# OmniNav:A Unified Framework for Prospective Exploration and Visual-Language Navigation

Xinda Xue[1,2*]  Junjun Hu[1*†✉]  Minghua Luo[1*]  Shichao Xie [1]  Jintao Chen[1,2]
Zixun Xie[2]  Kuichen Quan[1]  Wei Guo[1]  Mu Xu[1]  Zedong Chu[1]

[1]Amap, Alibaba Group    [2]Peking University

{xuexinda.xxd, hujunjun.hjj, luominghua.lmh, tenan.xsc, anyi.cjt, quankuichen.qkc,
weisheng.gw, xumu.xm, chuzedong.czd}@alibaba-inc.com

*Equal Contribution, †Project Lead., ✉Corresponding authors

Embodied navigation is a foundational challenge for intelligent robots, demanding the ability to comprehend visual environments, follow natural language instructions, and explore autonomously. However, existing models struggle to provide a unified solution across heterogeneous navigation paradigms, often yielding low success rates and limited generalization. We present OmniNav, a unified framework that handles instruct-goal, object-goal, point-goal navigation, and frontier-based exploration within a single architecture. First, we introduce a lightweight, low-latency policy that predicts continuous-space waypoints (coordinates and orientations) with high accuracy, outperforming action-chunk methods in precision and supporting real-world deployment with control frequencies up to 5 Hz. Second, at the architectural level, OmniNav proposes a fast-slow system design: a fast module performs waypoint generation from relatively short-horizon visual context and subtasks, while a slow module conducts deliberative planning using long-horizon observations and candidate frontiers to select the next subgoal and subtask. This collaboration improves path efficiency and maintains trajectory coherence in exploration and memory-intensive settings. Notably, we find that the primary bottleneck lies not in navigation policy learning per se, but in robust understanding of general instructions and objects. To enhance generalization, we incorporate large-scale general-purpose training datasets including those used for image captioning and referring/grounding into a joint multi-task regimen, which substantially boosts success rates and robustness. Extensive experiments demonstrate state-of-the-art performance across diverse navigation benchmarks, and real-world deployment further validates the approach. OmniNav offers practical insights for embodied navigation and points to a scalable path toward versatile, highly generalizable robotic intelligence.

**Date:** Jan. 6, 2026
**Project Page:** https://astra-amap.github.io/omninav.github.io/
**Github:** https://github.com/amap-cvlab/OmniNav/

## 1   INTRODUCTION

Embodied navigation (Gao et al., 2024; Gu et al., 2022) has emerged as a core problem in embodied intelligence: enabling robots to perceive, understand, and explore real-world environments without pre-built maps while following natural language instructions. To act reliably in dynamic, partially observable environments, an agent must not only ground instantaneous visual inputs but also maintain coherent spatiotemporal memory and perform active exploration. Application demands for real-time responsiveness further increase the requirements for low-latency decision-making and cross-environment generalization.

Current research largely revolves around three paradigms: point-goal (Liu et al., 2025), instruct-goal (Anderson et al., 2018; Ku et al., 2020), and object-goal (Yokoyama et al., 2024b). point-goal tasks are well-specified and straightforward to evaluate but rely on explicit coordinates rarely available in practice; instruction-goal aligns with human usage but often generalizes poorly to unseen instructions or environments; object-goal is the most practical but requires robust target recognition coupled with efficient path planning, making it the most challenging. Many existing methods remain customized, relying on task-specific data, which limits cross-task transfer and the potential for mutual enhancement. Uni-Navid (Zhang et al., 2024a) proposed a VLM-based discrete action predictor unifying vision-and-language navigation,

object-goal navigation, embodied question answering (Das et al., 2018), and following (Wang et al., 2025a), but its study of LLM long-horizon planning is not sufficiently developed. MTU3D (Zhu et al., 2025) advances a "move to understand" paradigm by coupling frontier exploration with visual localization in a single objective, yet requires constructing 3D object coordinates, leading to deployment complexity. Although recent Video-LLMs (Wei et al., 2025; Qi et al., 2025; Zhang et al., 2024b) and VLAs (Sapkota et al., 2025; Zitkovich et al., 2023; Ma et al., 2024) integrate vision, language, and action prediction end to end, they still face bottlenecks in streaming video input, long-context management, and low-latency inference: discretized action modeling sacrifices precision and flexibility; constrained LLM call frequency and frequent context resets lead to deployment difficulties; besides, in practice, the dominant failure mode often stems from inadequate understanding of generic instructions and open-vocabulary objects rather than policy learning itself. These gaps call for a unified, efficient framework that balances long/short-horizon reasoning with real-time responsiveness.

We present OmniNav, a unified embodied navigation framework that concurrently covers instruct-goal, object-goal, point-goal, and frontier-based exploration within a single architecture. Inspired by dual-system theory (Figure, 2024; Black et al., 2025), OmniNav coordinates a fast–slow system (Black et al., 2025): a fast system reacts to comparatively short-horizon perception and current tasks or subtasks, generating high-precision waypoints (coordinates and orientations) to support low-latency control up to 5 Hz; a slow system deliberates over long-horizon observations and frontier cues, leveraging a VLM's chain-of-thought (Wei et al., 2022) to decompose complex goals and select the next subgoal and subtask. The two are coupled through a central memory module that uses a key–value (KV) cache to provide essential spatiotemporal context, yielding decisions that are both locally agile and globally consistent.

OmniNav addresses the triad of real-time operation, fast–slow collaboration, and generalization. A lightweight flow-matching policy (Bjorck et al., 2025) avoids the precision degradation and latency accumulation inherent to action discretization; fast–slow collaboration ensures exploration efficiency and trajectory coherence in long-memory scenarios; more importantly, training unifies large-scale generic vision–language data (captioning, referring/grounding, etc.) with multiple navigation tasks, significantly strengthening instruction following and open-vocabulary object perception to improve success rates and robustness. Our contributions are threefold:

- A unified architecture that, under a single training framework and policy, supports multiple goal modalities (point, object, and instruction) as well as frontier-based exploration;

- An end-to-end fast–slow coordination with central memory that reconciles low-latency control and high-level deliberation;

- A principled strategy to incorporate generic vision–language data into joint training, systematically improving cross-task and cross-environment generalization.

Extensive experiments set new state-of-the-art results across multiple navigation benchmarks, with real-robot deployments further validating practicality. We contend that OmniNav charts a scalable path toward multifunctional, highly generalizable embodied navigation systems.
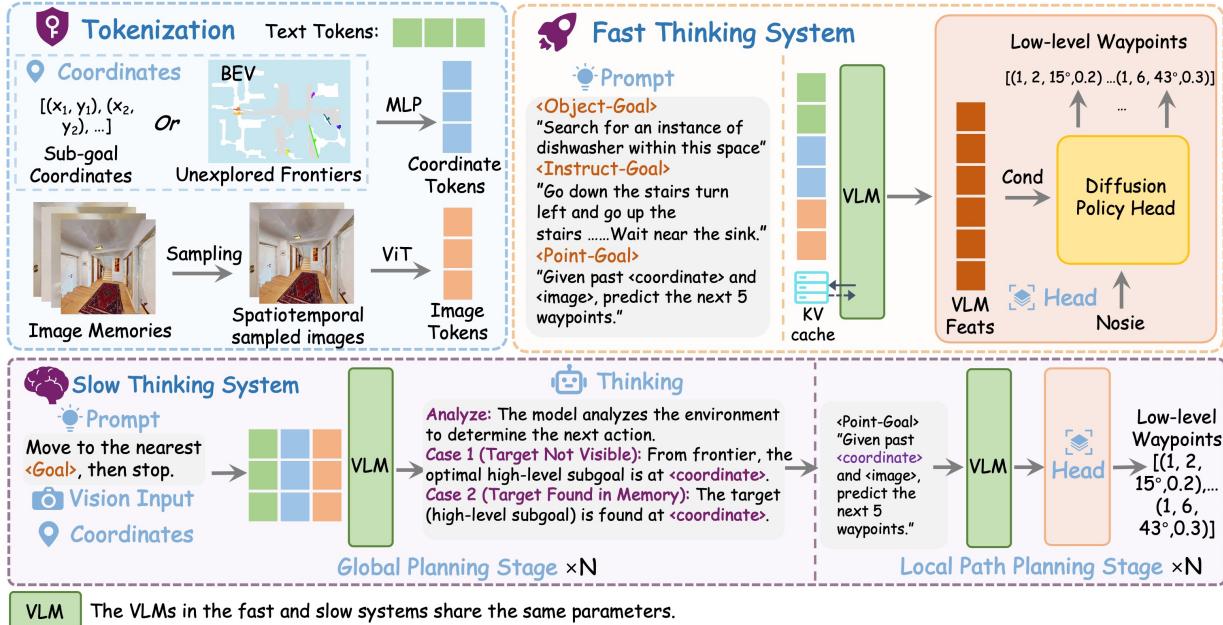
## 2 RELATED WORKS

**Vision Language Models for Navigation** Leveraging their powerful generalization capabilities in understanding and planning, Visual Language Models (VLMs) (Chiang et al., 2023; Liu et al., 2023; Zhu et al., 2023) have been increasingly applied to the domain of robotic navigation, achieving notable success. Prevailing methods (Dorbala et al., 2022; Zhou et al., 2024b; Long et al., 2024b) typically employ VLMs to process multimodal instructions and directly decode low-level actions in an autoregressive manner. However, this paradigm suffers from significant drawbacks: it is prone to compounding errors in sequential prediction and is often hampered by slow inference speeds. In contrast, our approach draws inspiration from recent advances in Vision-Language-Action (VLA) models (Zitkovich et al., 2023; Kim et al., 2024; Li et al., 2024). We introduce a novel architecture that appends a flow-matching policy (Zhao et al., 2024; Chen et al., 2024; Zhou et al., 2024a) to a VLM backbone. This design enables our model to generate entire action trajectories non-autoregressively, leading to substantially improved prediction accuracy and computational efficiency, especially when navigating unseen environments.

**Dual-System Design** Dual-system architectures have been widely adopted across various domains to meet diverse operational demands. In the realm of Vision-Language-Action (VLA) models, several works (Bjorck et al., 2025; Bu et al., 2025; Ge et al., 2024; Song et al., 2025) have implemented dual-system designs to balance fast control execution

and intelligent planning. Inspired by this paradigm and motivated by the specific requirements of embodied navigation, we propose a novel dual-system framework. Our framework consists of two complementary components. The first is a fast system, a purely visual, end-to-end policy designed for direct deployment and is highly effective in the majority of navigation scenarios. The second system is specifically engineered for challenging long-horizon tasks. To serve as a long-term memory mechanism, we employ a planning strategy that combines frontier-based exploration (Zhu et al., 2025) with images. This approach offers a more concise implementation compared to alternative memory structures such as scene graphs (Team et al., 2025) or complex semantic maps (Long et al., 2024a). These alternatives memory structures can also be implementations for the slow system. The idea stays the same: the slow system is responsible for global planning, while the fast system handles local execution. This synergistic design has proven its superiority by achieving State-of-the-Art performance on multiple benchmarks.

**Frontier-based Navigation** Recent studies on exploration and navigation adopt different strategies for selecting informative targets in unknown environments. GOAT (Chang et al., 2023) and its benchmark GOAT-Bench (Khanna et al., 2024) study lifelong navigation and object search using an object-instance memory and frontier exploration. Similarly, MTU3D (Zhu et al., 2025) keep an object-goal memory built from 3D point clouds and semantic segmentation, and combine this with frontier exploration. A different group of methods uses non-semantic frontier exploration (Chang et al., 2023; Sakamoto et al., 2024; Nayak et al., 2025), where the next target is usually just the closest frontier, sometimes adjusted by simple heuristics such as distance–heading scores. OmniNav instead uses a semantics- and reasoning-aware frontier selection: it links each frontier to its egocentric images, then uses explicit chain-of-thought reasoning over these views to decide which frontier is more informative or promising for the current task.

# 3   Approach



**Figure 1** The fast system can independently handle multi-task navigation, using the VLM backbone and a flow-matching policy to rapidly generate waypoints. Building on this, a slow thinking module is integrated to enable long-term memory and planning: it constructs long-range spatial and semantic memory using frontiers and images, and provides subgoal cues. The collaboration between the slow and fast proceeds as follows: the slow system uses frontiers or memory to generate high-level subgoals, once a subgoal is determined, the fast system takes over and progressively produces low-level waypoint sequences, ultimately reaching the target.

**Multimodal Input tokenizations** To handle all four task types through a unified interface, text, coordinates, and visual history are converted into a set of discrete tokens consumable by a Large-Language Model (LLM) see Fig. 1. We use Qwen2.5-VL-3B-Instruct (Bai et al., 2025) as the base model and extend it with a coordinate modality. During streaming inference, a key–value (KV) cache is maintained to reduce latency. Text tokens: Derived from natural-language task descriptions, object category labels, and point-goal commands, are all converted into a standardized instruction sequence.

Coordinate tokens: Candidate search regions are represented as sets of 2D coordinates and heading angles, sourced from point-goal inputs or subgoal positions generated by the slow system. These coordinates are processed via an MLP to dense embeddings that serve as coordinate tokens. Image tokens: The central memory maintains a ring buffer of pose-stamped images. For the fast system, it spatiotemporally samples from the historical image sequence, maintaining a maximum number of images (e.g., 20 frames). For the slow system, it samples images from the spatiotemporal neighborhood of candidate frontiers. All images are encoded with a ViT to produce image tokens.

**Fast Thinking System** OmniNav operates at a high frequency, designed to execute either subtasks provided by a slow system or end-to-end multi-task navigation, as shown in Fig. 1. It parallelly outputs a sequence of 5 continuous-space waypoints, $\mathbf{w}_{t:t+H} \in \mathbb{R}^{H \times 5}$ with $H = 5$. We formulate waypoint prediction as a conditional diffusion generation task. The input coordinates are first embedded by an MLP and then encoded together with the images and texts by the VLM. The VLM performs deep fusion over these features, and the resulting fused features are used as conditions for the diffusion model, guiding waypoint generation. This design preserves the VLM's semantic understanding while enabling rich interactions between language-guided context and the waypoint, leading to robust instruction following. Compared with conventional autoregressive methods, the policy head achieves an speedup and, with a history of 20 frames, supports an inference rate of 5 Hz for real-time closed-loop control. At the same time, it produces smoother and more precise trajectories.

We employ a variant of the Denoising Transformer (DiT) (Peebles & Xie, 2023) to model waypoint sequences. The policy network consists of self-attention blocks that operate on noised tokens to capture temporal and spatial dependencies within the waypoint sequence, and cross-attention blocks that attend to the vision-language context $\mathbf{O}_{VLM}$. The output is a sequence of $H = 5$ spatial-temporal waypoints $\mathbf{w}_t^{(i)} \in \mathbb{R}^5$, $i = 1, \ldots, H$, each encoding:

$$\mathbf{w}_t^{(i)} = \left( x^{(i)}, y^{(i)}, \sin \theta^{(i)}, \cos \theta^{(i)}, c^{(i)} \right), \tag{1}$$

where $(x^{(i)}, y^{(i)})$ denotes the 2D position, $\theta^{(i)}$ is the orientation (represented via sine-cosine embedding to avoid discontinuity at $\pi/ - \pi$), and $c^{(i)} \in \{0, 1\}$ is a binary completion flag indicating whether the "arrive" command should be triggered at the $i$-th waypoint.

Conditional flow matching policy is employed (Lipman et al., 2022). Given a ground-truth waypoint sequence $\mathbf{w}_{t:t+H}$, noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and a time parameter $\tau \in [0, 1]$, the input is constructed as:

$$\mathbf{w}_{t:t+H}^{\tau} = \tau \mathbf{w}_{t:t+H} + (1 - \tau)\epsilon, \tag{2}$$

and the policy $\pi$ is trained to estimate the denoising residual $\epsilon - \mathbf{w}_{t:t+H}$ by minimizing:

$$\mathbb{E}_{\tau, \epsilon} \left[ \left\| \pi(\mathbf{O}_{VLM}, \mathbf{w}_{t:t+H}^{\tau}) - (\epsilon - \mathbf{w}_{t:t+H}) \right\|^2 \right]. \tag{3}$$

At inference, waypoints are generated via $S = 5$ steps of Euler integration. Starting from initial noise $\mathbf{w}_{t:t+H}^0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we iteratively refine the sequence:

$$\mathbf{w}_{t:t+H}^{\tau + \Delta\tau} = \mathbf{w}_{t:t+H}^{\tau} + \frac{1}{S}\pi(\mathbf{O}_{VLM}, \mathbf{w}_{t:t+H}^{\tau}), \quad \Delta\tau = \frac{1}{S}, \tag{4}$$

with $\tau$ increasing from 0 to 1. The final denoised output $\mathbf{w}_{t:t+H}^1$ serves as the predicted waypoints.
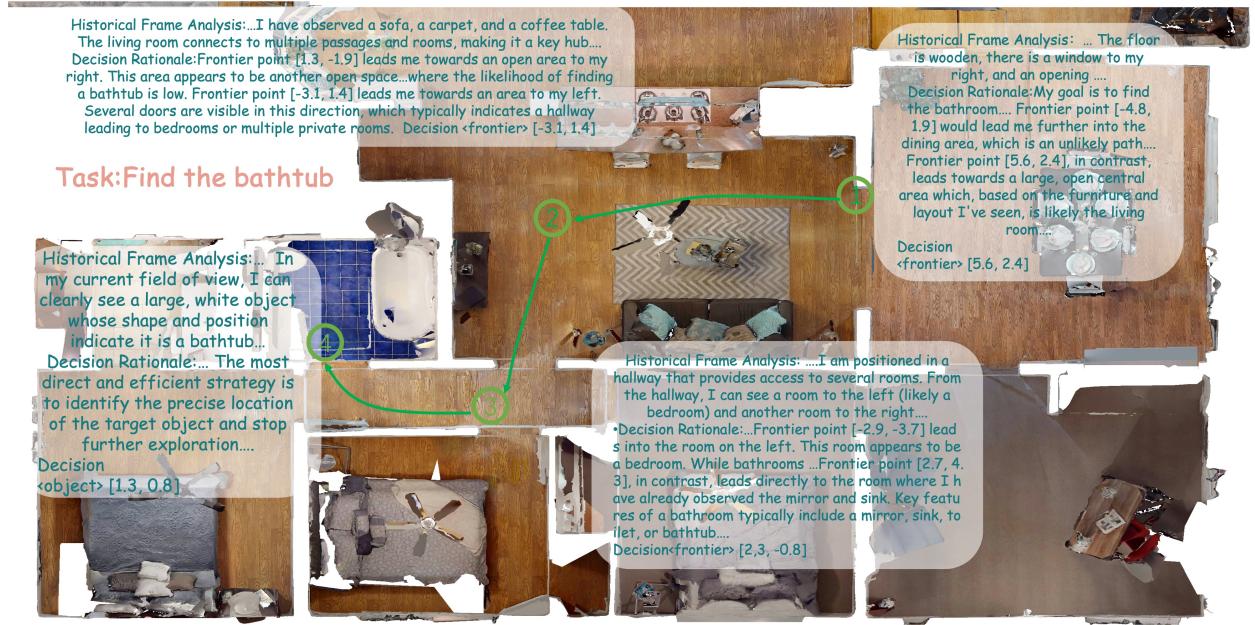
**Slow Thinking System** The slow thinking system is the deliberative planning module for hierarchical active exploration. Its core responsibilities are twofold: when the target appears in the current or historical field of view, it can quickly localize the target and generate the coordinates of the subgoal that drive the fast system to progressively approach it, when the target is not observed, it selects subgoal position with strong semantic relevance to the target to explore next. This process demands both exploration and environmental understanding. The former requires the model to navigate and discover the environment, while the latter involves a core task: predicting the spatial coordinates of the target based on input camera poses and images, this grounds the model's semantic understanding in a concrete, geometric output.

In the fast–slow system collaboration, the fast system is not just a low-level controller following preset coordinates or smoothing a pre-planned path. It must constantly use raw visual input to move toward its subgoal. For example, if it gets a coordinate from the slow system but a wall blocks the straight-line path, the fast system must use visual cues to find a path that avoids the obstacle. If it gets a target coordinate from object memory, it can adjust its final pose based on

what it currently sees, such as stopping precisely at the left, right, or center of the object. By continuously updating and refining its waypoints while moving, the fast system can reach targets more accurately.

Frontier (Zhu et al., 2025) is employed to guide the active exploration. We maintain a 3D occupancy map, which categorizes each region as explored or unknown, and frontiers are then identified as the boundary points between explored and unknown regions. In addition, to comprehend past temporal and spatial information, we construct a memory bank (Zhu et al., 2025; Olton, 1977; Xu et al., 2024; Zhou et al., 2023). This repository archives a history of observations, storing the visual data and corresponding pose information (coordinates and orientations) after every executed action. We then design a sampling strategy that connects this historical context to future exploration by collecting all historical images captured near the agent's current location. It then evaluates each frontier by iterating through these images, sampling the one whose original capture viewpoint is most suitably aligned with the frontier's spatial coordinates. This image thereby becomes a visual proxy for that frontier.
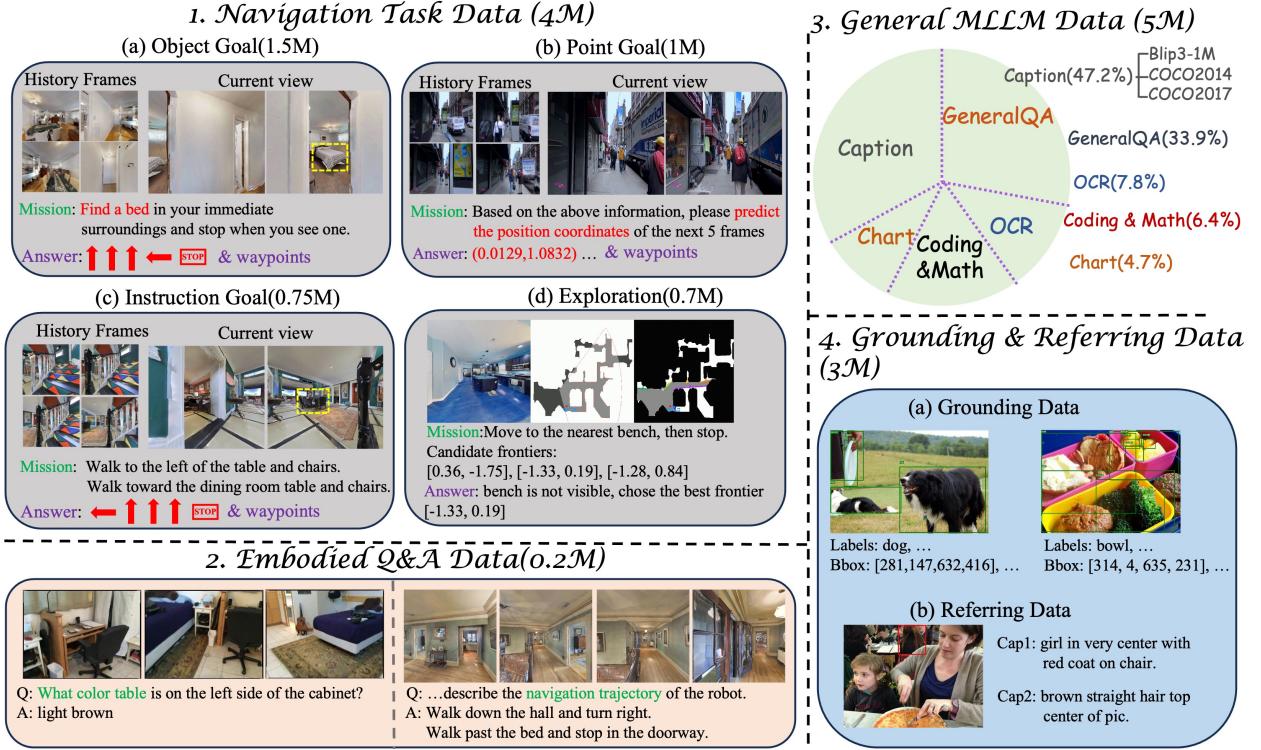
During frontier selection, the model engages in comprehensive spatial and content reasoning, reflecting its capability to actively explore unknown environments and propose subgoal locations related to the target object. For example, when searching for a toilet, it prioritizes exploration locations in the bathroom; when searching for a television, it seeks locations associated with the living room. Moreover, if the target object appears in memory or within the current view, it outputs the target's existing location. We incorporate explicit chain-of-thought (CoT) (Lin et al., 2025; Wang et al., 2025b; Zhao et al., 2025) reasoning into the slow system's prediction process to enable transparent process expression, achieving interpretability and self-correction. This also allows for richer textual outputs that strengthen the model's grasp of logic and improve complex reasoning performance. Fig.2 illustrates the slow system's reasoning process.



**Figure 2** Reasoning process by the slow system for exploration. For the "find the bathtub" task, the model reasons over the frontier set using memory and semantic priors, iteratively generating subgoals for the next exploration.

# 4 DATA and Training

Data in the embodied domain is typically organized as a data pyramid (Bjorck et al., 2025), with internet data and human video data at the bottom, simulation or synthetic data in the middle, and real-robot data at the top. Our dataset follows this rich composition as well, including general web data, simulation data, and a very small amount of real-robot data detailed in Fig. 3

**Figure 3** Data composition overview. Four data types are used for training: Navigation task data, Embodied Q&A data, General MLLM data and Grounding and referring data.

## 4.1 General Dataset

It is found in our experiment that models can learn the navigation paradigm relatively easily, whereas general-purpose capabilities remain challenging. To strengthen these abilities, we extended the training data with broad-coverage general-purpose datasets. General-purpose datasets (general QA, image captioning, OCR, chart understanding, coding, and math) complement Vision-and-Language Navigation (VLN) by supplying foundational skills in language understanding, visual semantics, OCR recognition, structured reasoning, and algorithmic planning. These capabilities improve instruction comprehension, path planning, and overall robustness. They also introduce commonsense and functional priors, such as "bath towels are commonly found in bathrooms". We draw on the open-source MAmmoTH-VL corpus (Guo et al., 2024), and subsampled examples to expand our general capabilities; the composition ratios are shown in Fig. 3.

We further incorporated grounding and referring data to more reliably ground linguistic targets and relations to concrete pixels, instances, and locations in the scene, thereby yielding policies that are more interpretable, robust, and generalizable. This includes fine-grained language-to-target mapping—for example, precisely localizing language like "red sofa," "door with a handle," or "the second chair" to the correct image regions to support instance-level disambiguation. It also encompasses spatial and relational understanding, such as learning to ground spatial prepositions and relations ("to the left/in front of/at the corner"), and extracting actionable turning and stopping cues from relational descriptions like "the door next to the painting" or "after passing the hallway, turn right at the first intersection." The referring and grounding data are sourced from RefCOCO series (Chen et al., 2025b) and Objects365 (Shao et al., 2019).

## 4.2 Multi-Task Navigation Dataset

**Point-goal data**. We use the open-source Citywalker corpus (Liu et al., 2025), which consists of first-person city-walking videos collected from public YouTube channels. Each trajectory comprises a single forward-view sequence. Long videos are segmented into 2-minute clips, then the camera pose for every frame are recovered by DPVO (Deep Patch Visual Odometry) (Teed et al., 2023).

**Instruct-goal data**. Under Habitat's VLN-CE (continuous environment) setting and Matterport3D scens (Chang et al.,

6

2017), based on the public instruction–path pairs from R2R (Anderson et al., 2018) and RxR (Ku et al., 2020). For each trajectory, we store a panoramic sequence composed of three-views (front, left, right), the action and waypoint sequences, and the natural-language instruction.

**Object-goal data**. Specifically, open-vocabulary object navigation (OVON) data (Yokoyama et al., 2024b). In Habitat-Matterport3D (HM3D) (Ramakrishnan et al., 2021) scenes, we randomly sample (start pose, target object category) pairs. The built-in shortest-path navigation algorithm is used to move the agent to the vicinity of the target. The same content are recorded as the instruct-goal data for each trajectory.

**Object-goal (with frontier-based exploration) data**. At each step the agent updates its occupancy map and visible region to identify all current frontiers. A policy then selects one frontier—favoring the shortest path while introducing limited randomness—as the current subgoal. Once the target object is found, the exploration episode terminates successfully, yielding a trajectory (Zhu et al., 2025). Each trajectory record includes a single forward-view frame sequence, a unexplored frontiers sequence and a natural-language instruction.

**Embodied QA**. ScanQA (Azuma et al., 2022) focuses on real-world indoor 3D scene understanding, with QA pairs centered on object locations, attributes, and spatial relations. R2R-EnvDrop (Tan et al., 2019) addresses continuous visual navigation: scenes are from Matterport3D, trajectories are from R2R-CE, and the navigation instruction–trajectory pairs are recast into a QA format to strengthen alignment between linguistic expressions and visual observations.

**Navigation data Process**.Each action step corresponds to a 3D continuous pose, described jointly by position and orientation. Taking the agent's current front-view coordinate frame as the origin, all other trajectory points are transformed into this local frame. Then, we project each 3D pose onto the ground plane, keeping only the planar coordinates $(x, y)$ and the heading angle $(\theta)$, and we additionally attach an "arrival" flag to each step. In this way, each trajectory point is ultimately represented as a quadruple $(x, y, \theta, arrive)$, which serves as the primary target for waypoint optimization. Under this abstraction, since the underlying simulator ensure that any valid 2D path passing through a stair region can be realized as a 3D trajectory across floors, this representation naturally supports multi-floor navigation in HM3D: moving from one floor to another is encoded as following a sequence of 2D poses from the stair entrance to the stair exit, while the vertical motion is handled implicitly by the simulator. At the same time, the agent always receives 3D visual observations with full geometric information, so the stair structure and floor changes are reflected in the visual features.

## 4.3 Discrete and Continuous Joint Training

We adopt a two-stage training paradigm to balance language–vision semantics and continuous control. In Stage 1, we use an autoregressive (AR) objective to predict discrete variables (e.g., navigation action chunks, general-purpose semantic data, Embodied QA, grounding and referring data; see the four data types in Fig. 3), to achieve alignment between language–vision and action. In Stage 2, we attach a flow-matching policy to the shared backbone to predict continuous waypoints, and perform joint training by including 20% of the Stage-1 discrete data to prevent degradation of the base VLM during continuous-control fine-tuning. The continuous waypoint coordinates are normalized using min-max normalization to ensure stable training and better convergence.

A joint training scheme is particularly critical for Vision-and-Language Navigation (VLN), which requires strong general knowledge and semantic understanding, leading to substantial improvements in success rates in open environments. Stage 1 is trained with 96 NVIDIA H20 GPUs for 120 hours, and Stage 2 is trained with 64 NVIDIA H20 GPUs for 48 hours with lower learning rate.

## 5 Experiments

**Metrics** We evaluate navigation performance using success rate (SR), oracle success rate (OS), success weighted by path length (SPL), and navigation error (NE). Our evaluation protocol is consistent with prior work (Zhang et al., 2024a; Zhu et al., 2025) and follows standard practice.

**Instruct goal** As shown in Table 1, On the R2R-CE and RxR-CE benchmarks, we compare our model against all relevant competitors, including both discrete and continuous prediction methods. Notably, using only its fast system and pure RGB inputs, OmniNav achieves state-of-the-art success rates on both benchmarks. It surpasses the previous leading model by improving the success rate by 4.4% on R2R-CE and 4.3% on RxR-CE.

**Table 1** Main comparison with prior methods on the Val-Unseen split of R2R-CE and RxR-CE.

| Method | Observation | | | | R2R-CE Val-Unseen | | | | RxR-CE Val-Unseen | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | S.RGB | Pano. | Depth | Odo. | NE↓ | OS↑ | SR↑ | SPL↑ | NE↓ | SR↑ | SPL↑ |
| HPN+DN* (Krantz et al., 2021) | | ✓ | ✓ | ✓ | 6.31 | 40.0 | 36.0 | 34.0 | - | - | - |
| CMA* (Hong et al., 2022) | | ✓ | ✓ | ✓ | 6.20 | 52.0 | 41.0 | 36.0 | 8.76 | 26.5 | 22.1 |
| Sim2Sim* (Krantz & Lee, 2022) | | ✓ | ✓ | ✓ | 6.07 | 52.0 | 43.0 | 36.0 | - | - | - |
| GridMM* (Wang et al., 2023b) | | ✓ | ✓ | ✓ | 5.11 | 61.0 | 49.0 | 41.0 | - | - | - |
| DreamWalker* (Wang et al., 2023a) | | ✓ | ✓ | ✓ | 5.53 | 59.0 | 49.0 | 44.0 | - | - | - |
| Reborn* (An et al., 2022) | | ✓ | ✓ | ✓ | 5.40 | 57.0 | 50.0 | 46.0 | 5.98 | 48.6 | 42.0 |
| ETPNav* (An et al., 2024) | | ✓ | ✓ | ✓ | 4.71 | 65.0 | 57.0 | 49.0 | 5.64 | 54.7 | 44.8 |
| HNR* (Wang et al., 2024) | | ✓ | ✓ | ✓ | 4.42 | 67.0 | 61.0 | 51.0 | 5.50 | 56.3 | 46.7 |
| AG-CMTP (Chen et al., 2021) | | ✓ | ✓ | ✓ | 7.90 | 39.0 | 23.0 | 19.0 | - | - | - |
| R2R-CMTP (Chen et al., 2021) | | ✓ | ✓ | ✓ | 7.90 | 38.0 | 26.0 | 22.0 | - | - | - |
| InstructNav (Long et al., 2024a) | | ✓ | ✓ | ✓ | 6.89 | - | 31.0 | 24.0 | - | - | - |
| LAW (Raychaudhuri et al., 2021) | ✓ | | ✓ | ✓ | 6.83 | 44.0 | 35.0 | 31.0 | 10.90 | 8.0 | 8.0 |
| CM2 (Georgakis et al., 2022) | ✓ | | ✓ | ✓ | 7.02 | 41.0 | 34.0 | 27.0 | - | - | - |
| WS-MGMap (Chen et al., 2022) | ✓ | | ✓ | ✓ | 6.28 | 47.0 | 38.0 | 34.0 | - | - | - |
| AO-Planner (Chen et al., 2025a) | | ✓ | ✓ | | 5.55 | 59.0 | 47.0 | 33.0 | 7.06 | 43.3 | 30.5 |
| Seq2Seq (Krantz et al., 2020) | ✓ | | ✓ | | 7.77 | 37.0 | 25.0 | 22.0 | 12.10 | 13.9 | 11.9 |
| CMA (Krantz et al., 2020) | ✓ | | ✓ | | 7.37 | 40.0 | 32.0 | 30.0 | - | - | - |
| NaVid (Zhang et al., 2024b) | ✓ | | | | 5.47 | 49.0 | 37.0 | 35.0 | - | - | - |
| Uni-NaVid (Zhang et al., 2024a) | ✓ | | | | 5.58 | 53.5 | 47.0 | 42.7 | 6.24 | 48.7 | 40.9 |
| NaVILA (Cheng et al., 2024) | ✓ | | | | 5.22 | 62.5 | 54.0 | 49.0 | 6.77 | 49.3 | 44.0 |
| StreamVLN (Wei et al., 2025) | ✓ | | | | 4.98 | 64.2 | 56.9 | 51.9 | 6.22 | 52.9 | 46.0 |
| CorrectNav (Yu et al., 2025) | ✓ | | | | 4.24 | 67.5 | 65.1 | 62.3 | 4.09 | 69.3 | 63.3 |
| OmniNav(w/o policy-head) | | ✓ | | | 4.36 | 65.0 | 59.8 | 57.5 | 3.87 | 64.1 | 53.9 |
| **OmniNav** | | ✓ | | | **3.74** | **74.6** | **69.5** | **66.1** | **3.77** | **73.6** | **62.0** |

**Object goal** To further validate open vocabulary generalization, we also compare OmniNav with prior methods on the HM3D-OVON benchmark. As shown in Table 2, under purely visual inputs, OmniNav already surpasses the best existing approach by 2.7%. However, given that the OVON task demands long-horizon and global planning, the limitations of a purely reactive fast system—such as getting trapped in local loops and exhibiting poor map coverage—become particularly pronounced. Therefore, we incorporate the slow system integrated with frontier-based reasoning, which necessitates the use of depth and odometry information to build and maintain an occupation map. This augmentation equips the agent with global spatial awareness and the capacity for proactive exploration, ultimately leading to superior overall performance(exceeding the strongest prior method by 18.4%).

**Point goal** We compare point-goal performance on the CityWalker benchmark, a state-of-the-art point-goal method that supports outdoor navigation scenarios. CityWalker adopts MAOE (Mean Average Orientation Error) as an open-set evaluation metric. On this open-set metric, our approach still outperforms the benchmark method (OmniNav:11.53% vs CityWalker:15.23%).

**Ablation Study** The independent and synergistic contributions of the four key components policies are evaluated as show in Table.1 and Table.3. 1) Fast system: autoregressively action chunks generation vs. continuous waypoints generation by flow-matching policy. On R2R-CE, RxR-CE, and OVON benchmarks, the degradation is substantial for action chunks. Their semantic tokens (e.g., left, right) align more easily with language, making them suitable for the first-stage training. However, because action chunks are coarse-grained motion control, continuous waypoints are better suited for fine-grained control. 2) Slow system (planning with frontier and long-term visual memory). We primarily compare with/without the slow system on long-horizon exploration in OVON and find the largest improvements here. Once previously explored areas are recorded, the agent reduces redundant exploration and improves efficiency. Moreover, decomposing active exploration into subgoals (e.g., "go to the bedroom first") and letting the fast system quickly approach each subgoal forms a hierarchical "plan–execute" loop, which better matches human reasoning and behavior in unfamiliar environments. 3) General data (general MLLM and referring/grounding). With the slow system enabled, adding the general-purpose datasets yields further stable gains. 4) CoT (explicit chain-of-thought outputs). Using CoT makes the basis for subgoal selection in the slow system transparent, enabling process-level self-check and correction. It reduces cumulative errors in long chains and complex semantic tasks, producing stable improvements.

**Table 2** Evaluation of object-goal navigation on HM3D-OVON, where * indicates OmniNav with slow thinking system.

| Method | Observation | | | Val-Seen | | Val-Seen-Synonyms | | Val-Unseen | |
|---|---|---|---|---|---|---|---|---|---|
| | S.RGB | Depth | Odo. | SR↑ | SPL↑ | SR↑ | SPL↑ | SR↑ | SPL↑ |
| BC | ✓ | | | 11.1 | 4.5 | 9.9 | 3.8 | 5.4 | 1.9 |
| DAgger | ✓ | | | 11.1 | 4.5 | 9.9 | 3.8 | 5.4 | 1.9 |
| RL | ✓ | | | 18.1 | 9.4 | 15.0 | 7.4 | 10.2 | 4.7 |
| DAgRL | ✓ | | | 41.3 | 21.2 | 29.4 | 14.4 | 18.3 | 7.9 |
| BCRL | ✓ | | | 39.2 | 18.7 | 27.8 | 11.7 | 18.6 | 7.5 |
| VLFM* (Yokoyama et al., 2024a) | ✓ | ✓ | ✓ | 35.2 | 18.6 | 32.4 | 17.3 | 35.2 | 19.6 |
| DAgRL+OD (Yokoyama et al., 2024b) | ✓ | ✓ | ✓ | 38.5 | 21.1 | 39.0 | 21.4 | 37.1 | 19.8 |
| Uni-NaVid* (Zhang et al., 2024a) | ✓ | | | 41.3 | 21.1 | 43.9 | 21.8 | 39.5 | 19.8 |
| MTU3D* (Zhu et al., 2025) | ✓ | ✓ | ✓ | 55.0 | 23.6 | 45.0 | 14.7 | 40.8 | 12.1 |
| OmniNav | ✓ | | | 46.6 | 23.3 | 50.4 | 28.5 | 43.5 | 27.3 |
| **OmniNav*(w/ cot)** | ✓ | ✓ | ✓ | **56.1** | **30.0** | **68.6** | **38.8** | **59.2** | **33.2** |

When all four are enabled, performance is best: the slow system provides semantically plausible long-horizon subgoals; the policy head executes with high precision and low latency; general data injects commonsense and language–vision alignment; and CoT provides auditable processes and self-correction.
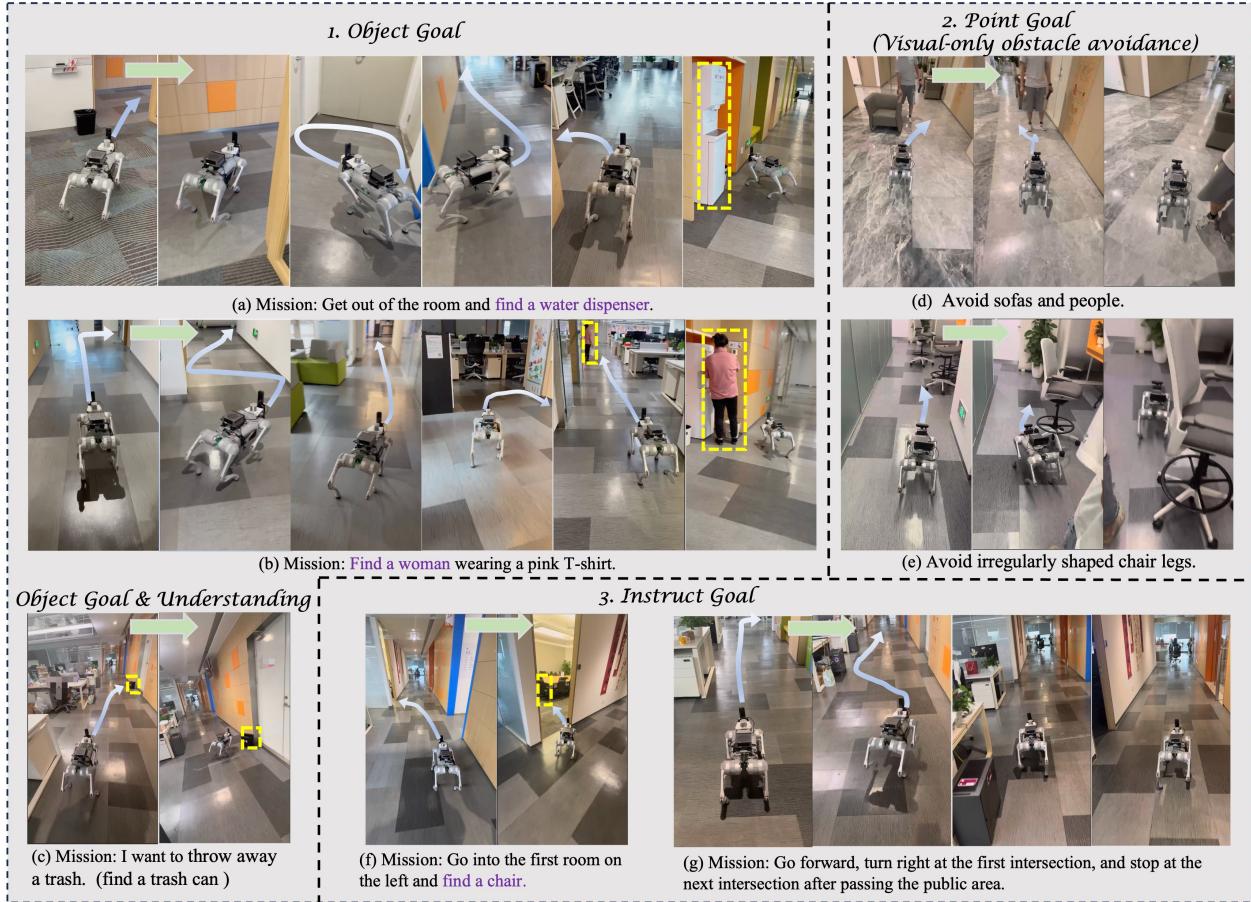
**Table 3** Ablation study on HM3D-OVON Val-Unseen

| Method | Module | | | | Val-Unseen | |
|---|---|---|---|---|---|---|
| | policy-head | slow-system | general data | COT | SR↑ | SPL↑ |
| OmniNav | | | | | 35.3 | 22.1 |
| | ✓ | | | | 43.5 | 27.3 |
| OmniNav* | ✓ | ✓ | | | 55.9 | 30.7 |
| | ✓ | ✓ | ✓ | | 57.7 | 32.9 |
| | ✓ | ✓ | ✓ | ✓ | **59.2** | **33.2** |

**Real-world deployment** In real-robot deployment see Fig. 4, we deploy the fast system component of the model architecture—comprising the VLM and the policy head on a cloud server with an RTX 3090 GPU. The history buffer holds up to 20 front-view frames with at a resolution of 120×106, while the current input is tri-view at 480×426. The system runs at over 5 Hz. For all tasks, it outputs waypoints, which are then fed into the onboard speed control module, which, based on the current speed and maximum acceleration, generates a set of candidate speeds. It then selects the optimal speed to approach the waypoint (selection criteria: minimize speed changes to maintain high speed as much as possible, and be closest to the target waypoint).

Deploying the full slow system in real-world settings requires additional engineering, such as robust real-time integration with LiDAR/depth estimation. This paper mainly focuses on validating the effectiveness of the dual-system collaboration framework in terms of navigation performance and behavior. A full physical deployment of the complete slow system, and systematic optimization of its performance under real-world constraints (including detailed latency–frequency trade-off analysis), constitutes an important avenue for future research.

# 6 CONCLUSION & FUTURE WORKS

**Conclusion** The core of OmniNav is a fast–slow dual-system architecture: the fast system, conditioned on VLM-fused multimodal context, employs a flow-matching policy to generate future continuous waypoints, achieving low latency, and high-precision closed-loop control; the slow system plans subgoals and subtasks supported by long-horizon visual memory and frontiers, and introduces explicit CoT for interpretability and self-correction. This architecture supports most basic tasks in embodied navigation. Through unified multimodal tokenization, different tasks (instruct goal, object goal, point goal) are seamlessly handled within a single model. On the training side, we adopt a two-stage scheme, where the second stage's joint training of discrete and continuous values prevents continuous-control fine-tuning

**Figure 4** Real-world deployment. It shows third-person view of the three different navigation tasks which are deployed in a zero-shot setting. The gradient blue arrows indicate the trajectory, and the yellow box marks the target location. Our model demonstrates highly effective navigation performance on the real quadruped robot.

from eroding the base VLM's capabilities—an approach that can be broadly applicable. We also incorporate sizable general-purpose data and referring/grounding data to bolster language understanding, visual semantics, structured reasoning, and commonsense priors for VLN, thereby improving generalization and robustness in embodied navigation. Experimentally, OmniNav improves success rates over the current best on R2R-CE and on RxR-CE, achieves the best performance on OVON, and benefits further from the slow-system design. Real-world quadruped robot deployment demonstrates the engineering feasibility of up to 5 Hz cloud inference with tri-view inputs and a 20-frame history buffer. Overall, the high spatial precision and low latency of continuous waypoints, the unified multimodal interface, the fast–slow system collaboration, and joint training collectively underpin OmniNav's strong performance across benchmarks, evidencing solid open-set generalization and practical deployment potential.

**Future Works** We aim to develop a more semantics-driven, learning-based subgoal selection strategy and realize purely visual memory capabilities—for example, remembering the semantic regions already explored. In addition, we plan to build a retrievable, lifelong spatiotemporal memory to enable lifelong navigation.

# References

Dong An, Zun Wang, Yangguang Li, Yi Wang, Yicong Hong, Yan Huang, Liang Wang, and Jing Shao. 1st place solutions for rxr-habitat vision-and-language navigation competition (cvpr 2022). *arXiv preprint arXiv:2206.11610*, 2022.

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.

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3674–3683, 2018.

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.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.

Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y Galliker, et al. $\pi0. 5$: a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.

Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.

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.

Matthew Chang, Theophile Gervet, Mukul Khanna, Sriram Yenamandra, Dhruv Shah, So Yeon Min, Kavit Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, et al. Goat: Go to any thing. *arXiv preprint arXiv:2311.06430*, 2023.

Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems*, 37:24081–24125, 2024.

Jiaqi Chen, Bingqian Lin, Xinmin Liu, Lin Ma, Xiaodan Liang, and Kwan-Yee K Wong. Affordances-oriented planning using foundation models for continuous vision-language navigation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 23568–23576, 2025a.

Jierun Chen, Fangyun Wei, Jinjing Zhao, Sizhe Song, Bohuai Wu, Zhuoxuan Peng, S-H Gary Chan, and Hongyang Zhang. Revisiting referring expression comprehension evaluation in the era of large multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 513–524, 2025b.

Kevin Chen, Junshen K Chen, Jo Chuang, Marynel Vázquez, and Silvio Savarese. Topological planning with transformers for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11276–11286, 2021.

Peihao Chen, Dongyu Ji, Kunyang Lin, Runhao Zeng, Thomas Li, Mingkui Tan, and Chuang Gan. Weakly-supervised multi-granularity map learning for vision-and-language navigation. *Advances in Neural Information Processing Systems*, 35:38149–38161, 2022.

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.

Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.

Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–10, 2018.

Vishnu Sashank Dorbala, Gunnar Sigurdsson, Robinson Piramuthu, Jesse Thomason, and Gaurav S Sukhatme. Clip-nav: Using clip for zero-shot vision-and-language navigation. *arXiv preprint arXiv:2211.16649*, 2022.

AI Figure. Helix: A vision-language-action model for generalist humanoid control. *Figure AI News*, 2024.

Peng Gao, Peng Wang, Feng Gao, Fei Wang, and Ruyue Yuan. Vision-language navigation with embodied intelligence: A survey. *arXiv preprint arXiv:2402.14304*, 2024.

Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.

Georgios Georgakis, Karl Schmeckpeper, Karan Wanchoo, Soham Dan, Eleni Miltsakaki, Dan Roth, and Kostas Daniilidis. Cross-modal map learning for vision and language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15460–15470, 2022.

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.

Jarvis Guo, Tuney Zheng, Yuelin Bai, Bo Li, Yubo Wang, King Zhu, Yizhi Li, Graham Neubig, Wenhu Chen, and Xiang Yue. Mammoth-vl: Eliciting multimodal reasoning with instruction tuning at scale. *arXiv preprint arXiv:2412.05237*, 2024.

Yicong Hong, Zun Wang, Qi Wu, and Stephen Gould. Bridging the gap between learning in discrete and continuous environments for vision-and-language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15439–15449, 2022.

Mukul Khanna, Ram Ramrakhya, Gunjan Chhablani, Sriram Yenamandra, Theophile Gervet, Matthew Chang, Zsolt Kira, Devendra Singh Chaplot, Dhruv Batra, and Roozbeh Mottaghi. Goat-bench: A benchmark for multi-modal lifelong navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16373–16383, 2024.

Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.

Jacob Krantz and Stefan Lee. Sim-2-sim transfer for vision-and-language navigation in continuous environments. In *European conference on computer vision*, pp. 588–603. Springer, 2022.

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.

Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18061–18070, 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*, 2025.

Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.

Xinhao Liu, Jintong Li, Yicheng Jiang, Niranjan Sujay, Zhicheng Yang, Juexiao Zhang, John Abanes, Jing Zhang, and Chen Feng. Citywalker: Learning embodied urban navigation from web-scale videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 6875–6885, 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.

Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. A survey on vision-language-action models for embodied ai. *arXiv preprint arXiv:2405.14093*, 2024.

Sharan Nayak, Grace Lim, Federico Rossi, Michael Otte, and Jean-Pierre de la Croix. Multi-robot exploration for the cadre mission. *Autonomous Robots*, 49(2):17, 2025.

David S Olton. Spatial memory. *Scientific American*, 236(6):82–99, 1977.

William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.

Zhangyang Qi, Zhixiong Zhang, Yizhou Yu, Jiaqi Wang, and Hengshuang Zhao. Vln-r1: Vision-language navigation via reinforcement fine-tuning. *arXiv preprint arXiv:2506.17221*, 2025.

Santhosh K Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alex Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, et al. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *arXiv preprint arXiv:2109.08238*, 2021.

Sonia Raychaudhuri, Saim Wani, Shivansh Patel, Unnat Jain, and Angel X Chang. Language-aligned waypoint (law) supervision for vision-and-language navigation in continuous environments. *arXiv preprint arXiv:2109.15207*, 2021.

Koya Sakamoto, Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Map-based modular approach for zero-shot embodied question answering. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10013–10019. IEEE, 2024.

Ranjan Sapkota, Yang Cao, Konstantinos I Roumeliotis, and Manoj Karkee. Vision-language-action models: Concepts, progress, applications and challenges. *arXiv preprint arXiv:2505.04769*, 2025.

Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8430–8439, 2019.

Haoming Song, Delin Qu, Yuanqi Yao, Qizhi Chen, Qi Lv, Yiwen Tang, Modi Shi, Guanghui Ren, Maoqing Yao, Bin Zhao, et al. Hume: Introducing system-2 thinking in visual-language-action model. *arXiv preprint arXiv:2505.21432*, 2025.

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

BAAI RoboBrain Team, Mingyu Cao, Huajie Tan, Yuheng Ji, Xiansheng Chen, Minglan Lin, Zhiyu Li, Zhou Cao, Pengwei Wang, Enshen Zhou, et al. Robobrain 2.0 technical report. *arXiv preprint arXiv:2507.02029*, 2025.

Zachary Teed, Lahav Lipson, and Jia Deng. Deep patch visual odometry. *Advances in Neural Information Processing Systems*, 36: 39033–39051, 2023.

Hanqing Wang, Wei Liang, Luc Van Gool, and Wenguan Wang. Dreamwalker: Mental planning for continuous vision-language navigation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10873–10883, 2023a.

Shaoan Wang, Jiazhao Zhang, Minghan Li, Jiahang Liu, Anqi Li, Kui Wu, Fangwei Zhong, Junzhi Yu, Zhizheng Zhang, and He Wang. Trackvla: Embodied visual tracking in the wild. *arXiv preprint arXiv:2505.23189*, 2025a.

Shuo Wang, Yongcai Wang, Wanting Li, Xudong Cai, Yucheng Wang, Maiyue Chen, Kaihui Wang, Zhizhong Su, Deying Li, and Zhaoxin Fan. Aux-think: Exploring reasoning strategies for data-efficient vision-language navigation. *arXiv preprint arXiv:2505.11886*, 2025b.

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, 2023b.

Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, Junjie Hu, Ming Jiang, and Shuqiang Jiang. Lookahead exploration with neural radiance representation for continuous vision-language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 13753–13762, 2024.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

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.

Xiuwei Xu, Huangxing Chen, Linqing Zhao, Ziwei Wang, Jie Zhou, and Jiwen Lu. Embodiedsam: Online segment any 3d thing in real time. *arXiv preprint arXiv:2408.11811*, 2024.

Naoki Yokoyama, Sehoon Ha, Dhruv Batra, Jiuguang Wang, and Bernadette Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 42–48. IEEE, 2024a.

Naoki Yokoyama, Ram Ramrakhya, Abhishek Das, Dhruv Batra, and Sehoon Ha. Hm3d-ovon: A dataset and benchmark for open-vocabulary object goal navigation. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5543–5550. IEEE, 2024b.

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.

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.

Chuyang Zhao, Yuxing Song, Wenhao Wang, Haocheng Feng, Errui Ding, Yifan Sun, Xinyan Xiao, and Jingdong Wang. Monoformer: One transformer for both diffusion and autoregression. *arXiv preprint arXiv:2409.16280*, 2024.

Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 1702–1713, 2025.

Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 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.

Hang Zhou, Junqing Yu, and Wei Yang. Dual memory units with uncertainty regulation for weakly supervised video anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 3769–3777, 2023.

Deyao Zhu, Jun Chen, Kilichbek Haydarov, Xiaoqian Shen, Wenxuan Zhang, and Mohamed Elhoseiny. Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions. *arXiv preprint arXiv:2303.06594*, 2023.

Ziyu Zhu, Xilin Wang, Yixuan Li, Zhuofan Zhang, Xiaojian Ma, Yixin Chen, Baoxiong Jia, Wei Liang, Qian Yu, Zhidong Deng, et al. Move to understand a 3d scene: Bridging visual grounding and exploration for efficient and versatile embodied navigation. *arXiv preprint arXiv:2507.04047*, 2025.

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  Appendix

## A.1  Analysis of Training Data and Model scale

**Analysis of Training Data** The three additional data components shown in Figure 3 are ablated individually. A qualitative analysis of specific examples reveals the following: 1) Removing Embodied Q&A or Grounding/referring noticeably degrades performance on small background objects (e.g., "picture", "flowerpot"), suggesting these data improve recognition of small objects. 2) Removing General MLLM data leads to failures on irregular objects (e.g., "handrail", "stair"), implying general vision–language data helps with such targets. These findings motivate constructing specific datasets to address failure modes. For example, the model still struggles with "carpet" and "clothes" where complex folds and textures are involved; we therefore can curate dedicated VQA-style data emphasizing these patterns to strengthen fine-grained visual reasoning on texture-heavy objects.

**Table 4**  Ablation Study of Data on HM3D-OVON Val-Unseen

| Embodied Q&A Data | Grounding and Referring Data | General MLLM Data | Ovon-Unseen |
|---|---|---|---|
|  |  |  | 55.9 |
|  | ✓ | ✓ | 56.5 |
| ✓ |  | ✓ | 56.7 |
| ✓ | ✓ |  | 57.0 |
| ✓ | ✓ | ✓ | 57.7 |

**Analysis of Model Scale** The influence of model size (Embodied Q&A, Grounding and Referring, and General MLLM data) is also analyzed. When including these additional data, the 3B and 7B models exhibit nearly identical navigation performance, indicating that once such data are incorporated, simply scaling up the model size brings little further improvement.

In contrast, in the absence of this additional data, a noticeable performance gap emerges, with the 7B model outperforming the 3B model. Our analysis suggests two key factors: 1) the performance of the 3B model appears to be constrained by data sufficiency rather than its inherent capacity. When provided with diverse and abundant data, its performance becomes comparable to that of the 7B model, indicating that model size itself is not the primary bottleneck in this data-rich scenario. 2) based on an analysis of failure cases, we find that further performance gains are limited by the intrinsic difficulty of the task: regardless of model size, the recognition of complex objects such as clothes and mirrors remains unstable.

Beyond these two ablation study, a more systematic study of scaling laws and optimal training configurations—e.g., how data quality, data composition, and more model size jointly affect performance—would also be highly valuable. Due to computational limits we have not yet conducted such a systems-level exploration, and we view this as an important direction for future work.

**Table 5**  Ablation Study of Model Size on HM3D-OVON Val-Unseen

| Model | Additional Data | Ovon-Unseen |
|---|---|---|
| 3B |  | 55.9 |
| 7B |  | 57.2 |
| 3B | ✓ | 57.7 |
| 7B | ✓ | 57.9 |