# AdaNav: Adaptive Reasoning with Uncertainty for Vision-Language Navigation

Xin Ding[1]*    Jianyu Wei[1]    Yifan Yang[2]    Shiqi Jiang[2]    Qianxi Zhang[2]    Hao Wu[3]    Fucheng Jia[4,6]

Liang Mi[6]    Yuxuan Yan[5]    Weijun Wang[6]    Yunxin Liu[6]    Zhibo Chen[1] †    Ting Cao[6] †

[1]University of Science and Technology of China    [2] Microsoft Research    [3]Nanjing University

[4] Central South University    [5]Zhejiang University

[6]Institute for AI Industry Research (AIR), Tsinghua University

## 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. The code is available at AdaNav.

## 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

---

*xinding64@mail.ustc.edu.cn

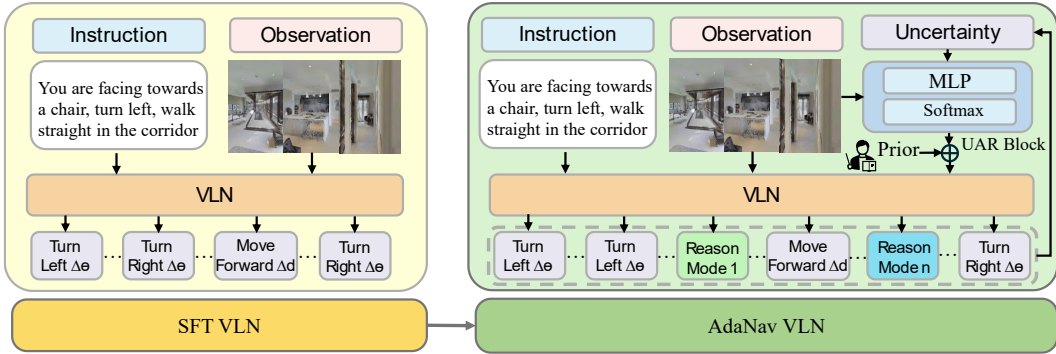†Corresponding Author. tingcao@mail.tsinghua.edu.cn

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.

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 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.,
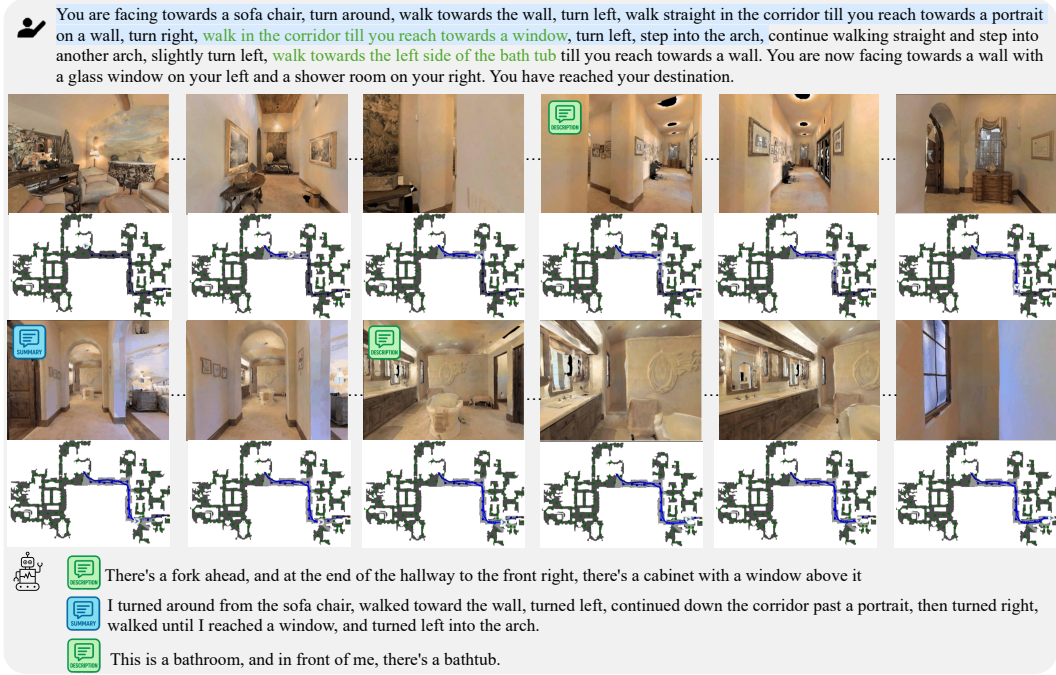
Figure 2: A visualization example of AdaNav. It autonomously invokes reasoning, e.g., summarization and description when necessary to maintain consistent temporal grounding and robust perception–action mapping.

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

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

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

While these rule-driven frameworks offer interpretability, they inherently restrict flexibility in open-ended environments, hinder efficiency, and may lead to overthinking (Fang et al., 2025; Dai et al., 2025). In contrast, our method will employ a learnable mechanism that enables agents to autonomously 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 $\mathcal{A} = \{\mathrm{turn\_left}(\Delta\theta), \mathrm{turn\_right}(\Delta\theta), \mathrm{move\_forward}(\Delta d), \mathrm{stop}\}$, where $\Delta\theta$ and $\Delta d$ 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:

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

**Adaptive Reasoning Navigation.** To improve VLN performance in long-horizon and complex environments, we allow explicit reasoning at step $t$ with a *mode* variable $m_t \in \{\varnothing\} \cup \mathcal{M}$ and reasoning content $r_t$. Here, $m_t = \varnothing$ denotes *no reasoning* (so $r_t = \varnothing$), 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 C for instances). The agent's policy then consists of two coupled processes: 1) a navigation policy $\pi_{\mathrm{nav}}(a_t \mid h_t^{\mathrm{nav}}, I, r_{\leq t})$, and 2) a reasoning policy $\pi_{\mathrm{rea}}(m_t, r_t \mid h_t^{\mathrm{rea}}, I)$ that jointly decides *when* to reason (via $m_t = \varnothing$ vs. $m_t \neq \varnothing$) and *which mode* to use (via $m_t \in \mathcal{M}$).

The overall joint policy is

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

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

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

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

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

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

where $R_{\mathrm{task}}(\tau)$ 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
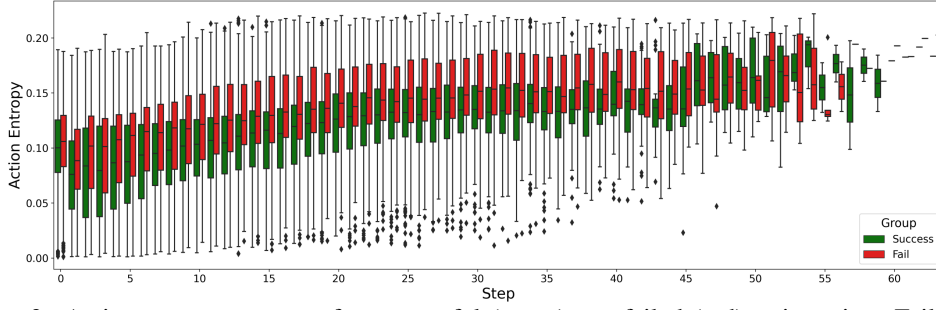
Figure 3: Action entropy per step for successful (green) vs. failed (red) trajectories. Failed trajectories show higher, especially in later steps, indicating that policy uncertainty correlates with navigation errors.

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.

### 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:

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

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:

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

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

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

### 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):

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

where $\{\tau_0, \tau_1, \ldots, \tau_{|\mathcal{M}|}\}$ are mode-specific entropy thresholds ($\tau_k = \tau_0 + k\delta$), and $\sigma$ 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:

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

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

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

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

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

**Algorithm 1** Heuristic-to-RL
1: Initialize navigation policy $\pi_{\text{nav}}$, reasoning selector $\pi_{\text{sel}}$, annealing schedule $\lambda_t$
2: **for** each training episode $\tau$ **do**
3:     **for** each step $t = 1 \ldots T$ **do**
4:         Observe state $o_t$ and entropy $H_{\leq t}$
5:         Compute control vector $p_{\text{mode}}$
6:         Estimate Uncertainty prior $p_{\text{prior}}$
7:         Fuse prior and model (Eq. 9)
8:         Sample reasoning mode $m_t \sim p_{\text{final}}$
9:         **if** $m_t \neq \varnothing$ **then**
10:             Generate reasoning $r_t$
11:         **end if**
12:         Select action $a_t \sim \pi_{\text{nav}}(a_t \mid o_t, r_{\leq t})$
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**

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 $L_{\text{window}}$ 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 $s_{\text{target}}$.

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

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

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 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 $1 \times 10^{-6}$. **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 A.

**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 MONODREAM (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 |

**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. 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 sample or complex instructions. Each instruction requires the agent to complete 5–10 sequential landmark-following sub-tasks (e.g., "After

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-Langauge-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 |

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

| Method | Meeting Room | | | | Home | | | | Office | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Sample | | Complex | | Sample | | Complex | | Sample | | 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 |

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*, and *Office*, following the protocol in prior works (Cheng et al., 2024; Zhang et al., 2024b). The results are summarized in Table 4.

## 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.
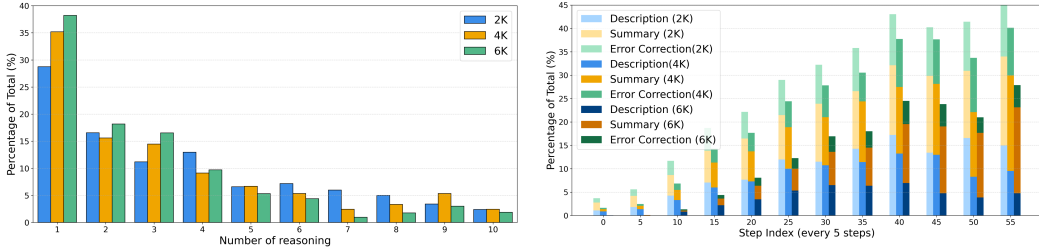
Table 6: Performance on R2R.

| 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 |

| 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.

(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.

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 B.

**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.

| Component | NE↓ | OS↑ | SR↑ | SPL↑ |
|---|---|---|---|---|
| Navid | 5.47 | 49.0 | 37.0 | 35.0 |
| w/o UAR* | 5.45 | 53.25 | 40.12 | 38.83 |
| w/o UAR† | 5.44 | 52.10 | 39.5 | 38.65 |
| w/o RL | 5.44 | 55.33 | 42.53 | 39.65 |
| w/o HP | 5.41 | 55.73 | 43.82 | 40.12 |
| AdaNav | 5.39 | 57.89 | 47.7 | 42.34 |

(a) Component ablation.

| $\tau_0$ | $\delta$ | NE↓ | OS↑ | SR↑ | SPL↑ |
|---|---|---|---|---|---|
| 80% | 0.1 | 5.43 | 57.72 | 48.82 | 43.56 |
| | 0.2 | 5.42 | 57.92 | 49.11 | 43.53 |
| | 0.3 | 5.42 | 58.01 | 49.05 | 43.57 |
| 85% | 0.1 | 5.40 | 58.75 | 49.61 | 43.87 |
| | 0.2 | 5.39 | 57.89 | 47.7 | 42.34 |
| | 0.3 | 5.39 | 58.81 | 49.53 | 43.85 |
| 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 |

(b) Effect of $(\tau_0, \delta)$.

| $\sigma$ | NE↓ | OS↑ | SR↑ | SPL↑ |
|---|---|---|---|---|
| 0.05 | 5.43 | 58.85 | 48.85 | 43.55 |
| 0.10 | 5.40 | 58.55 | 49.13 | 43.62 |
| 0.15 | 5.39 | 57.89 | 47.7 | 42.34 |
| 0.20 | 5.41 | 58.73 | 49.55 | 44.02 |
| 0.25 | 5.44 | 58.72 | 49.25 | 43.88 |
| 0.30 | 5.48 | 58.13 | 48.72 | 43.92 |

(c) Effect of $\sigma$.

**Hyperparameter Sensitivity.** The key hyperparameters in our framework lie in the *Heuristic-to-RL* stage, where we introduce mode-specific entropy thresholds: $(\tau_0, \delta)$ that govern reasoning triggers prior, and a scaling factor $\sigma$.

As shown in Table 7a, a well reasoning prior significantly facilitates training. Specifically, $\tau_0$ 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, $\delta$ incrementally shifts the thresholds for higher reasoning modes, thereby shaping the curriculum schedule. The corresponding results are summarized in Table 7b and Table 7c.

## 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.

## 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. Etp-nav: 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 Bıyık, 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.

Gongfan Fang, Xinyin Ma, and Xinchao Wang. Thinkless: Llm learns when to think. *arXiv preprint arXiv:2505.13379*, 2025.

Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence. *arXiv preprint arXiv:2508.15260*, 2025.

Tobias Groot and Matias Valdenegro-Toro. Overconfidence is key: Verbalized uncertainty evaluation in large language and vision-language models. *arXiv preprint arXiv:2405.02917*, 2024.

Jing Gu, Eliana Stefani, Qi Wu, Jesse Thomason, and Xin Eric Wang. Vision-and-language navigation: A survey of tasks, methods, and future directions. *arXiv preprint arXiv:2203.12667*, 2022.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

Yicong Hong, Cristian Rodriguez, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual entity relationship graph for agent navigation. *Advances in Neural Information Processing Systems*, 33:7685–7696, 2020.

Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy, Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment. 2024.

Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, pp. 104–120. Springer, 2020.

Jacob Krantz, Aaron Gokaslan, Dhruv Batra, Stefan Lee, and Oleksandr Maksymets. Waypoint models for instruction-guided navigation in continuous environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15162–15171, 2021.

Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint arXiv:2010.07954*, 2020.

Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In *European Conference on Computer Vision*, pp. 323–340. Springer, 2024.

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*, 2025a.

Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning. *arXiv preprint arXiv:2505.11917*, 2025b.

Haoran Liu, Weikang Wan, Xiqian Yu, Minghan Li, Jiazhao Zhang, Bo Zhao, Zhibo Chen, Zhongyuan Wang, Zhizheng Zhang, and He Wang. Na vid-4d: Unleashing spatial intelligence in egocentric rgb-d videos for vision-and-language navigation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 10607–10615. IEEE, 2025.

Yuxing Long, Wenzhe Cai, Hongcheng Wang, Guanqi Zhan, and Hao Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. *arXiv preprint arXiv:2406.04882*, 2024a.

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.

Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. *arXiv preprint arXiv:1901.03035*, 2019.

Sang-Min Park and Young-Gab Kim. Visual language navigation: A survey and open challenges. *Artificial Intelligence Review*, 56(1):365–427, 2023.

Yuankai Qi, Zizheng Pan, Shengping Zhang, Anton van den Hengel, and Qi Wu. Object-and-action aware model for visual language navigation. In *European conference on computer vision*, pp. 303–317. Springer, 2020a.

Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9982–9991, 2020b.

Yanyuan Qiao, Yuankai Qi, Zheng Yu, Jing Liu, and Qi Wu. March in chat: Interactive prompting for remote embodied referring expression. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 15758–15767, 2023.

Yi Shen, Jian Zhang, Jieyun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai Wang, Zhaoxiang Liu, and Shiguo Lian. Dast: Difficulty-adaptive slow-thinking for large reasoning models. *arXiv preprint arXiv:2503.04472*, 2025.

Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2998–3009, 2023.

Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, et al. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.

Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. Large language models are overconfident and amplify human bias. *arXiv preprint arXiv:2505.02151*, 2025.

Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. *arXiv preprint arXiv:1904.04195*, 2019.

Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. Vision-and-dialog navigation. In *Conference on Robot Learning*, pp. 394–406. PMLR, 2020.

Xiaoyu Tian, Sitong Zhao, Haotian Wang, Shuaiting Chen, Yunjie Ji, Yiping Peng, Han Zhao, and Xiangang Li. Think twice: Enhancing llm reasoning by scaling multi-round test-time thinking. *arXiv preprint arXiv:2503.19855*, 2025.

Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025a.

Shuo Wang, Yongcai Wang, Wanting Li, Yucheng Wang, Maiyue Chen, Kaihui Wang, Zhizhong Su, Xudong Cai, Yeying Jin, Deying Li, et al. Monodream: Monocular vision-language navigation with panoramic dreaming. *arXiv preprint arXiv:2508.02549*, 2025b.

Tian Wang, Junming Fan, and Pai Zheng. An llm-based vision and language cobot navigation approach for human-centric smart manufacturing. *Journal of Manufacturing Systems*, 75:299–305, 2024.

Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6629–6638, 2019.

Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, and Shuqiang Jiang. Gridmm: Grid memory map for vision-and-language navigation. In *Proceedings of the IEEE/CVF International conference on computer vision*, pp. 15625–15636, 2023.

Meng Wei, Chenyang Wan, Xiqian Yu, Tai Wang, Yuqiang Yang, Xiaohan Mao, Chenming Zhu, Wenzhe Cai, Hanqing Wang, Yilun Chen, et al. Streamvln: Streaming vision-and-language navigation via slowfast context modeling. *arXiv preprint arXiv:2507.05240*, 2025.

Bingbing Wen, Chenjun Xu, HAN Bin, Robert Wolfe, Lucy Lu Wang, and Bill Howe. Mitigating overconfidence in large language models: A behavioral lens on confidence estimation and calibration. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*, volume 1, 2024.

Yuyang Wu, Yifei Wang, Ziyu Ye, Tianqi Du, Stefanie Jegelka, and Yisen Wang. When more is less: Understanding chain-of-thought length in llms. *arXiv preprint arXiv:2502.07266*, 2025.

Chengguang Xu, Hieu T Nguyen, Christopher Amato, and Lawson LS Wong. Vision and language navigation in the real world via online visual language mapping. *arXiv preprint arXiv:2310.10822*, 2023.

Hang Yin, Xiuwei Xu, Zhenyu Wu, Jie Zhou, and Jiwen Lu. Sg-nav: Online 3d scene graph prompting for llm-based zero-shot object navigation. *Advances in neural information processing systems*, 37:5285–5307, 2024.

Minji Yoo. How much should we trust llm-based measures for accounting and finance research? *Available at SSRN*, 2024.

Zhuoyuan Yu, Yuxing Long, Zihan Yang, Chengyan Zeng, Hongwei Fan, Jiyao Zhang, and Hao Dong. Correctnav: Self-correction flywheel empowers vision-language-action navigation model. *arXiv preprint arXiv:2508.10416*, 2025.

Zhaohuan Zhan, Lisha Yu, Sijie Yu, and Guang Tan. Mc-gpt: Empowering vision-and-language navigation with memory map and reasoning chains. *arXiv preprint arXiv:2405.10620*, 2024.

Jiazhao Zhang, Kunyu Wang, Shaoan Wang, Minghan Li, Haoran Liu, Songlin Wei, Zhongyuan Wang, Zhizheng Zhang, and He Wang. Uni-navid: A video-based vision-language-action model for unifying embodied navigation tasks. *arXiv preprint arXiv:2412.06224*, 2024a.

Jiazhao Zhang, Kunyu Wang, Rongtao Xu, Gengze Zhou, Yicong Hong, Xiaomeng Fang, Qi Wu, Zhizheng Zhang, and He Wang. Navid: Video-based vlm plans the next step for vision-and-language navigation. *arXiv preprint arXiv:2402.15852*, 2024b.

Xiaoyun Zhang, Jingqing Ruan, Xing Ma, Yawen Zhu, Haodong Zhao, Hao Li, Jiansong Chen, Ke Zeng, and Xunliang Cai. When to continue thinking: Adaptive thinking mode switching for efficient reasoning. *arXiv preprint arXiv:2505.15400*, 2025.

Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and Liwei Wang. Towards learning a generalist model for embodied navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13624–13634, 2024.

Gengze Zhou, Yicong Hong, Zun Wang, Xin Eric Wang, and Qi Wu. Navgpt-2: Unleashing navigational reasoning capability for large vision-language models. In *European Conference on Computer Vision*, pp. 260–278. Springer, 2024a.

Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 7641–7649, 2024b.

Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pp. 2165–2183. PMLR, 2023.

## A  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.
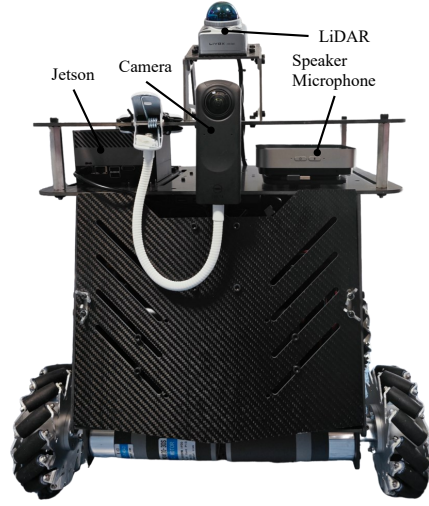
Building on this platform, we design a pipeline for vision-and-language navigation with AdaNav. We experiment with AdaNav on two base models, NAVID (Zhang et al., 2024b) and NAVILA (Cheng et al., 2024). In deployment, AdaNav runs on a server with a Jetson that receives compressed images from the robot and sends back high-level commands.



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

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.



You are in the hall, with a row of green plants on your right. Walk straight ahead until you reach the end of the plants, then turn right. After a few steps, you will see a row of tables. Walk to the end of the tables, then turn left. Continue straight, and after some distance you will see a row of bookshelves. Stop at the last bookshelf and that is your destination.

A wooden bar counter is in front, accompanied by high stools along its side, with flower arrangements displayed in the background.

I walk straight to the end of the plants, turn right, continue to the end of the tables and walk a few steps.
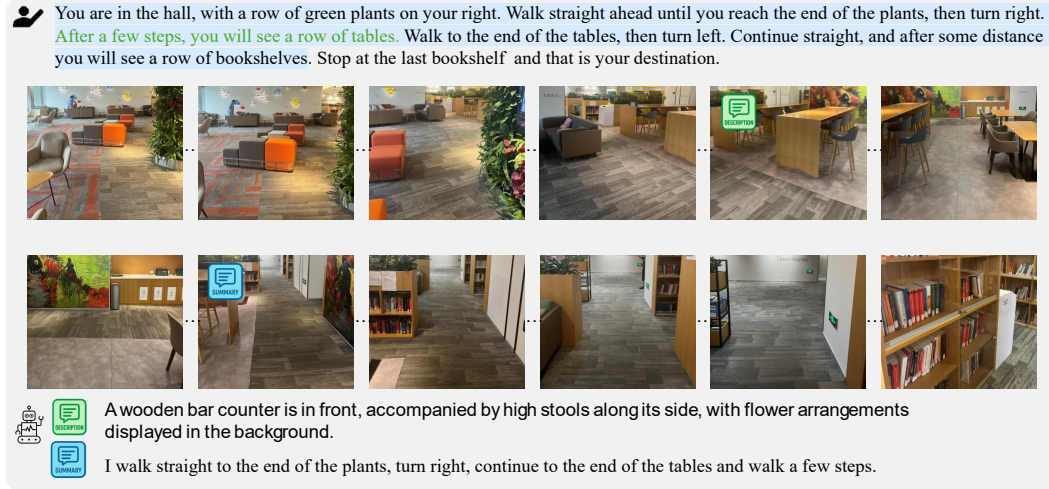
Figure 6: Visualization of adaptive reasoning navigation in real-world scene.

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

Table 8: Ablation results under different *reasoning modes*, where ①: Description only, ②: Summary only, ③: Error correction only, and ④: Generic prompt reasoning. All variants are trained with UAR Block and the heuristic-to-RL mechanism.

| 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 |

# C  VISUALIZATION RESULTS ON ADAPTIVE REASONING NAVIGATION

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

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

👤 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.