduct a systematic analysis using models trained with 2K, 4K, and 6K data. We focus on two aspects: (1) the frequency and distribution of reasoning invocations across different steps and reasoning modes, and (2) the tendency to trigger reasoning under different episode difficulty levels.

Frequency of Reasoning Figure 4a shows the distribution of reasoning steps under different scales of training data, while Table 6 reports the corresponding performance. As the training data increases, the model tends to reduce the frequency of reasoning, focusing more on triggering reasoning at critical moments, thereby balancing efficiency and effectiveness.

Step-wise Reasoning Statistics. Figure 4b shows the number of reasoning invocations at each navigation step, broken down by mode (*description*, *summary*, *error correction*) for the three training scales. We observe that as training data increases, the agent learns to concentrate reasoning on critical steps where uncertainty is higher, while reducing redundant reasoning in low-uncertainty steps. Additionally, the model increasingly favors *summary* and *error correction* modes at later steps, indicating adaptive mode selection based on task context.

Table 5: Proportion of reasoning triggers in hard episodes (success + failure = 100%, excl. Step 0). Agents tend to invoke more reasoning in harder episodes.

Step	0	5	10	15	20	25	30	35	40	45	50	55
failure	0	100.0	85.0	71.43	74.51	70.43	69.81	77.33	80.30	72.55	75.0	75.68

Table 6: Performance on R2R.

Data	2K	4K	6K
SR	44.8	46.5	47.7

Difficulty-conditioned Reasoning. To disentangle how reasoning adapts to task difficulty, we first categorize each episode by its outcome (success vs. failure) under a baseline agent without reasoning coordination. We treat successful episodes as relatively easy and failed episodes as harder ones. We then re-run the model with the coordination layer enabled and analyze reasoning triggers across these two difficulty groups.

As shown in Table 5, for hard episodes that the base model fails to solve, reasoning is triggered significantly more frequently. This indicates that the UAR Block adaptively allocates reasoning capacity, focusing on challenging scenarios rather than applying reasoning uniformly across all episodes.

Conclusions. These analyses demonstrate that the UAR Block effectively learns both *when* and *which mode* to reason. As training progresses, reasoning becomes more temporally focused, mode-adaptive, and difficulty-aware, enabling the agent to improve navigation performance while minimizing redundant reasoning overhead.

5 ABLATION

To examine the robustness of AdaNav and assess whether its performance is overly dependent on specific hyperparameter choices, we conduct a series of ablation studies. Our analyses focus on three aspects: (1) component ablation, (2) sensitivity to key hyperparameters. More detailed ablation results are provided in Appendix C.

Component Ablation. We use Navid as the base model and remove or replace major components to isolate their contributions. **(i) w/o UAR Block:** reasoning is invoked at a fixed step (5 step) interval or randomly, without adaptive control. **(ii) w/o Heuristic Prior:** the agent relies purely on reinforcement learning from scratch without uncertainty-based heuristic. **(iii) w/o RL Fine-tuning:** reasoning triggers are guided only by heuristic signals without further policy refinement. Results show that removing either coordination or RL fine-tuning leads to significant performance degradation, confirming that both adaptive gating and learned refinement are essential.

Table 7: Ablation on hyperparameter sensitivity and component effectiveness on R2R-CE Val-Unseen. Here, * denotes fixed-interval (5 steps) triggering, and † denotes random triggering.

		τ_0	δ	NE↓	OS↑	SR↑	SPL↑										
Component	Navid	80%	0.1	5.43	57.72	48.82	43.56	σ	NE↓	OS↑	SR↑	SPL↑					
			0.2	5.42	57.92	49.11	43.53										
			0.3	5.42	58.01	49.05	43.57										
	w/o UAR*	80%	0.1	5.40	58.75	49.61	43.87						0.05	5.43	58.85	48.85	43.55
			0.2	5.39	57.89	47.7	42.34						0.10	5.40	58.55	49.13	43.62
			0.3	5.39	58.81	49.53	43.85						0.15	5.39	57.89	47.7	42.34
	w/o UAR†	85%	0.1	5.40	58.75	49.61	43.87						0.20	5.41	58.73	49.55	44.02
			0.2	5.39	57.89	47.7	42.34						0.25	5.44	58.72	49.25	43.88
			0.3	5.39	58.81	49.53	43.85						0.30	5.48	58.13	48.72	43.92
	w/o RL	90%	0.1	5.47	57.98	48.95	43.42						(c) Effect of σ .				
			0.2	5.48	57.80	48.83	43.34										
			0.3	5.43	57.78	48.85	43.35										
w/o HP	90%	0.1	5.47	57.98	48.95	43.42											
		0.2	5.48	57.80	48.83	43.34											
		0.3	5.43	57.78	48.85	43.35											
AdaNav				5.39	57.89	47.7	42.34										

(a) Component ablation.

(b) Effect of (τ_0, δ) .(c) Effect of σ .

Hyperparameter Sensitivity. The key hyperparameters in our framework lie in the *Heuristic-to-RL* stage, where we introduce mode-specific entropy thresholds: (τ_0, δ) that govern reasoning triggers prior, and a scaling factor σ .

As shown in Table 7a, a well reasoning prior significantly facilitates training. Specifically, τ_0 is estimated from the base model by analyzing 1,000 validation episodes and selecting a percentile of the step-wise action entropy extrema. We experiment with percentiles at 80%, 85%, and 90%, which define progressively stricter confidence thresholds. On top of this, δ incrementally shifts the thresholds for higher reasoning modes, thereby shaping the curriculum schedule. The corresponding results are summarized in Table 7b and Table 7c.

Computational Efficiency of the UAR Block The UAR Block adds negligible latency while substantially improving navigation performance. To quantify its computational cost, we measure the latency of each component as summarized in Table 8.

Table 8: Latency of UAR Block and Navigation Policy Stages.

Component	Peak Latency	Average Latency	P90 Latency
UAR Block	0.0039 s	0.0029 s	0.0037 s
Base Action Policy	1.31 s	0.76 s	1.00 s
Text Reasoning	5.01 s	2.94 s	3.81 s

P90 Latency represents the time value below which 90% of the sample measurements fall.

As shown in Table 8, the UAR Block introduces **negligible additional computation**, accounting for less than 0.4% of the backbone inference time. Although a single text reasoning step is approximately $3.9\times$ slower than the base policy inference, the adaptive reasoning mechanism learned by the UAR Block reduces the number of reasoning steps per trajectory to only 2.5, for trajectories with an average length of 55 steps.

Effectiveness of the UAR Block Design As discussed in Section 3.2.1, *action entropy* serves as a reliable statistical indicator of model uncertainty, offering an interpretable and stable heuristic prior for reinforcement learning. At the same time, contextual features extracted from visual observations and language instructions remain essential, as they allow the policy to ground uncertainty in both the physical environment and the intended task objective.

Through Heuristic-to-RL Training, AdaNav progressively aligns three key elements—(i) environmental interaction dynamics, (ii) internal uncertainty signals captured by action entropy, and (iii) task intent from textual instructions—ultimately enabling the agent to learn a coherent and difficulty-aware reasoning policy.

Table 9: Ablation study on the contributions of action entropy and observation features.

Component	NE↓	OS↑	SR↑	SPL↑	Avg.Infer.Num
Navid	5.47	49.0	37.0	35.0	0
Only Action Entropy	5.41	56.73	45.23	41.55	4.5
Only Observation	5.43	55.23	43.12	40.74	8.7
Navid-AdaNav	5.39	57.89	47.7	42.34	2.5

As shown in Table 9, when trained with the same 6K samples, models that rely solely on action entropy or solely on observation features exhibit a $2\text{--}4\times$ increase in reasoning frequency, while simultaneously suffering a $5\text{--}10\%$ degradation in navigation success rate. These results indicate that combining contextual features with uncertainty signals is essential for achieving both efficient reasoning and high-performance navigation.

6 CONCLUSION

In this work, we tackled the long-standing challenges of consistent temporal grounding and robust perception–action mapping in Vision-Language Navigation. We proposed AdaNav, an uncertainty-based adaptive reasoning framework that integrates the UAR Block with a Heuristic-to-RL training mechanism. This design enables agents to invoke reasoning adaptively, guided first by interpretable heuristic priors and then refined through reinforcement learning, without relying on costly labeled reasoning data. Extensive experiments show that AdaNav delivers substantial improvements: surpassing million-scale closed-source models with only 6K samples, generalizing effectively across R2R and RxR, and demonstrating strong robustness in real-world deployment. Moreover, AdaNav reduces reasoning frequency while making it more difficulty-aware and mode-adaptive, striking a balance between efficiency and effectiveness. AdaNav provides a principled and practical step toward scalable, adaptive reasoning in embodied agents.

7 ETHICS STATEMENT

This work does not involve human subjects, personally identifiable information, or sensitive user data. All datasets used in our experiments are publicly available benchmark datasets (e.g., R2R and related VLN benchmarks) that have been widely adopted in the community. We carefully followed the respective dataset licenses and usage guidelines. Our proposed method aims to improve efficiency and robustness in embodied AI navigation, without introducing risks of harmful applications or discriminatory practices. To the best of our knowledge, this work does not pose ethical concerns regarding fairness, safety, or privacy.

8 REPRODUCIBILITY STATEMENT

We have made substantial efforts to ensure the reproducibility of our results. The paper provides detailed descriptions of the model architecture, training objectives, and evaluation protocols (see Sections 4.1). Hyperparameter settings are specified in the Section 5. All datasets used are publicly available. To further facilitate reproducibility, we will release the full source code and trained checkpoints in the near future, ensuring that the community can replicate and extend our findings.

REFERENCES

- Nvidia jetson nano developer kit. <https://developer.nvidia.com/embedded/jetson-nano-developer-kit>. Accessed: 2025-09-15.
- Dong An, Yuankai Qi, Yan Huang, Qi Wu, Liang Wang, and Tieniu Tan. Neighbor-view enhanced model for vision and language navigation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pp. 5101–5109, 2021.
- Dong An, Hanqing Wang, Wenguan Wang, Zun Wang, Yan Huang, Keji He, and Liang Wang. Etpnav: Evolving topological planning for vision-language navigation in continuous environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 19129–19139, 2022.
- 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. *arXiv preprint arXiv:1709.06158*, 2017.
- Jiaqi Chen, Bingqian Lin, Ran Xu, Zhenhua Chai, Xiaodan Liang, and Kwan-Yee K Wong. Mapgpt: Map-guided prompting with adaptive path planning for vision-and-language navigation. *arXiv preprint arXiv:2401.07314*, 2024a.
- Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. Think global, act local: Dual-scale graph transformer for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16537–16547, 2022.
- Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13320–13331, 2024b.
- An-Chieh Cheng, Yandong Ji, Zhaojing Yang, Zaitian Gongye, Xueyan Zou, Jan Kautz, Erdem Biyik, Hongxu Yin, Sifei Liu, and Xiaolong Wang. Navila: Legged robot vision-language-action model for navigation. *arXiv preprint arXiv:2412.04453*, 2024.
- Muzhi Dai, Chenxu Yang, and Qingyi Si. S-grpo: Early exit via reinforcement learning in reasoning models. *arXiv preprint arXiv:2505.07686*, 2025.

- 594 Gongfan Fang, Xinyin Ma, and Xinchao Wang. Thinkless: Llm learns when to think. *arXiv preprint*
595 *arXiv:2505.13379*, 2025.
- 596
- 597 Yichao Fu, Xuwei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence. *arXiv*
598 *preprint arXiv:2508.15260*, 2025.
- 599
- 600 Tobias Groot and Matias Valdenegro-Toro. Overconfidence is key: Verbalized uncertainty evaluation
601 in large language and vision-language models. *arXiv preprint arXiv:2405.02917*, 2024.
- 602
- 603 Jing Gu, Eliana Stefani, Qi Wu, Jesse Thomason, and Xin Eric Wang. Vision-and-language nav-
604 igation: A survey of tasks, methods, and future directions. *arXiv preprint arXiv:2203.12667*,
605 2022.
- 606
- 607 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
608 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 609
- 610 Yicong Hong, Cristian Rodriguez, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual en-
611 tity relationship graph for agent navigation. *Advances in Neural Information Processing Systems*,
33:7685–7696, 2020.
- 612
- 613 Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordani, Siva Reddy,
614 Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through
615 refined credit assignment. 2024.
- 616
- 617 Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-
618 graph: Vision-and-language navigation in continuous environments. In *European Conference on*
619 *Computer Vision*, pp. 104–120. Springer, 2020.
- 620
- 621 Jacob Krantz, Aaron Gokaslan, Dhruv Batra, Stefan Lee, and Oleksandr Maksymets. Waypoint
622 models for instruction-guided navigation in continuous environments. In *Proceedings of the*
623 *IEEE/CVF International Conference on Computer Vision*, pp. 15162–15171, 2021.
- 624
- 625 Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldrige. Room-across-room:
626 Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint*
arXiv:2010.07954, 2020.
- 627
- 628 Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language
629 models. In *European Conference on Computer Vision*, pp. 323–340. Springer, 2024.
- 630
- 631 Bingqian Lin, Yunshuang Nie, Ziming Wei, Jiaqi Chen, Shikui Ma, Jianhua Han, Hang Xu, Xi-
632 aojun Chang, and Xiaodan Liang. Navcot: Boosting llm-based vision-and-language navigation
633 via learning disentangled reasoning. *IEEE Transactions on Pattern Analysis and Machine Intelli-*
gence, 2025a.
- 634
- 635 Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwo: A
636 unified vision-language-action model with adaptive reasoning. *arXiv preprint arXiv:2505.11917*,
637 2025b.
- 638
- 639 Haoran Liu, Weikang Wan, Xiqian Yu, Minghan Li, Jiazhao Zhang, Bo Zhao, Zhibo Chen,
640 Zhongyuan Wang, Zhizheng Zhang, and He Wang. Na vid-4d: Unleashing spatial intelligence in
641 egocentric rgb-d videos for vision-and-language navigation. In *2025 IEEE International Confer-*
ence on Robotics and Automation (ICRA), pp. 10607–10615. IEEE, 2025.
- 642
- 643 Yuxing Long, Wenzhe Cai, Hongcheng Wang, Guanqi Zhan, and Hao Dong. Instructnav: Zero-
644 shot system for generic instruction navigation in unexplored environment. *arXiv preprint*
645 *arXiv:2406.04882*, 2024a.
- 646
- 647 Yuxing Long, Xiaoqi Li, Wenzhe Cai, and Hao Dong. Discuss before moving: Visual language
navigation via multi-expert discussions. In *2024 IEEE International Conference on Robotics and*
Automation (ICRA), pp. 17380–17387. IEEE, 2024b.

- 648 Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming
649 Xiong. Self-monitoring navigation agent via auxiliary progress estimation. *arXiv preprint*
650 *arXiv:1901.03035*, 2019.
- 651
- 652 Sang-Min Park and Young-Gab Kim. Visual language navigation: A survey and open challenges.
653 *Artificial Intelligence Review*, 56(1):365–427, 2023.
- 654
- 655 Yuankai Qi, Zizheng Pan, Shengping Zhang, Anton van den Hengel, and Qi Wu. Object-and-action
656 aware model for visual language navigation. In *European conference on computer vision*, pp.
657 303–317. Springer, 2020a.
- 658
- 659 Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton
660 van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environ-
661 ments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
662 pp. 9982–9991, 2020b.
- 663
- 664 Yanyuan Qiao, Yuankai Qi, Zheng Yu, Jing Liu, and Qi Wu. March in chat: Interactive prompting for
665 remote embodied referring expression. In *Proceedings of the IEEE/CVF international conference*
666 *on computer vision*, pp. 15758–15767, 2023.
- 667
- 668 Matt Schmittle, Anna Lukina, Lukas Vacek, Jnaneshwar Das, Christopher P Buskirk, Stephen Rees,
669 Janos Sztipanovits, Radu Grosu, and Vijay Kumar. Openuav: A uav testbed for the cps and
670 robotics community. In *2018 ACM/IEEE 9th International Conference on Cyber-Physical Sys-*
671 *tems (ICCPS)*, pp. 130–139. IEEE, 2018.
- 672
- 673 Yi Shen, Jian Zhang, Jieyun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai
674 Wang, Zhaoxiang Liu, and Shiguo Lian. Dast: Difficulty-adaptive slow-thinking for large reason-
675 ing models. *arXiv preprint arXiv:2503.04472*, 2025.
- 676
- 677 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
678 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
679 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2998–3009, 2023.
- 680
- 681 Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu,
682 Andrew Wen, Shaochen Zhong, Hanjie Chen, et al. Stop overthinking: A survey on efficient
683 reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.
- 684
- 685 Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. Large language models are overconfi-
686 dent and amplify human bias. *arXiv preprint arXiv:2505.02151*, 2025.
- 687
- 688 Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back transla-
689 tion with environmental dropout. *arXiv preprint arXiv:1904.04195*, 2019.
- 690
- 691 Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. Vision-and-dialog navi-
692 gation. In *Conference on Robot Learning*, pp. 394–406. PMLR, 2020.
- 693
- 694 Xiaoyu Tian, Sitong Zhao, Haotian Wang, Shuaiting Chen, Yunjie Ji, Yiping Peng, Han Zhao, and
695 Xiangang Li. Think twice: Enhancing llm reasoning by scaling multi-round test-time thinking.
696 *arXiv preprint arXiv:2503.19855*, 2025.
- 697
- 698 Shenzi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,
699 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive
700 effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025a.
- 701
- 702 Shuo Wang, Yongcai Wang, Wanting Li, Yucheng Wang, Maiyue Chen, Kaihui Wang, Zhizhong Su,
703 Xudong Cai, Yeying Jin, Deying Li, et al. Monodream: Monocular vision-language navigation
704 with panoramic dreaming. *arXiv preprint arXiv:2508.02549*, 2025b.
- 705
- 706 Tian Wang, Junming Fan, and Pai Zheng. An llm-based vision and language cobot navigation
707 approach for human-centric smart manufacturing. *Journal of Manufacturing Systems*, 75:299–
708 305, 2024.

- 702 Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang,
703 William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised im-
704 itation learning for vision-language navigation. In *Proceedings of the IEEE/CVF conference on*
705 *computer vision and pattern recognition*, pp. 6629–6638, 2019.
- 706 Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, and Shuqiang Jiang. Gridmm: Grid memory map
707 for vision-and-language navigation. In *Proceedings of the IEEE/CVF International conference on*
708 *computer vision*, pp. 15625–15636, 2023.
- 709 Meng Wei, Chenyang Wan, Xiqian Yu, Tai Wang, Yuqiang Yang, Xiaohan Mao, Chenming Zhu,
710 Wenzhe Cai, Hanqing Wang, Yilun Chen, et al. Streamvln: Streaming vision-and-language navi-
711 gation via slowfast context modeling. *arXiv preprint arXiv:2507.05240*, 2025.
- 712 Bingbing Wen, Chenjun Xu, HAN Bin, Robert Wolfe, Lucy Lu Wang, and Bill Howe. Mitigat-
713 ing overconfidence in large language models: A behavioral lens on confidence estimation and
714 calibration. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*, volume 1, 2024.
- 715 Yuyang Wu, Yifei Wang, Ziyu Ye, Tianqi Du, Stefanie Jegelka, and Yisen Wang. When more is
716 less: Understanding chain-of-thought length in llms. *arXiv preprint arXiv:2502.07266*, 2025.
- 717 Chengguang Xu, Hieu T Nguyen, Christopher Amato, and Lawson LS Wong. Vision and
718 language navigation in the real world via online visual language mapping. *arXiv preprint*
719 *arXiv:2310.10822*, 2023.
- 720 Hang Yin, Xiuwei Xu, Zhenyu Wu, Jie Zhou, and Jiwen Lu. Sg-nav: Online 3d scene graph prompt-
721 ing for llm-based zero-shot object navigation. *Advances in neural information processing systems*,
722 37:5285–5307, 2024.
- 723 Minji Yoo. How much should we trust llm-based measures for accounting and finance research?
724 Available at SSRN, 2024.
- 725 Zhuoyuan Yu, Yuxing Long, Zihan Yang, Chengyan Zeng, Hongwei Fan, Jiyao Zhang, and Hao
726 Dong. Correctnav: Self-correction flywheel empowers vision-language-action navigation model.
727 *arXiv preprint arXiv:2508.10416*, 2025.
- 728 Zhaohuan Zhan, Lisha Yu, Sijie Yu, and Guang Tan. Mc-gpt: Empowering vision-and-language
729 navigation with memory map and reasoning chains. *arXiv preprint arXiv:2405.10620*, 2024.
- 730 Jiazhao Zhang, Kunyu Wang, Shaoan Wang, Minghan Li, Haoran Liu, Songlin Wei, Zhongyuan
731 Wang, Zhizheng Zhang, and He Wang. Uni-navid: A video-based vision-language-action model
732 for unifying embodied navigation tasks. *arXiv preprint arXiv:2412.06224*, 2024a.
- 733 Jiazhao Zhang, Kunyu Wang, Rongtao Xu, Gengze Zhou, Yicong Hong, Xiaomeng Fang, Qi Wu,
734 Zhizheng Zhang, and He Wang. Navid: Video-based vlm plans the next step for vision-and-
735 language navigation. *arXiv preprint arXiv:2402.15852*, 2024b.
- 736 Xiaoyun Zhang, Jingqing Ruan, Xing Ma, Yawen Zhu, Haodong Zhao, Hao Li, Jiansong Chen,
737 Ke Zeng, and Xunliang Cai. When to continue thinking: Adaptive thinking mode switching for
738 efficient reasoning. *arXiv preprint arXiv:2505.15400*, 2025.
- 739 Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and Liwei Wang. Towards learning a generalist
740 model for embodied navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
741 *and Pattern Recognition*, pp. 13624–13634, 2024.
- 742 Gengze Zhou, Yicong Hong, Zun Wang, Xin Eric Wang, and Qi Wu. Navgpt-2: Unleashing naviga-
743 tional reasoning capability for large vision-language models. In *European Conference on Com-*
744 *puter Vision*, pp. 260–278. Springer, 2024a.
- 745 Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language naviga-
746 tion with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
747 volume 38, pp. 7641–7649, 2024b.
- 748 Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart,
749 Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge
750 to robotic control. In *Conference on Robot Learning*, pp. 2165–2183. PMLR, 2023.

A LARGE LANGUAGE MODELS (LLMs) IN PAPER WRITING

In this work, Large Language Models (LLMs) were only employed for proofreading, spelling, and formatting checks. LLMs did not contribute to any aspect of research idea generation, experimental design, or coding.

B IMPLEMENTATION DETAILS

Real-World Evaluation We provide a detailed description of our real-robot platform in Figure 5. The mobile robot is equipped with basic locomotion capabilities and augmented with a camera, microphone, speaker, and LiDAR sensors for user interaction and environment perception. Notably, our method relies only on RGB images and does not require LiDAR input. The system is powered by a Jetson AGX Orin (nvi) running Ubuntu 24.04 with ROS2 Jazzy. In addition, the platform integrates a 19V power regulator and a 110V/220V inverter (both rated at 500W+) to support the compute and sensor modules.

Building on this platform, we design a pipeline for vision-and-language navigation with AdaNav. We experiment with AdaNav on two base models, NAVID (Zhang et al., 2024b) and NAVILA (Cheng et al., 2024). In deployment, AdaNav runs on a server with a Jetson that receives compressed images from the robot and sends back high-level commands. The robot then executes these commands (e.g., “Turn right” or “Move forward”) through its locomotion system. During navigation, the robot continuously monitors its motion to ensure that rotations and forward movements remain aligned with the issued commands.

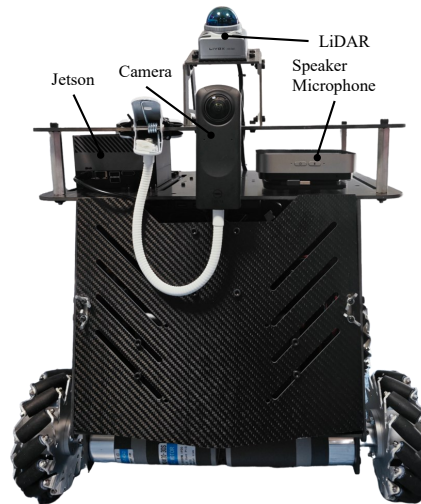


Figure 5: Hardware setup of the mobile robot platform used for real-world evaluation.



Figure 6: Visualization of adaptive reasoning navigation in real-world scene.

C DETAILED ABLATION ANALYSIS FOR ROBUSTNESS AND SENSITIVITY

To further analyze our method, we conduct ablations on different reasoning modes. As shown in Figure 4b, different modes contribute at different stages of navigation. This raises the question: *how*

810 *much would performance degrade if only a single reasoning mode were available?* To investigate,
 811 we construct variants where only one reasoning mode is enabled (i.e., description, summary, or
 812 error correction). In addition, we add a control setting, where the model is simply instructed with
 813 a generic prompt such as: *"At this point, perform some analysis based on the past trajectory."* For
 814 all these variants, we apply UAR Block together with the heuristic-to-RL training mechanism, and
 815 report the results as follows.

816 Table 10: Ablation results under different *reasoning modes*, where ①: Description only, ②: Sum-
 817 mary only, ③: Error correction only, and ④: Generic prompt reasoning. All variants are trained with
 818 UAR Block and the heuristic-to-RL mechanism.
 819

Method	R2R-CE Val-Unseen			
	NE↓	OS↑	SR↑	SPL↑
AdaNav	5.39	57.89	47.7	42.34
AdaNav ①	5.41	55.34	44.52	40.88
AdaNav ②	5.42	56.53	46.01	41.37
AdaNav ③	5.40	56.12	45.77	41.67
AdaNav ④	5.40	55.23	45.62	42.01

828 D VISUALIZATION RESULTS ON ADAPTIVE REASONING NAVIGATION

829 To better understand how AdaNav adaptively allocates reasoning during navigation, we visualize
 830 example trajectories in Figure 7. The figure illustrates both the agent’s path and the steps where
 831 reasoning is invoked. As shown, AdaNav selectively triggers reasoning at critical or challenging
 832 moments, while skipping unnecessary steps in simpler regions of the environment. This demon-
 833 strates that the Uncertainty-Adaptive Reasoning (UAR) Block effectively guides the agent to balance
 834 efficiency and accuracy.
 835

836 Figure 7 also highlights that reasoning is concentrated on hard trajectories: compared to the first,
 837 simpler scenario, the second, more complex instruction involves more reasoning steps. This ob-
 838 servation is consistent with our quantitative analysis in Table 5. These visualizations confirm that
 839 the agent’s reasoning behavior is both difficulty-aware and mode-adaptive, providing interpretability
 840 and insight into its decision-making process.
 841

842 E ADDITIONAL EXPERIMENTS WITH THE UNI-NAVID BACKBONE

843 To further evaluate the effectiveness of AdaNav, we conducted additional experiments using the
 844 newly released Uni-NaVid backbone. We also extended the evaluation to additional continuous-
 845 control tasks, including:
 846

- 847 • Object Goal Navigation on Matterport3D, and
- 848 • Object Goal Navigation on HM3D-OVON.

849 The results of these experiments are summarized in Tables 11, 12, and 13.
 850

851 Table 11: Results on Object Goal Navigation (Matterport3D and HM3D-OVON).
 852

Method	SR	SPL
PIRLNav-IL-RL	70.4	34.1
OVRL	62.0	26.8
OVRL-v2	64.7	28.1
Uni-NaVid	73.7	37.1
Uni-NaVid-AdaNav	80.15	42.2

860 These additional results demonstrate that AdaNav consistently improves navigation performance
 861 across different backbones, environments, and task settings, including the newly released Uni-NaVid
 862 backbone.
 863

Table 12: Results on OVON tasks (Val Seen, Seen Synonyms, Val Unseen).

Method	SR	SPL	SR	SPL	SR	SPL
	Val Seen		Seen Synonyms		Val Unseen	
DAgRL	41.3	21.2	29.4	14.4	18.3	7.9
VLFM	35.2	18.6	32.4	17.3	35.2	19.6
DAgRL+OD	38.5	21.1	39.0	21.4	37.1	19.8
Uni-NaVid	41.3	21.1	43.9	21.8	39.5	19.8
Uni-NaVid-AdaNav	50.62	25.62	52.61	27.8	47.8	26.2

Table 13: Results on R2R and RxR benchmarks.

Method	R2R					RxR				
	TL	NE↓	OS↑	SR↑	SPL↑	TL	NE↓	OS↑	SR↑	SPL↑
Uni-NaVid	9.71	5.58	53.3	47.0	42.7	15.8	6.24	55.5	48.7	40.9
Uni-NaVid-AdaNav	9.91	5.11	61.6	56.9	49.5	17.7	6.20	59.6	54.5	44.2

F COMPARISON OF ALTERNATIVE UNCERTAINTY MEASURES

To evaluate different uncertainty signals for guiding adaptive reasoning, we designed and compared several measures: **value uncertainty**, **value disagreement**, and **temporal variance**, in addition to **action entropy**, which serves as our core uncertainty signal. All signals were trained using the Heuristics-to-RL training method proposed in our work.

F.1 VALUE UNCERTAINTY

Following prior work (Wang et al., 2019), an additional critic head is trained to estimate the expected return. The variance across multiple value heads

$$\sigma_V^2 = \text{Var}(V_t^{(i)})$$

provides a measure of value-based uncertainty, capturing inconsistencies in future reward estimation.

F.2 VALUE DISAGREEMENT

Value disagreement quantifies diversity among predictions of multiple critic heads. Specifically, the training samples are split into three subsets, and one critic head is trained on each subset. The disagreement is computed as the average pairwise squared difference:

$$D_V = \frac{2}{N(N-1)} \sum_{i < j} (V_t^{(i)} - V_t^{(j)})^2,$$

where $V_t^{(i)}$ denotes the value predicted by the i -th head at time t . Larger D_V indicates greater disagreement among the value heads.

F.3 TEMPORAL VARIANCE

Temporal variance measures short-term fluctuations of action entropy over a sliding window of length w :

$$\sigma_H^2 = \frac{1}{w} \sum_{k=t-w+1}^t (H(\pi_k) - \bar{H})^2, \quad \bar{H} = \frac{1}{w} \sum_{k=t-w+1}^t H(\pi_k),$$

where $H(\pi_k)$ is the entropy of the action distribution at step k . Large σ_H^2 indicates unstable decision confidence.

As shown in Table 14, using action entropy as the core uncertainty signal provides better perception of task difficulty, achieves higher navigation performance, and improves reasoning efficiency compared with alternative uncertainty measures.

Table 14: Comparison of alternative uncertainty measures on navigation performance.

Component	NE↓	OS↑	SR↑	SPL↑	Avg. reasoning num
NaVid	5.47	49.0	37.0	35.0	0
Value Uncertainty	5.41	54.31	44.69	40.67	4.8
Value Disagreement	5.44	54.78	43.73	40.53	5.5
Temporal Variance	5.40	54.65	45.21	41.15	4.2
NaVid-AdaNav	5.39	57.89	47.7	42.34	2.5

G PRELIMINARY EXPLORATION IN OUTDOOR VLN TASKS

Current VLN base models perform poorly on outdoor tasks due to several challenges. First, training data for outdoor VLN is scarce, as most available datasets focus on indoor navigation. Second, generalization from indoor to outdoor environments is limited, even with existing indoor datasets. As a result, current models achieve low performance on outdoor VLN benchmarks.

To investigate this scenario, we conducted experiments on the **OpenUAV (Schmittle et al., 2018) outdoor VLN benchmark** as well as in **real-world environments**, including both simple and complex settings. Remarkably, even without access to outdoor training data, **AdaNav** was able to improve the performance of the base model. This improvement is attributed to AdaNav’s reasoning mechanism, which leverages knowledge acquired by the base VLM from large-scale internet datasets during pretraining, partially compensating for the scarcity of outdoor data during post-training.

Table 15: Performance of VLN models on OpenUAV and real-world outdoor tasks.

	OpenUAV				Real World-Simple		Real World-Complex	
	NE↓	SR↑	OSR↑	SPL↑	SR↑	NE↓	SR↑	NE↓
CMA	102.92	14.83	22.49	13.90	-	-	-	-
Navila	93	21.5	45.1	21.5	55.0	2.5	35.0	3.5
AdaNav	85	25.5	50.1	23.8	67.0	1.9	45.0	2.5

H ANALYSIS OF UAR BLOCK FAILURE CASES AND IMPROVEMENT

The UAR Block occasionally fails to trigger reasoning under transient ambiguities, such as vague instructions or sudden changes in lighting and viewpoint. In these cases, the action entropy may fluctuate briefly but does not remain high enough to surpass the reasoning threshold, resulting in missed reasoning opportunities.

To address this issue, we conducted RL fine-tuning on data specifically designed to cover such scenarios. During this process, the model learns to associate transient uncertainty with the need for reasoning. We collected 500 samples exhibiting instruction vagueness and lighting changes and conducted additional evaluations on R2R and RxR benchmarks.

Table 16: Performance comparison of UAR Block on R2R and RxR benchmarks after targeted RL fine-tuning.

Method	R2R				RxR			
	NE↓	OS↑	SR↑	SPL↑	NE↓	OS↑	SR↑	nDTW↑
NaVid	5.47	49.0	37.0	35.0	6.79	46.2	40.5	52.2
NaVid-AdaNav	5.39	57.89	47.7	42.34	6.38	58.1	47.01	56.8
NaVid-AdaNav*	5.35	58.72	48.85	42.88	6.12	59.62	48.55	57.2

972 These results indicate that tailoring the training environment to address transient ambiguity scenarios
 973 further improves navigation performance and reduces UAR Block failure cases.
 974

975 I COMPARISON WITH FIXED-INTERVAL REASONING STRATEGIES

976
 977 To investigate whether adaptive uncertain inference can result in a higher actual reasoning frequency
 978 while maintaining better performance, we additionally evaluated a series of fixed-interval reasoning
 979 strategies (1, 3, 5, 7, 9, 15 steps) without RL training.
 980

981 Table 17: Navigation performance of Navid and fixed-interval reasoning strategies without RL train-
 982 ing.
 983

984 Component	985 NE↓	986 OS↑	987 SR↑	988 SPL↑
989 Navid	5.47	49.0	37.0	35.0
990 fixed 1-step	5.49	51.2	38.2	36.5
991 fixed 3-step	5.47	50.35	39.51	37.87
992 fixed 5-step	5.45	53.25	40.12	38.83
993 fixed 7-step	5.45	52.56	40.1	37.55
994 fixed 9-step	5.43	51.25	40.02	38.35
995 fixed 15-step	5.45	52.23	39.2	37.48
996 Navid-AdaNav	5.39	57.89	47.7	42.34

997
 998 As shown in Table 17, Navid-AdaNav consistently outperforms all fixed-interval strategies, demon-
 999 strating that the RL-based UAR Block enables more effective reasoning, achieving both higher nav-
 1000 igation performance and efficient adaptive inference.
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074


1075

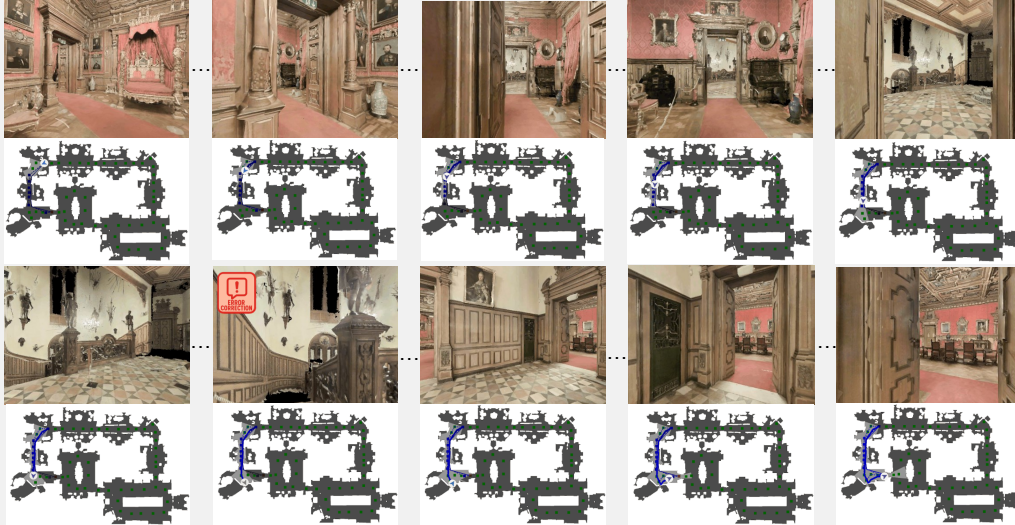
1076



1077


1078

1079



 Your few steps away the entrance, facing towards the pillar, turn slightly left and go straight near the another door, towards your left, enter in to the room, **walk through the walkway and enter in to the corridor, from there turn left and go near the door which is open**, now you are at the entrance of the room, in front of you there is a dinning table, this would be your final destination.



  I started near the entrance facing a pillar, then veered slightly left and headed toward a nearby door on my left. After entering the room, I moved along the path inside and made my way into the corridor. I have now reached a staircase landing, which does not match the given instructions

 You're facing towards the wooden cabinet onto your right side, now slightly turn to your left, and walk a few steps forward, you can see an ash color sofa onto your right and from there turn to your left, walk a few steps forward and towards your right you can see a wooden table with a plant on it. Walk a few steps forward, you can see an open entrance in front of you, **now from the open entrance, turn to your right** you can see an open wooden door, enter the door, walk a few steps forward and towards your left you can see another open door, enter it. Now you're in a bedroom, **go straight and from the bed, turn towards your left, you enter into a washroom**, you're facing towards the sink and that is your destination.



  A painting is hanging in the center of the wall. On the left side, there is a doorway leading to a kitchen area with dark cabinets. Wooden-framed glass doors are visible on both the left and right sides.


 I moved through the living area, passed a sofa and a wooden table, entered successive doors, and finally arrived in the bedroom, which match the given instructions.

Figure 7: Visualizations of adaptive reasoning navigation, where AdaNav autonomously invokes reasoning at high-uncertainty points to better align the trajectory with the instruction.