# UrbanNav: Learning Language-Guided Urban Navigation from Web-Scale Human Trajectories

**Yanghong Mei**[*1,5], **Yirong Yang**[*2], **Longteng Guo**[†1], **Qunbo Wang**[3], **Ming-Ming Yu**[2],
**Xingjian He**[1], **Wenjun Wu**[2,4], **Jing Liu**[1,5]

[1]Institute of Automation, Chinese Academy of Sciences
[2]Beihang University
[3]Beijing Jiaotong University
[4]Hangzhou International Innovation Institute
[5]School of Artificial Intelligence, University of Chinese Academy of Sciences
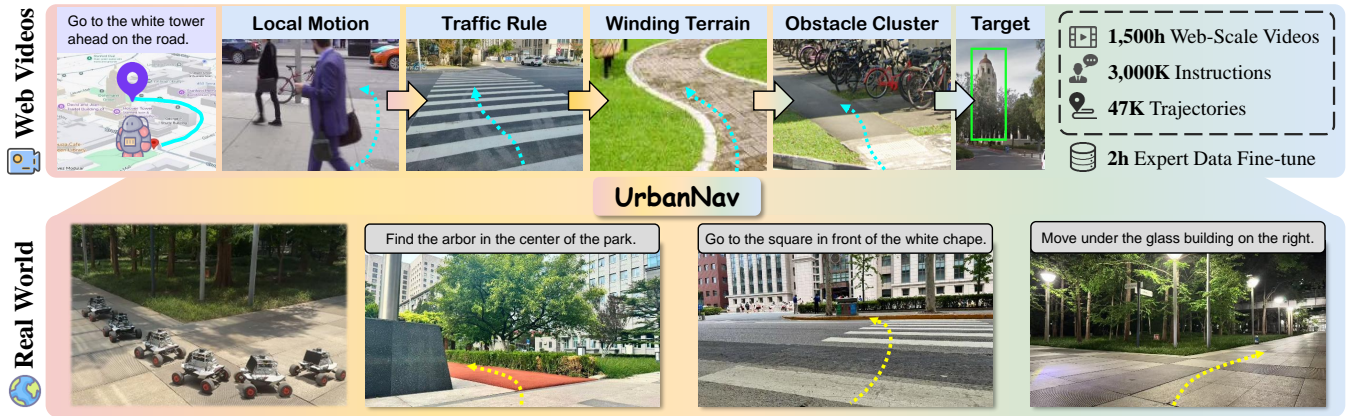meiyanghong2024@ia.ac.cn, yirongyang@buaa.edu.cn, longteng.guo@ia.ac.cn

Figure 1: **Overview of Our UrbanNav framework.** UrbanNav is designed to tackle the challenging task of language-guided urban navigation. Its scalable data pipeline constructs a large dataset from web-scale human walking videos. Our policy is trained on this dataset and fine-tuned with a small amount of real-world data, enabling it to interpret complex natural language instructions and navigate challenging, unseen urban environments.

## Abstract

Navigating complex urban environments using natural language instructions poses significant challenges for embodied agents, including noisy language instructions, ambiguous spatial references, diverse landmarks, and dynamic street scenes. Current visual navigation methods are typically limited to simulated or off-street environments, and often rely on precise goal formats, such as specific coordinates or images. This limits their effectiveness for autonomous agents like last-mile delivery robots navigating unfamiliar cities. To address these limitations, we introduce UrbanNav, a scalable framework that trains embodied agents to follow free-form language instructions in diverse urban settings. Leveraging web-scale city walking videos, we develop an scalable annotation pipeline that aligns human navigation trajectories with language instructions grounded in real-world landmarks. UrbanNav encompasses over 1,500 hours of navigation data and 3 million instruction-trajectory-landmark triplets, capturing a wide range of urban scenarios. Our model learns robust navigation policies to tackle complex urban scenarios, demonstrating superior spatial reasoning, robustness to noisy

instructions, and generalization to unseen urban settings. Experimental results show that UrbanNav significantly outperforms existing methods, highlighting the potential of large-scale web video data to enable language-guided, real-world urban navigation for embodied agents.

**Code** — https://github.com/Vigar0108M/UrbanNav

## 1 Introduction

Language-guided navigation in real-world urban environments is a cornerstone capability for autonomous agents, enabling applications such as last-mile delivery robots, autonomous vehicles, and assistive robotics. Unlike structured indoor spaces or synthetic simulation settings, urban scenes are inherently dynamic, featuring diverse terrains, unpredictable obstacles, and dense pedestrian interactions (Shah et al. 2023a). To operate effectively in such environments, embodied agents must not only reason about spatial layouts and adhere to implicit social norms but also interpret ambiguous human instructions, such as "go to the cafe by the old bridge" or "move to the bookstore opposite the park." These often vague or context-specific directives re-

---

[*]These authors contributed equally.
[†]Corresponding author (longteng.guo@ia.ac.cn)

quire sophisticated reasoning, making language-guided urban navigation a complex yet essential task for deploying autonomous agents in real-world cities (Gao et al. 2025).

Prior research (Shah et al. 2023b; Sridhar et al. 2024; Yu et al. 2025) on visual navigation has made substantial progress in simulation and indoor domains. Classical methods (Muhlbauer et al. 2009; Kümmerle et al. 2013) combine SLAM and modular planning to achieve goal-oriented navigation, while more recent works (Ehsani et al. 2024; Zeng et al. 2024) leverage reinforcement learning or imitation learning within high-fidelity simulators. These advances have enabled impressive results in point goal or object goal navigation tasks. However, current approaches are still constrained by their reliance on precise goal specifications, such as GPS coordinates or target images and limited diversity.

The complexity of language-guided urban navigation stems from the need to align natural language instructions with real-world spatial (Schumann and Riezler 2022). In real urban environments, users typically provide free-form directions such as "deliver to the building near the park fountain," which often reference salient visual landmarks. Human navigation often relies on such landmarks and contextual cues embedded in language, which are challenging to model in simulations or small-scale datasets. Although collecting expert trajectories via teleoperation is an adopted approach, in practice it is constrained by limited data diversity and high annotation costs, hindering generalizability across varied urban scenarios. To overcome these challenges, we propose **UrbanNav**, a scalable framework illustrated in Fig. 1, that leverages web-scale human navigation trajectories from city walking videos to train embodied agents for language-guided urban navigation. UrbanNav answers two key questions associated with learning from unstructured web videos:

*First, are all video segments suitable for training embodied agents?* Many clips exhibit viewpoint divergence, where camera orientation deviates from the direction of motion—contrary to the forward-facing perspective of robots—or capture unsafe behaviors, such as weaving though dense crowds, posing risks for robotic deployment (Bar et al. 2025). UrbanNav tackles this by employing a filtering pipeline, leveraging visual odometry to estimate camera pose and detect misalignment between viewpoint and trajectory, and integrating object detector to identify and exclude segments with unsafe interactions. This ensures only high-quality, robot-compatible data informs training.

*Second, how to obtain instruction-action supervision from in-the-wild videos for imitation learning?* Manual annotation of such videos is infeasible, necessitating an automated, scalable approach (He et al. 2025a). UrbanNav addresses this through a sophisticated pipeline. Using off-the-shelf visual odometry models, we extract egocentric trajectories with pose information from video sequences. Next, a large vision-language model (VLM) detects contextually relevant landmarks in urban scenes, followed by another VLM that produce natural, navigation-oriented instructions grounded in these landmarks. This pipeline yields a dataset of over 1,500 hours of navigation data and 3 million instruction-trajectory-landmark triplets, enabling robust, scalable language-guided navigation policies for di-

verse urban environments.

By training policy models on this large-scale dataset of real-world human trajectories, UrbanNav achieves superior performance in navigating complex urban environments. It demonstrates strong generalization to previously unseen cities and robust resilience to real-world challenges. Our key contributions are three-fold:

- We recognize language-guided embodied urban navigation as a complex and critical challenge and introduce UrbanNav, a scalable framework that harnesses web-scale human walking videos for robust navigation.

- We develop an automated data processing pipeline that filters and extracts egocentric trajectories from in-the-wild videos, generating instructions to enable large-scale imitation learning without requiring manual annotations.

- We demonstrate that training on web-scale data significantly enhances navigation performance in real-world experiments, empowering embodied agents to navigate complex urban environments effectively.

## 2 Related Works

**Navigation in Simulation.** The research on indoor navigation has greatly advanced with the increasing realism and diversity of simulators (Chang et al. 2017; Savva et al. 2019; Kolve et al. 2017; Deitke et al. 2022). In these environments, some works focus on improving robotic navigation capabilities through imitation learning using expert demonstration data (Ehsani et al. 2024) or online reinforcement learning (Zeng et al. 2024). Other works (Zhou, Hong, and Wu 2024; Qiao et al. 2024; Zhang et al. 2025a; Ding et al. 2025) have attempted to leverage the vast prior knowledge and strong generalization capabilities of LLM for zero-shot navigation reasoning. However, these methods suffer significant performance degradation (Gervet et al. 2023) and low efficiency (Zhu et al. 2024; Zhang et al. 2025b) when deployed on physical robotic platforms.

**Real-World Navigation.** Directly using real-world navigation data (Hirose et al. 2018; Shah et al. 2021; Karnan et al. 2022; Hirose et al. 2023) to construct training samples for supervised learning is an efficient approach. Some works (Shah et al. 2023a,b; Shah and Levine 2022) train models on a mix of real-robot navigation datasets, enabling direct deployment on different robotic platforms with strong generalization capabilities, while others (Sridhar et al. 2024; Bar et al. 2025; Dong et al. 2025) employ diffusion models to generate trajectories or leverage the model's imagination of the future to assist navigation. However, manually collecting such data is costly and lacks diversity. This work automates the construction of navigation datasets using web-scale data to overcome this limitation.

**Language-Guided Policies for Robotics.** The task of using text descriptions as navigation instructions has been widely studied. Vision-and-Language Navigation (Anderson et al. 2018; Qi et al. 2020; Krantz et al. 2020; Liu et al. 2025a) requires the agent to follow fine-grained navigation instructions to move within a scene. Object Goal Navigation
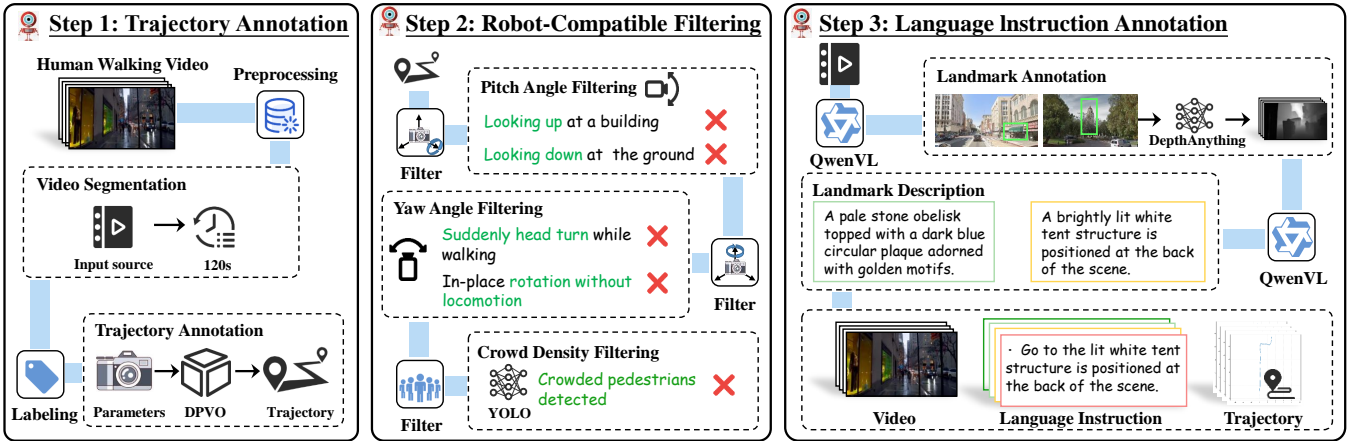
Figure 2: **UrbanNav Data Construction Pipeline.** The process is divided into three main steps: 1) Trajectory Annotation, where human walking videos are preprocessed, segmented, and then annotated with a camera pose estimator to generate trajectories. 2) Robot-Compatible Data Filtering, where low-quality segments with pitch, yaw, or dense crowd issues are automatically filtered out. 3) Language Instruction Annotation, where a large language model is used to generate rich, descriptive language instructions for each trajectory, along with landmark bounding boxes and depth maps.

(Chaplot et al. 2020; Wani et al. 2020; Cai et al. 2024; Sadek et al. 2023; Yokoyama et al. 2024) aims to enable agent to efficiently locate specific objects in the environment based on textual descriptions. In contrast, we focus on outdoor scenarios and train the language-guided navigation model using only real-world data, thereby avoiding this gap and enabling the model to learn navigation rules and environmental affordances inherent in urban environments.

**Learning from Web Videos.** Learning from web-scale video data has advanced language and vision tasks (Wang et al. 2024b; He et al. 2024, 2025b). However, a key challenge is the lack of action labels for navigation tasks. LeLaN (Hirose et al. 2024) uses vision-language model prompting and pretrained navigation models to generate action labels, while (Liu et al. 2025b) applies visual odometry to annotate video frames with pose information. However, LeLaN focuses on close-range object navigation within indoor static scenes, making it incapable of handling complex urban environments, and the constructed dataset does not encompass obstacle avoidance capabilities. In contrast, Urban-Nav distinguishes itself by implementing a robust filtering pipeline to select robot-compatible data, ensuring viewpoint alignment and safe navigation behaviors. Additionally, our approach generates language-grounded instructions tied to landmarks and trajectories, enabling scalable, cost-effective, and high-quality action label generation for imitation learning in complex urban navigation tasks.

# 3 Methodology

## 3.1 Problem Formulation

We address the challenge of last-mile navigation, a task where an embodied agent, equipped with an egocentric camera, must precisely navigate in an urban environment to a goal location specified by a free-form natural language instruction. This requires the agent to interpret vague or context-dependent directives and translate them into a sequence of low-level control actions. The central objective is to develop a control policy $\pi$. Formally, given a natural language instruction $g$ defining the goal, the agent receives its current RGB observation $o_t$ at each time step. To enable robust navigation, the policy also leverages a historical context of past $k$ observations $o_{(t-k):t}$ (we set $k = 8$). The policy is thus defined as a function that predicts the next action $a_t$: $\pi(a_t|o_{(t-k):t}, g)$.

## 3.2 Labeling In-the-Wild Videos

**Human Walking Videos.** As shown in Fig. 2, to train our language-guided urban navigation framework, we curate a large-scale dataset of over 2,000 hours of in-the-wild egocentric human walking videos sourced from YouTube. These videos capture first-person perspectives of pedestrians navigating diverse urban environments, such as bustling city streets, residential neighborhoods, and park pathways. The dataset encompasses a wide range of conditions, such as varying weather, lighting, and obstacle densities, reflecting the complexity and dynamism of real-world urban settings. Human walking trajectories in these videos closely resemble the egocentric, forward-facing motion of robotic agents, making them highly relevant for training navigation policies. Unlike teleoperated or robotic datasets, which are often limited in scale and diversity, these videos provide rich, naturalistic navigation behaviors, including adaptive maneuvers around obstacles and adherence to social norms like maintaining safe distances from pedestrians.

**Trajectory Annotation.** We uniformly segment the raw videos into 2-minute clips, each representing a candidate navigation trajectory. To extract pose information, we employ the state-of-the-art visual odometry model DPVO (Teed, Lipson, and Deng 2023). The first frame of each clip is taken as the origin of a local world coordinate system, and we annotate the camera pose of each subsequent frame
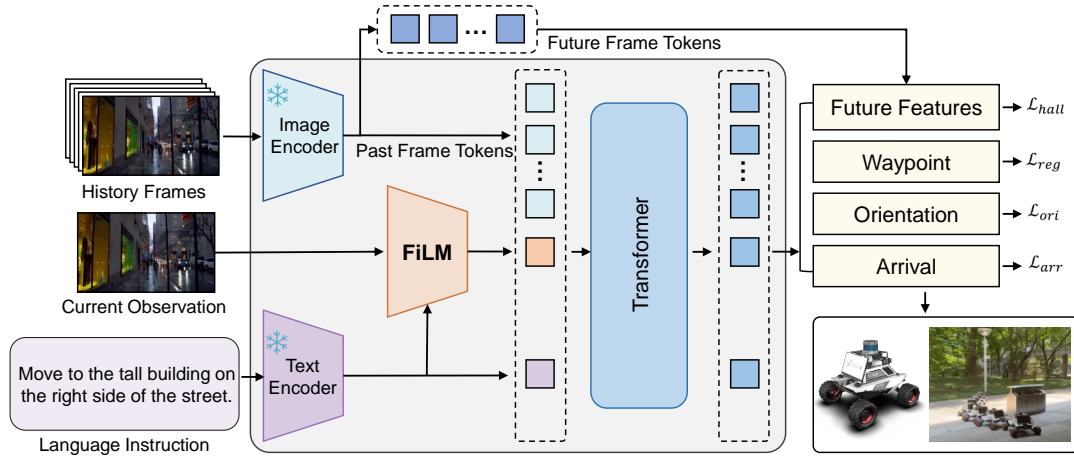
Figure 3: **Overall Illustration of UrbanNav.** The model takes historical images and a language instruction as input, fuses their features, and uses a Transformer to predict future frame features, waypoints, directions, and an arrival status.

relative to it. Although visual odometry may suffer from cumulative drift across long horizons, this effect is largely mitigated in our case: the policy is trained to predict future actions based on only the past 8 steps, operating within short temporal windows where VO remains locally consistent. This design, following the practice in (Liu et al. 2025b), enables reliable supervision from noisy egocentric videos. As a result, we obtained a total of 106,603 egocentric navigation trajectories with pose labels, spanning 3,553 hours.

**Robot-Compatible Data Filtering.** A crucial yet often overlooked issue in learning from in-the-wild walking videos is that not all segments are suitable for training embodied agents. Robotic platforms typically employ forward-facing, fixed-view cameras, requiring visual inputs that maintain consistent alignment with motion direction. However, Unlike robots, human walkers exhibit flexible body dynamics and frequently adjust their head orientation, resulting in significant pitch variations and misalignment between the camera viewpoint and the actual trajectory. Such discrepancies render many video clips incompatible with robotic requirements. Additionally, some segments feature densely crowded scenes with close pedestrian interactions, which pose safety risks for robotic deployment. Unlike prior work (Liu et al. 2025b), which often overlooks these issues and risks degrading policy performance, our approach employs a robust filtering pipeline to select predominantly robot-compatible data for training.

First, we estimate the per-frame camera pitch angle and reject trajectories exhibiting excessive vertical oscillations—specifically, those with pitch variation beyond $15°$. We further analyze the alignment between movement direction and viewing direction using a sliding window, and discard segments with significant directional divergence (over $60°$), which commonly results from abrupt head turns or side glances. Second, to eliminate trajectories captured in overly crowded areas, we apply YOLOv10 (Wang et al. 2024a) to detect pedestrians. Based on empirical observations, we discard any trajectory where more than five people appear in a single frame and such occurrences happen in more than

three frames, indicating sustained dense proximity. After applying these filtering steps, we retain 47,008 high-quality, robot-compatible trajectories spanning 1,566 hours, which serve as the core data for training.

**Language Instruction Annotation.** To enable language-guided navigation, we aim to identify feasible landmarks along walking trajectories and annotate them with natural language descriptions. Landmarks are selected based on the following criteria: (1) they must be located near the walking trajectory to ensure reachability; (2) they should possess clear, distinguishable visual features, including both large-scale structures (e.g., buildings or sculptures) and smaller but stable street objects (e.g., signboards or traffic lights); and (3) dynamic entities, including pedestrians and vehicles, are excluded to ensure stability and consistency.

We leverage Qwen2.5-VL-72B (Bai et al. 2025), a state-of-the-art VLM, to automatically detect and localize candidate landmarks in video frames via prompted queries aligned with the above criteria. The model outputs both bounding boxes and preliminary landmark names. To ensure quality, we manually review and filter out low-confidence or ambiguous annotations. For the retained landmarks, we prompt Qwen2.5-VL-72B to generate concise and descriptive natural language instructions. Through this process, we obtain a total of 3 million landmark annotations, each paired with a bounding box and a language description. On average, each trajectory contains 65 identified landmarks, with the mean description length being 17 words.

### 3.3 Policy Architecture and Training

**Architecture.** As illustrated in Fig. 3, our policy model predicts a sequence of future egocentric positions based on both language instructions and visual observations. The input to the model comprises four components: (1) the language instruction, (2) the current visual observation, (3) the past $k$ visual observations. It builds upon previous works (Hirose et al. 2024; Liu et al. 2025b). We use CLIP (Radford et al. 2021) to encode the language instruction and DINOv2 (Oquab et al. 2023) to extract features from all

| Method | Test Seen | | | | Test Unseen | | | |
|---|---|---|---|---|---|---|---|---|
| | AOE ↓ | MAOE ↓ | ADE ↓ | MADE ↓ | AOE ↓ | MAOE ↓ | ADE ↓ | MADE ↓ |
| Nomad + CLIP | 22.85 | 39.10 | 3.61 | 6.89 | 22.77 | 39.12 | 3.65 | 6.96 |
| ViNT + CLIP | 13.37 | 19.58 | 1.32 | 2.37 | 13.69 | 20.08 | 1.39 | 2.50 |
| LeLaN | 10.14 | 16.25 | 0.93 | 1.77 | 10.36 | 16.49 | 0.98 | 1.84 |
| **UrbanNav (Ours)** | **8.88** | **14.62** | **0.83** | **1.57** | **9.22** | **14.99** | **0.88** | **1.67** |

Table 1: **UrbanNav Benchmark.** We evaluated the model performance in both seen and unseen environments using the Ur-BanNav offline data. AOE and MAOE are used to evaluate the angular difference between the predictions and ground truth (in degrees). ADE and MADE assess the distance difference (in meters).

visual frames. Both encoders are kept frozen during training. To ground the current observation in the instruction, we apply a FiLM module (Perez et al. 2018) to modulate the current visual embedding using the language features. All input tokens—including language embeddings, FiLM-modulated current visual features and historical visual observation—are concatenated and fed into a Transformer encoder. This encoder captures both temporal dynamics and cross-modal interactions, and outputs a predicted future trajectory in egocentric coordinates. To enable smoother and more anticipatory navigation, we adopt a multi-step prediction framework following prior work such as (Sridhar et al. 2024). At each time step, the model predicts the agent's way-points for the next $k = 8$ steps using a lightweight action head.

**Training Objective.** The training objective integrates four complementary loss terms that jointly supervise spatial accuracy, directional correctness, goal awareness, and predictive understanding of future observations. The total loss is a weighted sum of these terms, defined as:

$$L_{\text{total}} = \lambda_{\text{reg}}L_{\text{reg}} + \lambda_{\text{ori}}L_{\text{ori}} + \lambda_{\text{arr}}L_{\text{arr}} + \lambda_{\text{hall}}L_{\text{hall}},$$

where the weights $\lambda$ are selected to ensure all loss terms operate within comparable magnitude ranges. The individual loss terms are formulated as follows: The *waypoint regression loss* $L_{\text{reg}}$ minimizes the L2 distance between the predicted and ground-truth positions over the future time horizon, ensuring accurate spatial forecasting. The *orientation loss* $L_{\text{ori}}$ penalizes angular discrepancy between the predicted and ground-truth motion directions. It is calculated using negative cosine similarity over the future horizon $k$.

The *arrival prediction loss* $L_{\text{arr}}$ is a binary cross-entropy loss that supervises the model's ability to judge whether the navigation goal has been reached. The *feature hallucination loss* $L_{\text{hall}}$ acts as an auxiliary loss that guides the model to anticipate high-level visual features of future scenes, thereby encouraging an internal modeling of scene dynamics. This loss is defined as the L1 distance between the predicted and ground-truth features over the future horizon $k$:

$$L_{\text{hall}} = \frac{1}{k} \sum_{f=1}^{k} \|\hat{h}_{t+f} - h_{t+f}\|_1,$$

where $\hat{h}_{t+f}$ is the ground-truth visual feature vector extracted from the future observation and $h_{t+f}$ is the corresponding predicted feature vector.

**Training Details.** For each training sample, we randomly select a trajectory and one of its annotated landmarks, along with the corresponding language instruction $g$, as the navigation goal. We then sample a starting time step $t$ between 10 and 60 frames before the goal time $t_g$, so that the agent learns to navigate toward the target from diverse initial distances and directions. The segment from $t$ to $t_g$, combined with the instruction $g$, forms the input for the policy. To help the model learn when to stop, we also simulate goal-reached scenarios by sampling time steps very close to $t_g$, and label these as arrival cases. This enables the model to distinguish between approaching and already-at-goal situations.

## 4 Experiments

### 4.1 Experimental Setup

**Baselines.** We compare our model against several prominent policies previously for real-world navigation. To ensure a fair comparison, we adapt these baselines to our language-guided task. For instance, NoMaD (Sridhar et al. 2024) and ViNT (Shah et al. 2023b) were originally designed for image-goal navigation, and lacked native support for textual instructions. To adapt them, we augmented their architectures by encoding the natural language instruction with CLIP (Radford et al. 2021), concatenating the resulting text features with the visual features and passing them through a fusion layer, while the core network remained unchanged. Similarly, LeLaN (Hirose et al. 2024), originally developed for indoor object-goal navigation with textual inputs, predicts robot linear and angular velocities. To align its output space with our task formulation for consistent evaluation, we adapted its output head to directly regress the future navigation waypoints.

**Metrics.** Given the challenges of evaluating end-to-end task completion in real-world environments without autoregressive execution, we designed a comprehensive evaluation protocol that includes both offline benchmarking and real-world deployment. **(1) Offline Evaluation.** For offline evaluation, we assess the model's performance on a held-out validation set. The model is given a sequence of historical observations, positions, and a natural language instruction. Its task is to predict a future trajectory, and we measure the deviation between this prediction and the ground-truth trajectory. Following prior work (Liu et al. 2025b), we use the average orientation error (AOE) and maximum average orientation error (MAOE) to measure the directional alignment between predicted and ground-truth actions. To
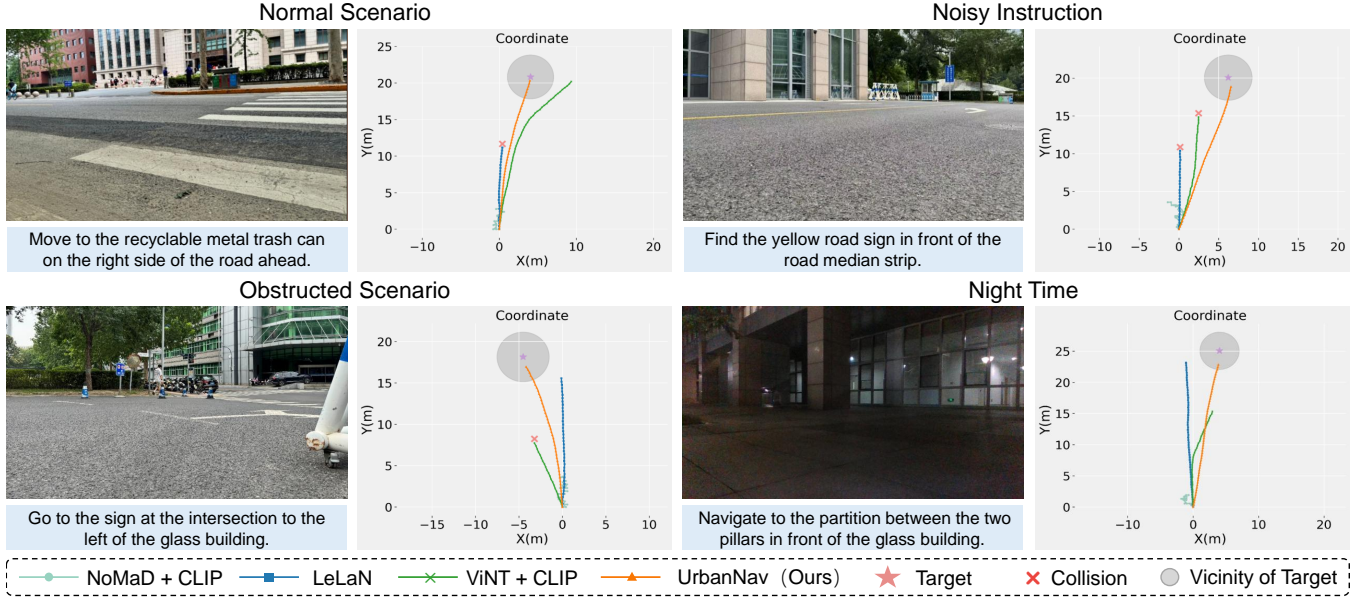
Figure 4: **Qualitative Results.** The figures show trajectory visualizations in four different scenarios. For each set of images, the left side represents the initial observation and instruction, while the right side shows the real-robot trajectories and target positions in the world coordinate system for different methods.

better capture local path-following capabilities essential for obstacle avoidance, we introduce two additional metrics: average distance error (ADE), which calculates the mean L2-distance between predicted and ground-truth positions, and maximum average distance error (MADE), which represents the Fréchet distance (Alt and Godau 1995) in a discrete setting. For these metrics, a lower value indicates better performance. A more detailed definition of these metrics is provided in the appendix. **(2) Real-World Deployment.** To validate the policy's real-world generalization, we deploy the model on a physical robot platform. The primary metric for this evaluation is the navigation success rate, which measures the percentage of trials where the robot successfully reaches the specified goal landmark without any collisions.

## 4.2 Performance Benchmarking

**Offline Evaluation.** We evaluated UrbanNav's performance on our validation set, which is divided into "seen" and "unseen" components. The seen portion contains scenes present in the training set, while the unseen portion consists of entirely new environments. As shown in Table 1, our approach achieves state-of-the-art results across all metrics on both seen and unseen data. UrbanNav's trajectories demonstrate superior alignment with the ground truth, outperforming baselines in terms of both directional accuracy (AOE and MAOE) and precise path-following (ADE and MADE). This robust performance on unseen environments confirms that our framework effectively learns generalizable navigation policies from web-scale human trajectory data.

**Real-World Deployment.** As shown in Table 2, our ablation study, UrbanNav* (trained exclusively on real-world data), achieved a notably lower overall success rate com-

pared to models pre-trained on web-scale data. This result highlights the crucial benefits of our pre-training approach, which provides a strong, generalizable foundation for effective navigation. The full UrbanNav model achieves a superior overall success rate of 83.3%, a significant margin over the second-best performing baseline, LeLaN (62.5%). While all methods experienced performance degradation in nighttime scenarios due to visual noise from the robot's cameras, UrbanNav maintained a high success rate (75.0%), demonstrating its robustness and strong generalization to completely unseen real-world environments.

## 4.3 Robustness Analysis in Challenging Scenarios

To validate our method's robustness against real-world complexities, we designed several challenging scenarios for testing, with results presented in Table 3. We categorized these scenarios into three types. The "Normal" case involves clear language instructions where the target is within the initial field of view. In contrast, "Noisy" scenarios use ambiguous or misleading instructions, while "Obstructed" cases indicate the target is initially outside the field of view or occluded.

**Noisy Language Instructions.** Our method achieved a 100% success rate in the normal case and maintained a high success rate of 87.5% in noisy conditions. The full policy's ability to handle variations in language is a direct benefit of pre-training on our diverse, web-scale dataset, a finding further validated by the relatively poor performance of UrbanNav* (the variant without web-scale pre-training) in these same scenarios.

**Obstructed Targets.** For the obstructed case, where targets were initially out of view or occluded, we observed a

| Method | Overall | Day Time | Night Time |
|---|---|---|---|
| *Trained on Real-World Data Only* | | | |
| UrbanNav* | 33.4 | 41.7 | 25.0 |
| *Pre-trained with Web-Scale Data* | | | |
| Nomad + CLIP | 29.2 | 33.4 | 25.0 |
| ViNT + CLIP | 45.8 | 50.0 | 41.7 |
| LeLaN | 62.5 | 75.0 | 58.3 |
| **UrbanNav (Ours)** | **83.3** | **91.7** | **75.0** |

Table 2: **Real-World Navigation Results.** The table shows the success rate in unseen real-world environments, with results separately shown for daytime and nighttime conditions.

performance degradation with a 62.5% success rate. This is an expected outcome as our policy is primarily designed for local navigation, not long-term exploration for initially invisible targets. However, our approach consistently outperformed all other methods in this challenging condition, confirming a clear advantage. The robustness to environmental and instructional changes is a direct result of training on our diverse, web-scale dataset.

**Qualitative Results.** Figure 4 presents a visual comparison of UrbanNav and other baselines across various scenarios. UrbanNav successfully navigates to the target even in challenging conditions, while the baseline methods frequently fail, either misinterpreting the language instructions or resulting in collisions. This superior performance demonstrates the efficacy of our model's ability to leverage environmental affordances and strong instruction-following capabilities.

| Method | Normal | Noisy | Obstructed |
|---|---|---|---|
| *Trained on Real-World Data Only* | | | |
| UrbanNav* | 62.5 | 25.0 | 12.5 |
| *Pre-trained with Web-Scale Data* | | | |
| Nomad + CLIP | 50.0 | 25.0 | 12.5 |
| ViNT + CLIP | 62.5 | 37.5 | 25.0 |
| LeLaN | 75.0 | 62.5 | 37.5 |
| **UrbanNav (Ours)** | **100.0** | **87.5** | **62.5** |

Table 3: **Robustness in Challenging Scenarios.** The table presents the success rates under different difficulty conditions. "Normal" refers to simple scenarios, "Noisy" indicates noisy language instructions, and "Obstructed" denotes scenarios with obstacles or occlusions.

### 4.4 Ablation Studies

**Impact of Model Components.** The results of our ablation study on key architectural components are shown in Table 4. The findings indicate that the FiLM feature fusion module is crucial for performance, as its removal incurs a substantial performance degradation. We hypothesize that using language instructions to modulate visual features allows the agent to better attend to semantic cues pertinent to the navigation goal, thereby enhancing directional guidance.

| Components | | Test Unseen | | | |
|---|---|---|---|---|---|
| Feature Hall. | FiLM | AOE↓ | MAOE↓ | ADE↓ | MADE↓ |
| ✓ | | 11.35 | 17.54 | 1.07 | 1.94 |
| | ✓ | 9.56 | 15.51 | 0.92 | 1.71 |
| ✓ | ✓ | **9.22** | **14.99** | **0.88** | **1.67** |

Table 4: **Ablation Study of Model Components.** The table shows the results of our ablation study on the feature fusion and feature hallucination loss in unseen environments.
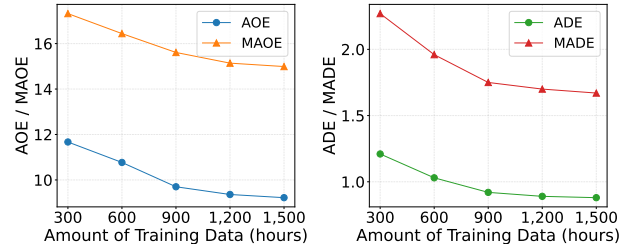


Figure 5: **Impact of Web data scaling.** The figures show the performance of UrbanNav on unseen environments as a function of the training data size.

Furthermore, we observe that the feature hallucination loss provides a clear performance benefit in unseen scenarios, which contrasts with the findings of some prior works that it has a negative impact on zero-shot inference.. We attribute this success to our use of high-quality, robot-compatible data. By training on a clean dataset that avoids behavioral discrepancies like viewpoint inconsistency, our framework allows the auxiliary loss to effectively enable the model to predict future observations, a capability that directly contributes to robust navigation.

**Impact of Scaling Up Web Data.** To validate the effectiveness of our web-scale data, we conducted an ablation study on the impact of training data quantity on model performance. As shown in Figure 5, we observe a consistent and significant decrease in all error metrics as the training data size increases from 300 to 1,500 hours. This trend provides strong empirical evidence that larger, more diverse datasets enable the model to learn a more effective policy. The performance improvement begins to plateau around 1,200 hours, highlighting the benefits and scalability of our UrbanNav framework.

## 5 Conclusion

In this work, we introduced UrbanNav, a novel framework for language-guided urban navigation that overcomes data scarcity by leveraging web-scale human walking videos. Our approach uses a scalable data pipeline to create a substantial dataset for large-scale imitation learning. By training on this diverse data, UrbanNav achieves superior performance with strong generalization and remarkable resilience in dynamic urban environments. Ultimately, UrbanNav offers a practical path toward real-world deployment, proving that agents trained on web-scale human trajectories can robustly handle the complexities of last-mile navigation.

# References

Alt, H.; and Godau, M. 1995. Computing the Fréchet distance between two polygonal curves. *International Journal of Computational Geometry & Applications*, 5(01n02): 75–91.

Anderson, P.; Wu, Q.; Teney, D.; Bruce, J.; Johnson, M.; Sünderhauf, N.; Reid, I.; Gould, S.; and Van Den Hengel, A. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3674–3683.

Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; et al. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Bar, A.; Zhou, G.; Tran, D.; Darrell, T.; and LeCun, Y. 2025. Navigation world models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 15791–15801.

Cai, W.; Huang, S.; Cheng, G.; Long, Y.; Gao, P.; Sun, C.; and Dong, H. 2024. Bridging zero-shot object navigation and foundation models through pixel-guided navigation skill. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 5228–5234. IEEE.

Chang, A.; Dai, A.; Funkhouser, T.; Halber, M.; Niessner, M.; Savva, M.; Song, S.; Zeng, A.; and Zhang, Y. 2017. Matterport3D: Learning from RGB-D Data in Indoor Environments. *International Conference on 3D Vision (3DV)*.

Chaplot, D. S.; Gandhi, D. P.; Gupta, A.; and Salakhutdinov, R. R. 2020. Object goal navigation using goal-oriented semantic exploration. *Advances in Neural Information Processing Systems*, 33: 4247–4258.

Deitke, M.; VanderBilt, E.; Herrasti, A.; Weihs, L.; Ehsani, K.; Salvador, J.; Han, W.; Kolve, E.; Kembhavi, A.; and Mottaghi, R. 2022. ProcTHOR: Large-Scale Embodied AI Using Procedural Generation. *Advances in Neural Information Processing Systems*, 35: 5982–5994.

Ding, H.; Xu, Z.; Fang, Y.; Wu, Y.; Chen, Z.; Shi, J.; Huo, J.; Zhang, Y.; and Gao, Y. 2025. LaViRA: Language-Vision-Robot Actions Translation for Zero-Shot Vision Language Navigation in Continuous Environments. *arXiv preprint arXiv:2510.19655*.

Dong, Y.; Wu, F.; Chen, G.; Cheng, Z.-Q.; Hu, Q.; Zhou, Y.; Sun, J.; He, J.-Y.; Dai, Q.; and Hauptmann, A. G. 2025. Unified World Models: Memory-Augmented Planning and Foresight for Visual Navigation. arXiv:2510.08713.

Ehsani, K.; Gupta, T.; Hendrix, R.; Salvador, J.; Weihs, L.; Zeng, K.-H.; Singh, K. P.; Kim, Y.; Han, W.; Herrasti, A.; et al. 2024. Spoc: Imitating shortest paths in simulation enables effective navigation and manipulation in the real world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16238–16250.

Gao, Y.; Li, C.; You, Z.; Liu, J.; Li, Z.; Chen, P.; Chen, Q.; Tang, Z.; Wang, L.; Yang, P.; et al. 2025. OpenFly: A Comprehensive Platform for Aerial Vision-Language Navigation. *arXiv preprint arXiv:2502.18041*.

Gervet, T.; Chintala, S.; Batra, D.; Malik, J.; and Chaplot, D. S. 2023. Navigating to objects in the real world. *Science Robotics*, 8(79): eadf6991.

He, H.; Yang, C.; Lin, S.; Xu, Y.; Wei, M.; Gui, L.; Zhao, Q.; Wetzstein, G.; Jiang, L.; and Li, H. 2025a. Cameractrl ii: Dynamic scene exploration via camera-controlled video diffusion models. *arXiv preprint arXiv:2503.10592*.

He, T.; Gao, J.; Xiao, W.; Zhang, Y.; Wang, Z.; Wang, J.; Luo, Z.; He, G.; Sobanbab, N.; Pan, C.; et al. 2025b. ASAP: Aligning Simulation and Real-World Physics for Learning Agile Humanoid Whole-Body Skills. *arXiv preprint arXiv:2502.01143*.

He, T.; Luo, Z.; He, X.; Xiao, W.; Zhang, C.; Zhang, W.; Kitani, K.; Liu, C.; and Shi, G. 2024. Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning. *arXiv preprint arXiv:2406.08858*.

Hirose, N.; Glossop, C.; Sridhar, A.; Shah, D.; Mees, O.; and Levine, S. 2024. Lelan: Learning a language-conditioned navigation policy from in-the-wild videos. *arXiv preprint arXiv:2410.03603*.

Hirose, N.; Sadeghian, A.; Vázquez, M.; Goebel, P.; and Savarese, S. 2018. Gonet: A semi-supervised deep learning approach for traversability estimation. In *2018 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, 3044–3051. IEEE.

Hirose, N.; Shah, D.; Sridhar, A.; and Levine, S. 2023. Sacson: Scalable autonomous control for social navigation. *IEEE Robotics and Automation Letters*, 9(1): 49–56.

Karnan, H.; Nair, A.; Xiao, X.; Warnell, G.; Pirk, S.; Toshev, A.; Hart, J.; Biswas, J.; and Stone, P. 2022. Socially compliant navigation dataset (scand): A large-scale dataset of demonstrations for social navigation. *IEEE Robotics and Automation Letters*, 7(4): 11807–11814.

Kolve, E.; Mottaghi, R.; Han, W.; VanderBilt, E.; Weihs, L.; Herrasti, A.; Deitke, M.; Ehsani, K.; Gordon, D.; Zhu, Y.; et al. 2017. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*.

Krantz, J.; Wijmans, E.; Majumdar, A.; Batra, D.; and Lee, S. 2020. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, 104–120. Springer.

Kümmerle, R.; Ruhnke, M.; Steder, B.; Stachniss, C.; and Burgard, W. 2013. A navigation system for robots operating in crowded urban environments. In *2013 IEEE International Conference on Robotics and Automation*, 3225–3232. IEEE.

Liu, Q.; Zhang, S.; Qiao, Y.; Zhu, J.; Li, X.; Guo, L.; Wang, Q.; He, X.; Wu, Q.; and Liu, J. 2025a. GroundingMate: Aiding Object Grounding for Goal-Oriented Vision-and-Language Navigation. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 1775–1784. IEEE.

Liu, X.; Li, J.; Jiang, Y.; Sujay, N.; Yang, Z.; Zhang, J.; Abanes, J.; Zhang, J.; and Feng, C. 2025b. Citywalker: Learning embodied urban navigation from web-scale videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 6875–6885.

Muhlbauer, Q.; Sosnowski, S.; Xu, T.; Zhang, T.; Kuhnlenz, K.; and Buss, M. 2009. Navigation through urban environments by visual perception and interaction. In *2009 IEEE In-*

ternational Conference on Robotics and Automation, 3558–3564. IEEE.

Oquab, M.; Darcet, T.; Moutakanni, T.; Vo, H.; Szafraniec, M.; Khalidov, V.; Fernandez, P.; Haziza, D.; Massa, F.; El-Nouby, A.; et al. 2023. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*.

Perez, E.; Strub, F.; De Vries, H.; Dumoulin, V.; and Courville, A. 2018. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Qi, Y.; Wu, Q.; Anderson, P.; Wang, X.; Wang, W. Y.; Shen, C.; and Hengel, A. v. d. 2020. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9982–9991.

Qiao, Y.; Lyu, W.; Wang, H.; Wang, Z.; Li, Z.; Zhang, Y.; Tan, M.; and Wu, Q. 2024. Open-Nav: Exploring Zero-Shot Vision-and-Language Navigation in Continuous Environment with Open-Source LLMs. *arXiv preprint arXiv:2409.18794*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PmLR.

Sadek, A.; Bono, G.; Chidlovskii, B.; Baskurt, A.; and Wolf, C. 2023. Multi-Object Navigation in real environments using hybrid policies. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 4085–4091. IEEE.

Savva, M.; Kadian, A.; Maksymets, O.; Zhao, Y.; Wijmans, E.; Jain, B.; Straub, J.; Liu, J.; Koltun, V.; Malik, J.; et al. 2019. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9339–9347.

Schumann, R.; and Riezler, S. 2022. Analyzing generalization of vision and language navigation to unseen outdoor areas. *arXiv preprint arXiv:2203.13838*.

Shah, D.; Eysenbach, B.; Kahn, G.; Rhinehart, N.; and Levine, S. 2021. Rapid exploration for open-world navigation with latent goal models. *arXiv preprint arXiv:2104.05859*.

Shah, D.; and Levine, S. 2022. Viking: Vision-based kilometer-scale navigation with geographic hints. *arXiv preprint arXiv:2202.11271*.

Shah, D.; Sridhar, A.; Bhorkar, A.; Hirose, N.; and Levine, S. 2023a. Gnm: A general navigation model to drive any robot. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 7226–7233. IEEE.

Shah, D.; Sridhar, A.; Dashora, N.; Stachowicz, K.; Black, K.; Hirose, N.; and Levine, S. 2023b. ViNT: A foundation model for visual navigation. *arXiv preprint arXiv:2306.14846*.

Sridhar, A.; Shah, D.; Glossop, C.; and Levine, S. 2024. Nomad: Goal masked diffusion policies for navigation and exploration. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 63–70. IEEE.

Teed, Z.; Lipson, L.; and Deng, J. 2023. Deep patch visual odometry. *Advances in Neural Information Processing Systems*, 36: 39033–39051.

Wang, A.; Chen, H.; Liu, L.; Chen, K.; Lin, Z.; Han, J.; et al. 2024a. Yolov10: Real-time end-to-end object detection. *Advances in Neural Information Processing Systems*, 37: 107984–108011.

Wang, Y.; Zheng, S.; Cao, B.; Wei, Q.; Zeng, W.; Jin, Q.; and Lu, Z. 2024b. Scaling Large Motion Models with Million-Level Human Motions. *arXiv preprint arXiv:2410.03311*.

Wani, S.; Patel, S.; Jain, U.; Chang, A.; and Savva, M. 2020. Multion: Benchmarking semantic map memory using multi-object navigation. *Advances in Neural Information Processing Systems*, 33: 9700–9712.

Yokoyama, N.; Ramrakhya, R.; Das, A.; Batra, D.; and Ha, S. 2024. Hm3d-ovon: A dataset and benchmark for open-vocabulary object goal navigation. In 2024 IEEE. In *RSJ International Conference on Intelligent Robots and Systems (IROS)*, 5543–5550.

Yu, M.-M.; Zhu, F.; Liu, W.; Yang, Y.; Wang, Q.; Wu, W.; and Liu, J. 2025. C-NAV: Towards Self-Evolving Continual Object Navigation in Open World. arXiv:2510.20685.

Zeng, K.-H.; Zhang, Z.; Ehsani, K.; Hendrix, R.; Salvador, J.; Herrasti, A.; Girshick, R.; Kembhavi, A.; and Weihs, L. 2024. Poliformer: Scaling on-policy rl with transformers results in masterful navigators. *arXiv preprint arXiv:2406.20083*.

Zhang, S.; Qiao, Y.; Wang, Q.; Guo, L.; Wei, Z.; and Liu, J. 2025a. Flexvln: Flexible adaptation for diverse vision-and-language navigation tasks. *arXiv preprint arXiv:2503.13966*.

Zhang, S.; Qiao, Y.; Wang, Q.; Yan, Z.; Wu, Q.; Wei, Z.; and Liu, J. 2025b. COSMO: Combination of Selective Memorization for Low-cost Vision-and-Language Navigation. *arXiv preprint arXiv:2503.24065*.

Zhou, G.; Hong, Y.; and Wu, Q. 2024. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 7641–7649.

Zhu, J.; Qiao, Y.; Zhang, S.; He, X.; Wu, Q.; and Liu, J. 2024. MiniVLN: Efficient Vision-and-Language Navigation by Progressive Knowledge Distillation. *arXiv preprint arXiv:2409.18800*.