

NaVILA: LEGGED ROBOT VISION-LANGUAGE-ACTION MODEL FOR NAVIGATION

An-Chieh Cheng^{1*}, Yandong Ji^{1*}, Zhaojing Yang^{2*}, Xueyan Zou¹,
 Jan Kautz³, Erdem Biyik², Hongxu Yin^{3†}, Sifei Liu^{3†}, Xiaolong Wang^{13†}
¹UC San Diego ²USC ³NVIDIA



Walk forward, when seeing the stair bars, turn right and walk around the stairs until reaching the hallway. Turn right and walk along the hallway, stop in front of a bathroom.



Turn right immediately, walk along the hallway, turn left at the end and enter the most left bedroom.



Walk all the way down, turn left at the intersection and find a bookshelf.



Walk to the other end of the room, turn left and find a toy kitchen set.

Figure 1: Real-world demonstration of NaVILA: Upon receiving human instructions, NaVILA uses a vision-language model to process RGB video frames and employs locomotion skills to execute the task on a robot. The robot successfully handles long-horizon navigation tasks and operates safely in challenging environments.

ABSTRACT

This paper proposes to solve the problem of Vision-and-Language Navigation with legged robots, which not only provides a flexible way for humans to command but also allows the robot to navigate through more challenging and cluttered scenes. However, it is non-trivial to translate human language instructions all the way to low-level leg joint actions. We propose NaVILA, a 2-level framework that unifies a Vision-Language-Action model (VLA) with locomotion skills. Instead of directly predicting low-level actions from VLA, NaVILA first generates mid-level actions with spatial information in the form of language, (e.g., “moving forward 75cm”), which serves as an input for a visual locomotion RL policy for execution. NaVILA substantially improves previous approaches on existing benchmarks. The same advantages are demonstrated in our newly developed benchmarks with IsaacLab, featuring more realistic scenes, low-level controls, and real-world robot experiments.

* Equal contribution, ordered alphabetically. † Equal advising.
 Codes and trained models will be made publicly available. More results are provided at .

1 INTRODUCTION

The ability to perform Vision-and-Language Navigation (VLN) has become a foundational component in modern robotics systems. With VLN, a robot is expected to navigate around unseen environments without a provided map following a language instruction (Anderson et al., 2018; Wang et al., 2019; Chaplot et al., 2020a;b;c; Ramrakhya et al., 2022). This not only offers a better interface for humans, but also strengthen cross-scene generalization through languages. In this paper, we further extend the study of VLN with legged robots (e.g., quadruped or humanoid). Using legs instead of wheels allows robots to navigate in more challenging and cluttered scenarios. As the examples shown in Figure 1, our robot can navigate through a crowded office with narrow walkways and scattered desks, or a messy home with toys and other objects on the floor.

To translate language to action, the robot needs to reason about the input language, and perform closed-loop planning as well as low-level control. With the recent advancement in Large Language Models (LLMs) and Vision-Language Models (VLMs), several end-to-end Vision-Language-Action (VLA) systems have been developed (Brohan et al., 2023; Kim et al., 2024; Padalkar et al., 2024). These systems fine-tune a general-purpose VLM with large-scale robot manipulation demonstrations to produce low-level actions for control. While unifying reasoning and execution in a single model is fascinating and shows encouraging results, it is worthy to dive deeper into the question: Is there a better way to represent actions beyond the quantized low-level commands? After all, LLMs and VLMs were primarily trained with natural language. Unifying reasoning and execution becomes challenging when we need to convert that reasoning into precise, non-verbal actions.

Inspired by the recent progress on VLM (Chen et al., 2024a; Cheng et al., 2024) for spatial location and distance reasoning, we propose **NaViLA**, a two-level framework for legged robot VLN: A VLM is fine-tuned to output a **mid-level action** (VLA) in the form of language such as “turn right 30 degrees”, and a low-level visual locomotion policy is trained to follow this instruction for execution. The mid-level action output of the VLA conveys the location and direction information without the low-level commands. The advantages of this framework are two-fold: (i) We disentangle the low-level execution from the VLA so the same VLA can be deployed across robots by switching the low-level policy, and we can train the model with different robot data or even human data; (ii) Having the action represented in the form of language allows the VLA to be trained with other reasoning QA tasks at the same time, which improves the model’s reasoning ability instead of overfitting the outputs to low-level commands.

To train the VLA, we demonstrate how to: (i) Integrate historical context and current observations in VLN within existing VLM frameworks; (ii) Create a specialized navigation prompt tailored for VLN tasks; (iii) Introduce a carefully curated dataset blend designed to enhance VLN generalizability. These strategies allow us to fine-tune a general-purpose image-based VLM into a navigation-focused agent while simultaneously training it on general vision-language datasets, thereby maintaining its broad generalization capabilities.

During training the locomotion skills, we employ a single-stage approach to learn vision-based locomotion policy. We construct a heightmap from raw LIDAR point clouds and introduce randomization to bridge the sim-to-real gap. This controller takes the output from our VLA model, converts it into command velocities, and tracks these velocities by controlling the positions of the joints. To our knowledge, this is the first end-to-end approach for training visual locomotion skills that are both robust and safe, enabling deployment in real-world, challenging environments (e.g., strong sunlight or near certain transparent surfaces).

Importantly, our framework operates on two distinct timescales: the VLA, typically a large and computationally intensive model, runs at a lower frequency, providing high-level navigation commands; while the locomotion policy operates in real-time. This dual-frequency approach allows the locomotion policy to handle instant obstacle avoidance and increases overall robustness.

In our experiments, we show that our VLA significantly outperforms the state-of-the-arts on classic VLN benchmarks, with over 12% improvement in success rate. Notably, our VLA performs comparably to state-of-the-art methods that rely on pre-computed maps. To better simulate the challenges of locomotion navigation in VLN, we introduce a new benchmark, VLN-CE-Isaac, using Isaac Sim. This benchmark considers detailed robotic joint movements and interactions with environments, which prior VLN works have not explored. In our VLN-CE-Isaac experiments, our vision-based policy outperforms the blind policy by a significant margin, showing a 14% improvement in success rate. We also demonstrate that our VLA can be deployed across different robots (Unitree Go2 and Unitree H1), each using distinct locomotion skills. Finally, we deploy our VLA model in the real

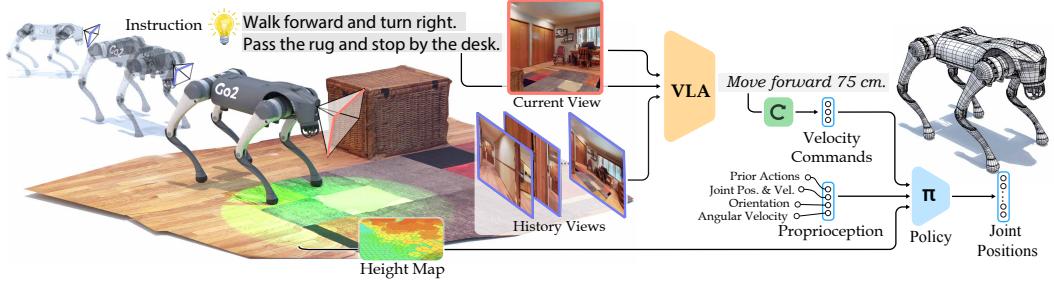


Figure 2: NaVILA is a two-level framework combining high-level visual language understanding with low-level locomotion control. Our VLA model processes single-view images to produce mid-level actions in natural language, which are then converted into precise joint movements by an advanced low-level locomotion policy. This integration allows for strong generalization and adaptability across different real-world environments, and can operate the robot in real-time.

world, exhibiting impressive robustness and achieving 55% success rate on 20 instructions, including a 70% success rate on complex instructions, across diverse scenes.

2 METHOD

Our VLA model (NaVILA) integrates high-level visual language understanding and action with low-level locomotion control (Figure 2). NaVILA employs a VLM that processes single-view images to generate waypoint instructions in natural language. These instructions are then interpreted by a low-level locomotion policy, which translates them into precise joint movements for real-time robot control. The synergy between the VLM’s high-level reasoning and the locomotion policy’s execution capabilities enables our method to demonstrate remarkable generalization and adaptability across diverse real-world environments. In the following sections, we detail the components of our approach. We begin by describing how we tame VLMs for high-level VLN in Sec. 2.1, followed by an overview of our robot configuration and low-level locomotion policy in Sec. 2.2.

2.1 TAMING VLMs FOR HIGH-LEVEL VISION LANGUAGE NAVIGATION

VLN requires processing video inputs as observations. A common approach to handling video inputs in VLMs is through video encoders. However, recent progress in VLMs has largely been driven by the availability of image-text data. While there have been efforts to extend this success to video encoders, the lack of large, high-quality video-text datasets has limited their pre-training. To address this challenge, we opt for image-based vision-language models in our approach. These models exhibit stronger generalization abilities and possess broader knowledge, making them more suitable for tackling the generalization challenges in VLN. Specifically, we built our approach upon VILA, an image-based VLM pre-trained with interleaved image-text corpus. VILA’s pre-training has proven particularly effective for multi-image reasoning, making it especially suitable for VLN tasks where understanding sequential image relationships is critical.

VILA Preliminary. VILA consists of three main components: a vision encoder, a projector, and an LLM. The vision encoder processes the input images, converting them into a sequence of visual tokens. These tokens are then downsampled and mapped into the language domain via an MLP projector. Afterward, the projected tokens, along with text tokens, are sent to the LLM for auto-regressive generation. When handling videos, VILA uniformly sampled frames at regular intervals. It puts all the frame information before any text. A typical prompt for describing a video might look like “⟨frame3⟩⟨frame6⟩⟨frame9⟩...Tell me about this video.” VILA undergoes a 3-stage training process: first, it pre-trains a connector between the frozen LLM and vision backbones using alignment data (Liu et al., 2023); then it pre-trains both the connector and the LLM using text-image interleaved corpus (Byeon et al., 2022; Zhu et al., 2024); and finally, it fine-tunes all modules (vision encoder, connector, LLM) with instruction tuning data (Liu et al., 2023; 2024).

Navigation Prompts. In vision-language navigation tasks, images from different time steps serve two distinct purposes. The image at time step t represents the current observation, which is crucial for a VLN agent to make immediate decisions (e.g., turning right at an intersection or stopping when

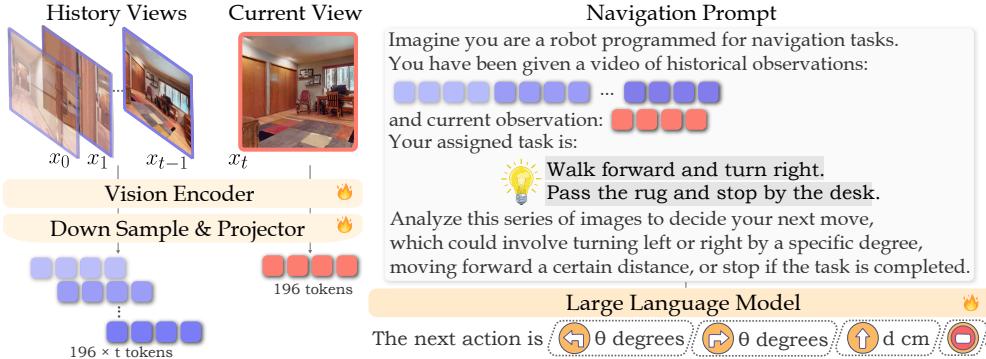


Figure 3: Overview of our VLA framework. We denote the purple blocks (●) as memory tokens sampled from historical frames, and the red blocks (●) as the current observation tokens. 🔥 denotes trainable parameters.

the goal is reached). On the other hand, frames before time step t are historical frames that function as a memory bank, helping the agent track overall progress (e.g., remembering the starting location, reasoning about places already visited and planning the next step). Uniformly sampling frames at regular intervals, as done in VILA, is not ideal because it doesn't differentiate between these two types of representations. Therefore, we first extract the most recent frame t as the current observation and then uniformly sample frames from the preceding $t - 1$ frames, ensuring the first frame is always included. Additionally, since current and historical observations serve different roles, we distinguish them in our task prompt using textual cues like a `video of historical observations`: for memory frames and `current observation`: for the latest frame. Unlike (Zhang et al., 2024), we avoid introducing additional special tokens that could complicate the LLM's learning process. Instead, we adhere to our design principle of keeping both the input and output of LLM in the language domain to fully leverage the reasoning capabilities of the pre-trained LLM. By integrating these tokens for historical and current observations with the navigation instruction, we construct a navigation task prompt, as shown in Figure 2.

Supervised Fine-tuning Data Blend. Effective Supervised Fine-tuning (SFT) data is crucial for developing a robust vision-language action model. Such a model should be specialized for an embodied task yet avoid overfitting to specific actions. It should also generalize well to real-world scenarios while retaining broad-world knowledge. With these goals in mind, we designed our SFT data blend from three perspectives.

First, we focus on navigation-specific data. Currently, there are limited options for VLN datasets in continuous environments, with only R2R-CE (Krantz et al., 2020b) and RxR-CE (Ku et al., 2020) available. These datasets provide sparse path points converted from discrete VLN versions. We utilize both datasets within the Habitat simulator, employing a shortest path follower to generate action sequences that adhere to the geodesic shortest path. This results in step-wise navigation videos, where each sample in our dataset comprises a $t + 1$ frames video and the corresponding oracle action at time step t . To encourage the LLM to generate continuous value labels for distances and angles, we merge consecutive actions (e.g., combining two forward 25 cm steps into a single forward 50 cm step), with a maximum of three consecutive actions. This merging process has two key advantages: it reduces the dataset size for more efficient processing, and it helps prevent overfitting by introducing greater diversity in the actions. Additionally, to address label imbalance, particularly the underrepresentation of the stop action, we apply a rebalancing technique for a more even distribution. For all navigation-specific data, we apply the previously described frame extraction strategy and navigation task prompt.

Second, to improve scene understanding and address the limitations of R2R-CE and RxR-CE, which have few scenes and limited instructions, we incorporate auxiliary navigational datasets. Following (Zhang et al., 2024), we use augmented instructions from EnvDrop (Tan et al., 2019) and introduce an auxiliary task of navigation trajectory summarization. Specifically, given a trajectory video, we sample frames by retaining the first frame and uniformly sampling the remaining ones as historical frames, and use the annotated instructions as labels. The LLM is then tasked with describing the robot's navigation trajectory based on these frames. To further encourage spatial scene understanding, we integrate the ScanQA (Azuma et al., 2022) dataset, which features real-world 3D scan QA



Figure 4: Height map reconstruction from the point cloud. (a) Go2 robot tracks linear and angular velocity commands while avoiding collision with obstacles in simulation. The red dots represent the LiDAR point cloud, raycasting from the sensor’s center towards the terrain mesh. The left image is a preprocessed height map from the LiDAR data, with values clipped according to real-world sensor constraints. Darker colors represent lower values. (b) Safe robot locomotion near a transparent glass surface. The top-down height map clearly detects the glass, while the forward-facing depth and RGB images struggle to capture it.

pairs involving human-edited questions and free-form answers grounded to 3D objects within each scene. For training, we utilize the multi-view RGB images from the raw scans to support this task.

Finally, to maintain the model’s general capabilities, we include general video/image VQA datasets from (Liu et al., 2024; Chen et al., 2024d; Maaz et al., 2024). This comprehensive dataset design enables NaVILA to generalize effectively to novel scenes and real-world environments.

Training and Inference Paradigm. Our training process begins with the stage two model of VILA, which has already undergone visual language corpus pre-training. We then apply our SFT data blend to train the entire VLM for one epoch, following standard practices. During this training, all three components—vision encoder, connector, and LLM—are unfrozen. For the inference phase, we implement a regular expression parser (Kearns, 1991), to extract action types (such as forward or turn left) and their corresponding arguments (like specific distance or angles) from the LLM output. This method has demonstrated effectiveness in both simulated environments and real-world experiments, where we empirically found that all actions throughout all experiments are successfully matched and mapped.

2.2 VISUAL LOCOMOTION POLICY

In this section, we begin with a brief overview of the Go2 robot dog, the experimental platform used in this work. Next, we describe the development of the end-to-end vision-based control policy, which interprets high-level language navigation commands from the VLM and converts them into precise joint movements. This control policy is trained in the Isaac Sim simulator using Isaac Lab (Mittal et al., 2023) and then directly deployed to the real-world robot.

Go2 Robot. As shown in Figure 4, the robot is equipped with a LiDAR sensor mounted at the base of its head, broadcasting point clouds at a frequency of 15Hz. The robot features 18 degrees of freedom (DoFs), comprising 6 DoFs for its base and 3 DoFs for each of its four legs. In the policy training process, we left the 6 DoFs on the base unconstrained so that the policy only controls the 12 joint motors on the legs.

Interpreting High-level Commands As in our formulation, VLM outputs a fixed set of actionable words, such as {move forward, turn left, turn right, stop}, we casts these instructions to fixed command velocities $\{0.5 \text{ m s}^{-1}, \frac{\pi}{6} \text{ rad s}^{-1}, -\frac{\pi}{6} \text{ rad s}^{-1}, 0\}$ and execute with corresponding time durations to align with the specific VLM value.

Low-level Action and Observation Space. The action space a of the control policy is defined as the desired joint position $q^d \in \mathbb{R}^{12}$, which are converted into torque input for the simulator using the stiffness and dampness. We adopt PPO algorithm (Schulman et al., 2017) to train the policy. During training, the critic observes the privileged environment and generates a value function to update the actor. The actor then only receives sensor data available in the real world. The observation space of the critic o^c contains the proprioception and velocity command at the current time step t and a privileged terrain height scan around the robot. The proprioceptive data includes robot linear and angular velocity, orientation, joint positions, joint velocities, and the previous action. In the actor’s observation space o^a , linear velocity is excluded, as it is unavailable in the real world, and instead,

Table 1: Comparison with state-of-the-art methods on the Val-Unseen split of R2R-CE (Krantz et al., 2020b) and RxR-CE (Ku et al., 2020). * indicates methods using the waypoint predictor from Hong et al. (2022). NaVILA achieves remarkable performance among methods using only single-view RGB input and shows competitive results compared to those using panorama, depth, or odometry sensors.

	Observation				R2R Val-Unseen				RxR Val-Unseen			
	S.RGB	Pano.	Depth	Odo.	NE ↓	OS ↑	SR ↑	SPL ↑	NE ↓	SR ↑	SPL ↑	nDTW ↑
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	47.0
VLN \cup BERT* (Hong et al., 2022)	✓	✓	✓		5.74	53.0	44.0	39.0	8.98	27.0	22.6	46.7
Sim2Sim* (Krantz & Lee, 2022)	✓	✓	✓		6.07	52.0	43.0	36.0	-	-	-	-
GridMM* (Wang et al., 2023c)	✓	✓	✓		5.11	61.0	49.0	41.0	-	-	-	-
Ego ² -Map* (Hong et al., 2023a)	✓	✓	✓		5.54	56.0	47.0	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	63.3
ETPNav* (An et al., 2024)	✓	✓	✓		4.71	65.0	57.0	49.0	5.64	54.7	44.8	61.9
HNR* (Wang et al., 2024)	✓	✓	✓		4.42	67.0	61.0	51.0	5.50	56.3	46.7	63.5
BEVBert* (An et al., 2023)	✓	✓	✓		4.57	67.0	59.0	50.0	4.00	68.5	-	69.6
HAMT+ScaleVLN* (Wang et al., 2023d)	✓	✓	✓		4.80	-	55.0	51.0	-	-	-	-
AG-CMTP (Chen et al., 2021a)	✓	✓	✓		7.90	39.0	23.0	19.0	-	-	-	-
R2R-CMTP (Chen et al., 2021a)	✓	✓	✓		7.90	38.0	26.0	22.0	-	-	-	-
LAW (Raychaudhuri et al., 2021)	✓	✓	✓		6.83	44.0	35.0	31.0	10.90	8.0	8.0	38.0
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., 2024b)		✓	✓		5.55	59.0	47.0	33.0	7.06	43.3	30.5	50.1
Seq2Seq (Krantz et al., 2020a)	✓	✓			7.77	37.0	25.0	22.0	12.10	13.9	11.9	30.8
CMA (Krantz et al., 2020a)	✓	✓			7.37	40.0	32.0	30.0	-	-	-	-
RGB-Seq2Seq (Krantz et al., 2020a)	✓				10.10	8.0	0.0	0.0	-	-	-	-
RGB-CMA (Krantz et al., 2020a)	✓				9.55	10.0	5.0	4.0	-	-	-	-
NaVid (Zhang et al., 2024)	✓				5.47	49.0	37.0	35.0	-	-	-	-
NaVILA	✓				5.37	57.6	49.7	45.5	6.23	50.8	45.1	60.3

a history of proprioceptive data is used to infer this information implicitly. The robot perceives the surrounding terrain using a heightmap from the LiDAR sensor.

Incorporating Height Map from LiDAR Point Cloud. Given LiDAR’s superior ability to detect transparent objects and robust performance under strong sunlight, we chose the manufacturer-provided LiDAR as the primary sensor for perceiving the robot’s surroundings and ensuring safe navigation. The Unitree L1 generates point clouds with a wide field of view of $360^\circ \times 90^\circ$, from which we create a 2.5D height map based on the parameters listed in Table 9. For each voxel grid, the lowest value within the range is selected, and a maximum filter is then applied over the last 5 lidar point clouds to smooth the resulting height map.

Training. Different from most existing works (Lee et al., 2020; Miki et al., 2022; Margolis et al., 2022; Lee et al., 2024) that utilize the two-stage teacher-student training paradigm, we adopt a single-stage manner to train the locomotion policy. Compared to two-stage training, single-stage RL is more time-efficient as it eliminates the need for policy distillation. Additionally, the policy interacts directly with the environment, allowing it to explore and potentially discover novel strategies. Leveraging Isaac Lab’s high-speed simulation, our vision-based RL policy training achieves a high throughput over 60K FPS on an RTX 4090 GPU. Training details are included in Appx. E.

3 EXPERIMENTS

We conduct experiments to answer the following questions: (1) How does our VLA’s performance compare to state-of-the-art methods in VLN-CE benchmarks and general spatial scene understanding tasks? (Section 3.1) (2) How to evaluate locomotion navigation in simulators, and how effective and flexible is NaVILA in these scenarios? (Section 3.2) (3) Can NaVILA pipeline be successfully deployed in real robot VLN experiments? (Section 3.3)

3.1 VLM PERFORMANCE

VLN-CE Benchmarks. We evaluate our VLM on the VLN-CE benchmarks, which provide continuous environments for executing navigational actions in reconstructed photorealistic indoor

Table 2: Cross-dataset performance on the RxR-CE (Ku et al., 2020) Val-Unseen split. All results are obtained without training on the RxR-CE training set. NaVILA significantly outperforms NaVid (Zhang et al., 2024), the current single-view state-of-the-art.

	Observation			RxR Val-Unseen			
	S.RGB	Depth	Odo.	NE ↓	OS ↑	SR ↑	SPL ↑
LAW (Raychaudhuri et al., 2021)	✓	✓	✓	10.87	21.0	8.0	8.0
CM2 (Georgakis et al., 2022)	✓	✓	✓	8.98	25.3	14.4	9.2
WS-MGMap (Chen et al., 2022)	✓	✓	✓	9.83	29.8	15.0	12.1
Seq2Seq (Krantz et al., 2020a)	✓	✓		11.8	5.02	3.51	3.43
CMA (Krantz et al., 2020a)	✓	✓		11.7	10.7	4.41	2.47
RGB-Seq2Seq (Zhang et al., 2024)	✓			11.2	12.2	0.0	0.0
RGB-CMA (Zhang et al., 2024)	✓			9.55	14.8	0.0	0.0
A ² NAV (Chen et al., 2023)	✓			-	-	16.8	6.3
NaVid (Zhang et al., 2024)	✓			8.41	34.5	23.8	21.2
NaVILA	✓			8.78	46.8	34.3	28.2

Table 3: Evaluation of spatial scene understanding performance on the ScanQA dataset (Azuma et al., 2022) Validation split. NaVILA outperforms current state-of-the-art VLA models and demonstrates comparable or superior performance to other 3D LMMs (LMMs) that require additional input, such as depth or camera pose. Note that * indicates 3D LMMs that require task-specific fine-tuning on the ScanQA dataset. [†] denotes model trained without navigation-specific and auxiliary data.

	ScanQA Validation				
	Bleu-4 ↑	Rouge ↑	Cider ↑	Meteor ↑	EM ↑
Task-specific Specialist					
VoteNet+MCAN (Yu et al., 2019)	6.2	29.8	54.7	11.4	17.3
ScanRefer+MCAN (Yu et al., 2019)	7.9	30.0	55.4	11.5	18.6
ScanQA (Azuma et al., 2022)	10.1	33.3	64.9	13.1	21.0
3D-VisTA (Zhu et al., 2023)	10.4	35.7	69.6	13.9	22.4
3D Large Multi-modal Models					
3D-LLM(<i>flamingo</i>) [*] (Hong et al., 2023b)	7.2	32.3	59.2	12.2	20.4
3D-LLM(<i>BLIP2-flants5</i>) [*] (Hong et al., 2023b)	12.0	35.7	69.4	14.5	20.5
LL3DA [*] (Chen et al., 2024e)	13.5	37.3	76.8	15.9	-
Chat-3DV2 [*] (Huang et al., 2024a)	14.0	-	87.6	-	-
Scene-LLM [*] (Fu et al., 2024)	12.0	40.0	80.0	16.6	27.2
LEO (Huang et al., 2024b)	13.2	49.2	101.4	20.0	24.5
2D Vision-Language-Action Models					
NavILM (Zheng et al., 2024)	12.0	38.4	75.9	15.4	23.0
NaVILA [†]	9.6	31.9	61.3	12.8	18.4
NaVILA	14.8	46.4	95.1	18.7	27.0

scenes. We focus on the val-unseen split in both R2R (Room-to-Room) and RxR (Room-across-Room) datasets within VLN-CE, as these are the two most recognized benchmarks in VLN. We employ the following widely used evaluation metrics for VLN tasks: Navigation Error (NE), Oracle Success Rate (OS), Success Rate (SR), Success-weighted Path Length (SPL), and normalize dynamic time wrapping (nDTW). We show results in Table 1, where NaVILA significantly surpasses all baselines under identical conditions (i.e., single-view RGB) in both benchmarks using a single model. Notably, this also marks the first time a VLN agent, trained solely on single-view RGB input, achieves comparable or superior results to models that use panoramic views, odometry, or simulator-pretrained waypoint predictors. This suggests that NaVILA’s strong generalization capabilities can effectively compensate for the limited observations in RGB views or odometry.

To evaluate the cross-dataset performance, we follow (Zhang et al., 2024) by training NaVILA exclusively on R2R samples and all other auxiliary datasets, but leaving out the RxR training set. We then evaluate its zero-shot performance on the RxR Val-Unseen data split. As shown in Table 2, our method significantly outperforms NaVid, the current state-of-the-art model, with a substantial 10% improvement in success rate.

Spatial Scene Understanding Benchmarks. As a general navigation agent, robust spatial scene understanding (e.g., object localization, referring, and spatial reasoning) is crucial. To evaluate NaVILA’s capabilities in scene understanding, we conduct evaluations on the ScanQA Val benchmark, a

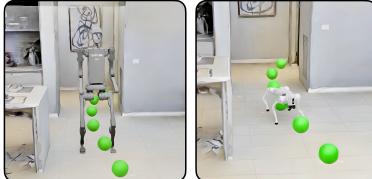


Figure 5: VLN-CE-Isaac Benchmark

Table 4: VLN-CE-Isaac Evaluation

	Low-level Observation			VLN-CE-Isaac			
	Proprio.	LiDAR	Height Scan	NE ↓	OS ↑	SR ↑	SPL ↑
NaVILA-Go2 blind	✓			6.03	49.0	36.2	33.3
NaVILA-Go2 vision	✓	✓		5.49	58.7	50.2	45.5
NaVILA-H1 blind	✓			7.67	33.3	24.4	21.0
NaVILA-H1 vision	✓			5.86	54.6	45.3	40.3

widely used dataset for 3D Question Answering. ScanQA is based on real-world scans, and we use multi-view images from these scans as input to query NaVILA for answers. As shown in Table 3, NaVILA significantly outperforms the previous state-of-the-art model, NaviLLM, by a substantial margin (20 points higher on the CIDEr score). Moreover, our method’s performance is comparable to, if not better than, the state-of-the-art 3D-based large multi-modal models. This is particularly noteworthy as these other models require either 3D scans or RGBD data with camera poses as inputs, while our method achieves similar or better results with less observation.

3.2 LEGGED ROBOT NAVIGATION PERFORMANCE IN SIMULATION

High-fidelity VLN-CE-Isaac Benchmark. Though the VLN-CE benchmark within the Habitat simulator is widely used for navigation policy evaluation, Habitat’s limited simulation fidelity restricts the ability to fully evaluate systems that integrate high-level planning with precise low-level robotic control. For example, the agent is able to navigate through a narrow 10cm gap between two sofas, which is unrealistic for legged robots such as quadrupeds or humanoids. To address this limitation, we introduce an evaluation benchmark VLN-CE-Isaac within Isaac Sim. Leveraging Isaac Sim’s high-fidelity simulation of detailed robotic joint movements and interactions with the environment, VLN-CE-Isaac enables a comprehensive evaluation of the entire pipeline, from high-level navigation planning to precise robotic execution. The same scenes from VLN-CE are imported, with robots spawned into the environment, as demonstrated in Figure 5. Given the variation in mesh quality across scenes, we select 1077 trajectories from scenes with high-quality meshes, out of the original 1839 trajectories in the R2R Val-Unseen split, to make sure all trajectories are traversable. To maintain consistency with the VLN-CE, the same metrics are used for evaluation: Navigation Error (NE), Oracle Success Rate (OS), Success Rate (SR), and Success-weighted Path Length (SPL).

In our experiments, we test our NaVILA model on a Unitree Go2 robot within the proposed VLN-CE-Isaac benchmark. To highlight the effectiveness of the vision-based policy, we performed an ablation study comparing it against a proprioception-only (blind) policy. As shown in Table 4, the vision-based policy outperforms the blind policy by 14% in Success Rate, owing to its superior obstacle avoidance capability of the vision-based policy. More details are attached in Appx. C.

Notably, VLN-CE-Isaac is compatible with a variety of robotic platforms. To demonstrate this flexibility, we also conduct evaluations using a Unitree H1 robot as in Figure 5. A blind policy with proprioception observation was trained and tested alongside our NaVILA model. The results of this evaluation are presented in Table 4. We observe that the success rate of NaVILA on the H1 robot is significantly lower than on the Go2, which is expected due to the larger size of the humanoid robot.

3.3 REAL WORLD EVALUATION

Experiment Setup. We used a total of 10 simple and 10 complex instructions across two environments: a lab space and a home setting. Simple instructions consisted of one or two navigation commands, where the robot did not need to navigate between rooms (e.g., ‘Go to the chair and stop’). In contrast, complex instructions involved three or four commands, requiring the robot to traverse multiple rooms (e.g., ‘Walk out of the room, turn right, enter the room in front of you, and stop at the table’). These tasks are similar to the trajectories in R2R-CE and RxR-CE. We compared our approach to two baselines: one with our NaVILA model where the VLA was trained only on R2R data, and the other using GPT-4o.

Results Analysis. As shown in Table 5, the NaVILA model trained on mixed data achieves the best performance on both simple and complex commands, significantly outperforming the commonly-

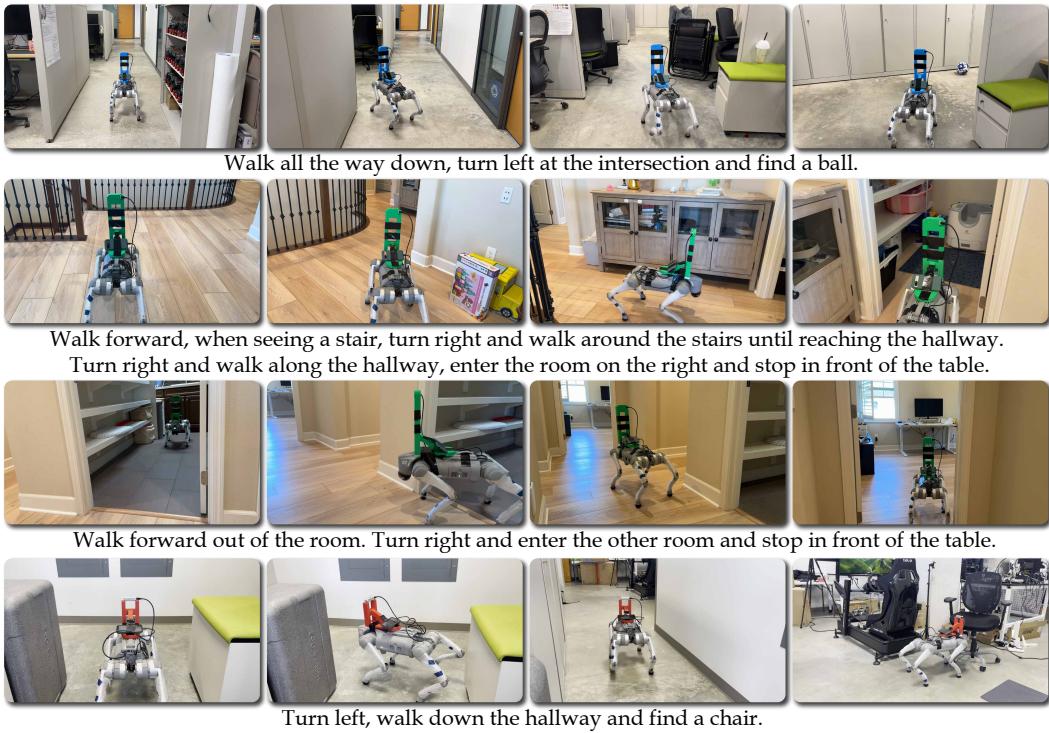


Figure 6: Qualitative results from the real-world deployment of NaVILA. The robot demonstrates the ability to solve long-horizon navigation tasks, following human instructions across various real-world environments.

used GPT-4o. Notably, the success rate on complex commands is higher than on simple commands, likely because the R2R and RxR instructions are lengthy, allowing the model to become proficient at following long instructions. Additionally, we observe that the model trained solely on R2R navigation data performed the worst among the three, indicating that the auxiliary data is essential to the model’s generalization capabilities in real-world.

4 RELATED WORK

Vision-Language Navigation. Visual navigation has been a long-standing research topic in robotics for decades (Moravec, 1980; Elfes, 1987; Thrun et al., 2001; Gervet et al., 2023). Classical navigation approaches often rely on pre-computed maps (Thrun et al., 1999) or build geometric maps of the environment using depth sensors (Newcombe et al., 2011) or monocular RGB cameras while localizing the robot simultaneously (SLAM) (Davison et al., 2007; Jones & Soatto, 2011). In recent years, learning-based approaches with Imitation Learning (Chaplot et al., 2018; Codevilla et al., 2018) and Reinforcement Learning (Mnih et al., 2015; Lillicrap, 2015) have not only shown impressive generalization results but also enabled wider applications such as vision-and-language navigation (VLN).

Vision-Language Navigation (VLN) is a fundamental challenge in embodied AI, where agents navigate complex environments using visual cues and natural language instructions. The field has evolved significantly over time. Early research (Anderson et al., 2018; Ku et al., 2020; Qi et al., 2020) focused on discrete navigation in simulated environments like MP3D (Chang et al., 2017), where agents teleport between predefined nodes on a navigation graph (Fried et al., 2018; Ma et al., 2019; Tan et al., 2019; Ke et al., 2019; Hong et al., 2020; Chen et al., 2021b; 2024c; Zhou et al., 2024). As foundation models advanced, many VLN systems improved dramatically by leveraging large-scale pre-trained models (Li et al., 2019; Majumdar et al., 2020) and pre-training techniques (Guhur et al., 2021; Wang et al., 2023d; Kamath et al., 2023), approaching human-level performance in this setting. However, this setup emphasized high-level decision-making while neglecting the challenges of underlying motion control. Recently, research (Raychaudhuri et al., 2021; Chen et al., 2022; Georgakis et al., 2022; Chen et al., 2024b; Zhang et al., 2024) has shifted towards continuous environments (VLN-CE (Krantz et al., 2020a)) using simulators like Habitat (Savva

Table 5: Real-world experiments conducted in three environments (Laboratory, House, and Outdoor). Simple and Complex refer to simple and complex instruction-following tasks, respectively.

	Laboratory				House				Outdoor			
	Simple		Complex		Simple		Complex		Simple		Complex	
	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↑
GPT-4o Krantz et al. (2020a)	2.01	0.67	2.38	0.33	1.49	0.53	3.00	0.00	-	0.67	-	0.50
NaVILA	2.00	0.60	1.81	0.73	2.17	0.47	2.32	0.40	-	0.00	-	0.00
NaVILA (w/ Real Videos)	1.29	1.00	1.76	0.80	1.15	1.00	1.76	0.67	-	1.00	-	0.83

et al., 2019). This introduces greater complexity, as agents must perform mid-level actions such as moving forward or rotating, rather than teleporting between nodes. To bridge the gap between discrete and continuous navigation, some approaches (Irshad et al., 2021; Krantz & Lee, 2022; An et al., 2023; 2024) use simulator pre-trained waypoint models (Hong et al., 2022; Krantz et al., 2021) that predict candidate positions around the agent and have shown significant performance gains. However, they often struggle to generalize due to their reliance on simulator-specific data. Additionally, the candidate positions predicted by these models only cover nearby locations and do not account for low-level motion planning or obstacle avoidance. In this paper, we aim to advance VLN towards real-world robotics applications, particularly for challenging legged robots. We propose a model that handles both high-level decision-making and generates low-level actions to control the robot’s full motion. Additionally, we introduce a new VLN benchmark built on Isaac Sim, offering a more realistic simulation environment, which we believe will benefit future work in VLN.

Robot Foundation Models. Robot foundation models aim to provide a unified framework that processes inputs from various modalities, such as vision and language, and directly outputs actions to enable robots to perform complex tasks. Existing works (Brohan et al., 2023; Team et al., 2024; Kim et al., 2024) trained on large-scale robotic dataset to get general robot policies, but mainly focusing on manipulation tasks. Doshi et al. (2024) and Yang et al. (2024) proposed end-to-end visual-language cross-embodiment models for different robotic tasks. As for legged robots, Ding et al. (2024) proposed a unified model to leverage vision and language inputs and generate executable low-level actions. However, these methods struggle to understand complex instructions which are crucial for navigation tasks. Based on this, we propose a VLA model specifically designed for navigation tasks. Our model generates high-level action commands, which are then executed by a low-level policy. This approach enables the robot to interpret complex instructions and navigate effectively towards the goals.

Legged Robot Locomotion Learning. Legged robot learning for locomotion navigation focuses on enabling robots to traverse various terrains. Previous works (Wang et al., 2023b; Long et al., 2024) rely solely on robot’s proprioceptive information struggle in scenarios like obstacle avoidance. Other end-to-end vision-based approaches (Kareer et al., 2023; Yang et al., 2021; Imai et al., 2022; Yang et al., 2023) are vulnerable to extreme environmental conditions, such as intense sunlight, due to the limitations of sensors. Lee et al. (2020) and Miki et al. (2022) incorporate LiDAR sensors in addition to depth cameras to improve terrain sensing, but rely on time-inefficient two-state training. To overcome these limitations, we propose a single-stage RL framework that integrates LiDAR sensing inputs, allowing the robot to directly learn from interacting with the environments for more efficient learning and robust performance in complex scenarios.

5 CONCLUSION

We introduce NaVILA, a powerful two-level framework that unifies vision-language models (VLMs) with locomotion skills for generic navigation tasks. NaVILA generates high-level, language-based commands, while a real-time locomotion policy handles obstacle avoidance. This dual-frequency design improves robustness and flexibility across different robots. By preserving reasoning capabilities through language-based actions, NaVILA avoids overfitting and can be trained on broader tasks. In experiments, NaVILA shows a 12% improvement on classic VLN benchmarks, outperforms vision-blind policies in our new VLN-CE-Isaac1K benchmark, and demonstrates strong real-world performance across diverse scenes. The source code will be released upon publication.

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. In *CVPRW*, 2022. 6
- Dong An, Yuankai Qi, Yangguang Li, Yan Huang, Liang Wang, Tieniu Tan, and Jing Shao. Bevbert: Multimodal map pre-training for language-guided navigation. In *ICCV*, 2023. 6, 10
- Dong An, Hanqing Wang, Wenguan Wang, Zun Wang, Yan Huang, Keji He, and Liang Wang. Etp-nav: Evolving topological planning for vision-language navigation in continuous environments. *IEEE TPAMI*, 2024. 6, 10
- 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 *CVPR*, 2018. 2, 9
- Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *CVPR*, 2022. 4, 7
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint*, 2023. 2, 10
- Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022. 3
- Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgbd data in indoor environments. In *3DV*, 2017. 9
- Devendra Singh Chaplot, Kanthashree Mysore Sathyendra, Rama Kumar Pasumarthi, Dheeraj Rajagopal, and Ruslan Salakhutdinov. Gated-attention architectures for task-oriented language grounding. In *AAAI*, 2018. 9
- Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. In *ICLR*, 2020a. 2
- Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*, 2020b. 2
- Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *CVPR*, 2020c. 2
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvilm: Endowing vision-language models with spatial reasoning capabilities. In *CVPR*, 2024a. 2
- Jiaqi Chen, Bingqian Lin, Xinmin Liu, Xiaodan Liang, and Kwan-Yee K Wong. Affordances-oriented planning using foundation models for continuous vision-language navigation. *arXiv preprint*, 2024b. 6, 9
- Jiaqi Chen, Bingqian Lin, Ran Xu, Zhenhua Chai, Xiaodan Liang, and Kwan-Yee Wong. Mapgpt: Map-guided prompting with adaptive path planning for vision-and-language navigation. In *ACL*, 2024c. 9
- Kevin Chen, Junshen K Chen, Jo Chuang, Marynel Vázquez, and Silvio Savarese. Topological planning with transformers for vision-and-language navigation. In *CVPR*, 2021a. 6
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *ECCV*, 2024d. 5

- 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. In *NeurIPS*, 2022. 6, 7, 9
- Peihao Chen, Xinyu Sun, Hongyan Zhi, Runhao Zeng, Thomas H. Li, Gaowen Liu, Mingkui Tan, and Chuang Gan. A²nav: Action-aware zero-shot robot navigation by exploiting vision-and-language ability of foundation models. *arXiv preprint*, 2023. 7
- Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, and Ivan Laptev. History aware multimodal transformer for vision-and-language navigation. In *NeurIPS*, 2021b. 9
- Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. Ll3da: Visual interactive instruction tuning for omni-3d understanding reasoning and planning. In *CVPR*, 2024e. 7
- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. In *NeurIPS*, 2024. 2
- Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *ICRA*, 2018. 9
- Andrew J Davison, Ian D Reid, Nicholas D Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE TPAMI*, 2007. 9
- Pengxiang Ding, Han Zhao, Zhitao Wang, Zhenyu Wei, Shangke Lyu, and Donglin Wang. Quar-vla: Vision-language-action model for quadruped robots. In *ECCV*, 2024. 10
- Ria Doshi, Homer Rich Walke, Oier Mees, Sudeep Dasari, and Sergey Levine. Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation. In *CoRL*, 2024. 10
- Alberto Elfes. Sonar-based real-world mapping and navigation. *IEEE Journal on Robotics and Automation*, 1987. 9
- Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower models for vision-and-language navigation. In *NeurIPS*, 2018. 9
- Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. Scene-llm: Extending language model for 3d visual understanding and reasoning. *arXiv preprint*, 2024. 7
- 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 *CVPR*, 2022. 6, 7, 9
- Theophile Gervet, Soumith Chintala, Dhruv Batra, Jitendra Malik, and Devendra Singh Chaplot. Navigating to objects in the real world. *SR*, 2023. 9
- Pierre-Louis Guhur, Makarand Tapaswi, Shizhe Chen, Ivan Laptev, and Cordelia Schmid. Airbert: In-domain pretraining for vision-and-language navigation. In *ICCV*, 2021. 9
- Yicong Hong, Cristian Rodriguez, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual entity relationship graph for agent navigation. In *NeurIPS*, 2020. 9
- 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 *CVPR*, 2022. 6, 10
- Yicong Hong, Yang Zhou, Ruiyi Zhang, Franck Dernoncourt, Trung Bui, Stephen Gould, and Hao Tan. Learning navigational visual representations with semantic map supervision. In *ICCV*, 2023a. 6
- Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. In *NeurIPS*, 2023b. 7

- Haifeng Huang, Zehan Wang, Rongjie Huang, Luping Liu, Xize Cheng, Yang Zhao, Tao Jin, and Zhou Zhao. Chat-scene: Bridging 3d scene and large language models with object identifiers. In *NeurIPS*, 2024a. 7
- Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puha Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. In *ICML*, 2024b. 7
- Chioko Sarah Imai, Minghao Zhang, Yuchen Zhang, Marcin Kierebiński, Ruihan Yang, Yuzhe Qin, and Xiaolong Wang. Vision-guided quadrupedal locomotion in the wild with multi-modal delay randomization. In *IROS*, 2022. 10
- Muhammad Zubair Irshad, Chih-Yao Ma, and Zsolt Kira. Hierarchical cross-modal agent for robotics vision-and-language navigation. In *ICRA*, 2021. 10
- Eagle S Jones and Stefano Soatto. Visual-inertial navigation, mapping and localization: A scalable real-time causal approach. *The International Journal of Robotics Research*, 2011. 9
- Aishwarya Kamath, Peter Anderson, Su Wang, Jing Yu Koh, Alexander Ku, Austin Waters, Yinfei Yang, Jason Baldridge, and Zarana Parekh. A new path: Scaling vision-and-language navigation with synthetic instructions and imitation learning. In *CVPR*, 2023. 9
- Simar Kareer, Naoki Yokoyama, Dhruv Batra, Sehoon Ha, and Joanne Truong. Vinl: Visual navigation and locomotion over obstacles. In *ICRA*, 2023. 10
- Liyiming Ke, Xiujun Li, Yonatan Bisk, Ari Holtzman, Zhe Gan, Jingjing Liu, Jianfeng Gao, Yejin Choi, and Siddhartha Srinivasa. Tactical rewind: Self-correction via backtracking in vision-and-language navigation. In *CVPR*, 2019. 9
- Steven M Kearns. Extending regular expressions with context operators and parse extraction. *Software: Practice and Experience*, 1991. 5
- 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*, 2024. 2, 10
- Jacob Krantz and Stefan Lee. Sim-2-sim transfer for vision-and-language navigation in continuous environments. In *ECCV*, 2022. 6, 10
- Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *ECCV*, 2020a. 6, 7, 9, 10
- Jacob Krantz, Erik Wijmans, Arjun Majundar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision and language navigation in continuous environments. In *ECCV*, 2020b. 4, 6
- Jacob Krantz, Aaron Gokaslan, Dhruv Batra, Stefan Lee, and Oleksandr Maksymets. Waypoint models for instruction-guided navigation in continuous environments. In *CVPR*, 2021. 6, 10
- Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In *EMNLP*, 2020. 4, 6, 7, 9
- Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning quadrupedal locomotion over challenging terrain. *SR*, 2020. 6, 10
- Joonho Lee, Marko Bjelonic, Alexander Reske, Lorenz Wellhausen, Takahiro Miki, and Marco Hutter. Learning robust autonomous navigation and locomotion for wheeled-legged robots. *SR*, 2024. 6
- Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah Smith, and Yejin Choi. Robust navigation with language pretraining and stochastic sampling. In *EMNLP*, 2019. 9
- TP Lillicrap. Continuous control with deep reinforcement learning. *arXiv preprint*, 2015. 9

- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. In *MLSys*, 2024a. 20
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models. In *CVPR*, 2024b. 20
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 3
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *CVPR*, 2024. 3, 5
- Junfeng Long, Zirui Wang, Quanyi Li, Liu Cao, Jiawei Gao, and Jiangmiao Pang. Hybrid internal model: Learning agile legged locomotion with simulated robot response. In *ICLR*, 2024. 10
- Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. In *ICLR*, 2019. 9
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *ACL*, 2024. 5
- Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra. Improving vision-and-language navigation with image-text pairs from the web. In *ECCV*, 2020. 9
- Gabriel B Margolis, Ge Yang, Kartik Paigwar, Tao Chen, and Pulkit Agrawal. Rapid locomotion via reinforcement learning. In *Robotics: Science and Systems*, 2022. 6
- Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*, 2022. 6, 10
- Mayank Mittal, Calvin Yu, Qinxi Yu, Jingzhou Liu, Nikita Rudin, David Hoeller, Jia Lin Yuan, Ritvik Singh, Yunrong Guo, Hammad Mazhar, Ajay Mandlekar, Buck Babich, Gavriel State, Marco Hutter, and Animesh Garg. Orbit: A unified simulation framework for interactive robot learning environments. In *RAL*, 2023. 5
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 2015. 9
- Hans Peter Moravec. *Obstacle avoidance and navigation in the real world by a seeing robot rover*. Stanford University, 1980. 9
- Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *International Symposium on Mixed and Augmented Reality*, 2011. 9
- Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. In *ICRA*, 2024. 2
- Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *CVPR*, 2020. 9
- Ram Ramrakhyा, Eric Undersander, Dhruv Batra, and Abhishek Das. Habitat-web: Learning embodied object-search strategies from human demonstrations at scale. In *CVPR*, 2022. 2
- Sonia Raychaudhuri, Saim Wani, Shivansh Patel, Unnat Jain, and Angel Chang. Language-aligned waypoint (law) supervision for vision-and-language navigation in continuous environments. In *EMNLP*, 2021. 6, 7, 9

- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 9
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint*, 2017. 5
- Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. In *NAACL*, 2019. 4, 9
- Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. In *RSS*, 2024. 10
- Sebastian Thrun, Maren Bennewitz, Wolfram Burgard, Armin B Cremers, Frank Dellaert, Dieter Fox, Dirk Hahnel, Charles Rosenberg, Nicholas Roy, Jamieson Schulte, et al. Minerva: A second-generation museum tour-guide robot. In *ICRA*, 1999. 9
- Sebastian Thrun, Dieter Fox, Wolfram Burgard, and Frank Dellaert. Robust monte carlo localization for mobile robots. *Artificial intelligence*, 2001. 9
- Hanqing Wang, Wei Liang, Luc Van Gool, and Wenguan Wang. Dreamwalker: Mental planning for continuous vision-language navigation. In *ICCV*, 2023a. 6
- Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *CVPR*, 2019. 2
- Yikai Wang, Zheyuan Jiang, and Jianyu Chen. Learning robust, agile, natural legged locomotion skills in the wild. In *CoRL*, 2023b. 10
- Zihan Wang, Xiangyang Li, Jiahao Yang, Yeqi Liu, and Shuqiang Jiang. Gridmm: Grid memory map for vision-and-language navigation. In *ICCV*, 2023c. 6
- 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 *CVPR*, 2024. 6
- Zun Wang, Jialu Li, Yicong Hong, Yi Wang, Qi Wu, Mohit Bansal, Stephen Gould, Hao Tan, and Yu Qiao. Scaling data generation in vision-and-language navigation. In *ICCV*, 2023d. 6, 9
- Jonathan Yang, Catherine Glossop, Arjun Bhorkar, Dhruv Shah, Quan Vuong, Chelsea Finn, Dorsa Sadigh, and Sergey Levine. Pushing the limits of cross-embodiment learning for manipulation and navigation. *arXiv preprint*, 2024. 10
- Ruihan Yang, Minghao Zhang, Nicklas Hansen, Huazhe Xu, and Xiaolong Wang. Learning vision-guided quadrupedal locomotion end-to-end with cross-modal transformers. *arXiv preprint*, 2021. 10
- Ruihan Yang, Ge Yang, and Xiaolong Wang. Neural volumetric memory for visual locomotion control. In *CVPR*, 2023. 10
- Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *CVPR*, 2019. 7
- Jiazhao Zhang, Kunyu Wang, Rongtao Xu, Gengze Zhou, Yicong Hong, Xiaomeng Fang, Qi Wu, Zhizheng Zhang, and Wang He. Navid: Video-based vlm plans the next step for vision-and-language navigation. In *RSS*, 2024. 4, 6, 7, 9, 19
- Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and Liwei Wang. Towards learning a generalist model for embodied navigation. In *CVPR*, 2024. 7
- Gengze Zhou, Yicong Hong, Zun Wang, Xin Eric Wang, and Qi Wu. Navgpt-2: Unleashing navigational reasoning capability for large vision-language models. In *ECCV*, 2024. 9

Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Young-jae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale corpus of images interleaved with text. In *NeurIPS*, 2024. [3](#)

Ziyu Zhu, Xiaojian Ma, Yixin Chen, Zhidong Deng, Siyuan Huang, and Qing Li. 3d-vista: Pre-trained transformer for 3d vision and text alignment. In *ICCV*, 2023. [7](#)

APPENDIX TABLE OF CONTENTS

A Ablation Study on Different Data Blends for Training VLA	18
B Ablation Study on Different Memory Size	18
C More Results on VLN-CE-Isaac	18
D Implementation Details for Video Navigation Trajectory Summarization	19
E Implementation Details for Locomotion Motion Policy	19
F Experiments Compute Resources	20
G Parameter-efficient Quantization	20
H Limitations	21

A ABLATION STUDY ON DIFFERENT DATA BLENDS FOR TRAINING VLA

We perform an ablation study to assess the impact of different data blends on training VLA. As shown in Table 6, training navigation data on VLN-CE without label rebalancing leads to a significant drop in performance. Additionally, training VLA exclusively on RxR data demonstrates reasonable cross-dataset performance on R2R-CE, supporting our observations in Table 2. Lastly, we investigate whether excluding RxR dataset degrades R2R-CE performance. We find that the RxR dataset does not significantly contribute to the R2R-CE performance.

Table 6: Results on R2R-CE using different data blends.

	R2R-CE Val Unseen			
	NE ↓	OSR ↑	SR ↑	SPL ↑
without label balancing	7.82	47.5	30.0	25.1
with RxR only	7.57	40.8	31.5	27.8
without RxR	6.11	57.0	47.7	42.4
NaVILA <i>Full Model</i>	5.37	57.6	49.7	45.5

B ABLATION STUDY ON DIFFERENT MEMORY SIZE

We conduct an ablation study to evaluate the impact of memory size (number of history frames) on two tasks: the navigation task using R2R and the spatial understanding task using ScanQA. The results show that for R2R, 8 frames are sufficient to cover most instruction horizons, with limited performance gains from increasing the memory size. In contrast, ScanQA requires a finer-grained memory to recall details such as spatial information of previously seen objects, and performance consistently improves with a larger memory size. Notably, with a memory size of 64 frames, NaVILA outperforms state-of-the-art 3D-LLM (LEO) across all metrics. For real-world experiments, we currently use an 8-frame memory size due to latency constraints. While larger memory sizes could potentially improve performance, we leave it as future work.

C MORE RESULTS ON VLN-CE-ISAAC

Here we show a visualization example highlighting why the Go2 vision policy significantly outperforms the blind policy. As demonstrated in Fig. 7, when encountering an obstacle, the VLA, which is not specifically trained for obstacle avoidance, failed to navigate around it effectively. The blind policy, following the VLA’s commands without additional sensory input, became stuck at the obstacle. In contrast, the vision-based policy, trained to handle obstacles using LiDAR input, can autonomously avoid dangers even when the high-level VLA model does not detect them.



Figure 7: Comparison between Go2 blind policy and vision policy. The blind policy failed to avoid the obstacles and got stuck. The vision policy detected the obstacle and got around to avoid it.

D IMPLEMENTATION DETAILS FOR VIDEO NAVIGATION TRAJECTORY SUMMARIZATION

We provide the data prompts for our auxiliary task of video navigation trajectory summarization. Following the approach in (Zhang et al., 2024), we construct prompt templates that characterize the LLM as a robot designed for navigation. We process the trajectory videos into history frames, insert the frame tokens into the prompt, and ask the LLM to infer the navigation instructions from the video. This task is designed to enhance the robot’s scene understanding and its familiarity with the instruction format.

Assume you are a robot designed for navigation. You are provided with captured image sequences: {frame3}{frame6}{frame9} Based on this image sequence, please describe the navigation trajectory of the robot.

D.1 VLA HYPERPARAMETERS

Please refer to VILA’s paper for details on the hyperparameters used in the first two stages. In the instruction fine-tuning stage, we use a learning rate of $1e^{-4}$ with cosine decay and a warm-up ratio of 0.03. We will release both our training code and data upon paper publication.

E IMPLEMENTATION DETAILS FOR LOCOMOTION MOTION POLICY

Reward and randomization: The reward functions and domain randomization used during Go2 locomotion policy training are listed in Tab. 7 and Tab. 8. The robust policy is trained on flat, rough, slope and obstacle terrains shown in Fig. 8. LiDAR and height map setting are detailed in Tab. 9.

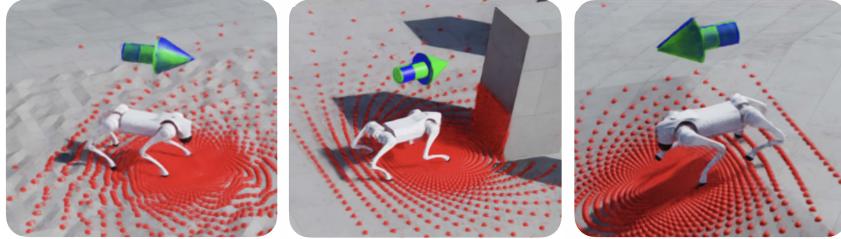


Figure 8: Random rough, obstacle and slope terrain.

Table 7: Rewards

Reward	Expression	Weight
Linear velocity tracking	$\exp(-\ v_{xy}^{\text{cmd}} - v_{xy}\ _2^2)$	1.5
Angular velocity tracking	$\exp(-(\omega_{\text{yaw}}^{\text{cmd}} - \omega_{\text{yaw}})^2)$	1.5
Linear velocity penalty (z)	v_z^2	-2.0
Angular velocity penalty (xy)	$\ \omega_{xy}\ _2^2$	-0.05
Flat orientation	$\ g\ _2^2$	-2.0
Joint accelerations	$\ \ddot{\theta}\ ^2$	-2.5×10^{-7}
Energy	$-\ \tau \dot{q}\ _2^2$	-2×10^{-5}
Body height	$(h^{\text{target}} - h)^2$	-5.0
Feet slipping	$-\ v_{\text{feet}} \cdot \mathbf{1}[F_{\text{feet}} > 1]\ _2$	0.05

Table 8: Domain Randomizations

Parameter	Value
Body Mass	[-3.0, 3.0]
Ground Static Friction	[0.4, 4.0]
Ground Dynamic Friction	[0.4, 4.0]
Motor Strength	[0.9, 1.1]
System Delay	[\$\Delta_t\$, \$\Delta_t\$]

Table 9: Simulation LiDAR and Height Map Settings

Parameter	Value
Channels	32
Vertical Range (degrees)	(0, 90)
Horizontal Range (degrees)	(-180, 180)
Horizontal Resolution (degrees)	4
Voxel Size (m)	0.06
X Range (m)	[-0.8, 0.2]
Y Range (m)	[-0.8, 0.8]
Z Range (m)	[0.05, 0.5]



Figure 9: Obstacle avoidance screenshots. Locomotion policy can ensure collision free in the face of high grass, certain transparent glass and large objects under strong sun lights. The policy presents robustness on sand and grass terrains.

F EXPERIMENTS COMPUTE RESOURCES

NaVILA Training. The first two stages of NaVILA are inherited from VILA (Lin et al., 2024b), which is trained on 16 A100 GPU nodes, with each node having 8 GPUs. The training times for each stage of our 8B model are as follows: connector initialization takes 4 hours, visual language pre-training takes 30 hours. The final visual instruction-tuning stage is experimented on 4 A100 GPU nodes, taking 18 hours.

G PARAMETER-EFFICIENT QUANTIZATION

Currently, NaVILA experiences a noticeable delay of approximately one second during real-world deployment. These delays arise from two parts: image transmission time from Go2 to the server, and the VLA inference time. The transmission time largely depends on the network conditions, while the VLA inference time is approximately 0.6 seconds per sample. To optimize this pipeline, we explore quantization techniques that could reduce both the memory footprint and inference time of NaVILA without sacrificing performance. By applying AWQ (Lin et al., 2024a) to convert the FP16 NaVILA-8B model to W4A16 (low-bit weight-only quantization) format, we achieved impressive gains: memory requirements dropped by half, and processing speed improved by about 40%. Most importantly, navigation capabilities remained robust. These optimizations make NaVILA deployable directly on the robot, which will eliminate image transmission time significantly. We leave on-device deployment as future work.

Table 10: NaVILA quantization results. Tested on RTX 4090 with 1737 context tokens and 10 generated tokens, using a sample from R2R as the test case.

	Computational Cost		RxR Val-Unseen			
	Total Latency (ms) ↓	GPU Memory (GB) ↓	NE ↓	OS ↑	SR ↑	SPL ↑
NaVILA (FP16)	594.58	18.5	5.37	57.6	49.7	45.5
NaVILA (W4A16)	367.80	8.6	5.66	56.8	48.2	43.6

H LIMITATIONS



Walk along the hallway and enter the bedroom.

Figure 10: Failure case of NaVILA.

While our method shows strong performance, we highlight a failure case in the real world where the robot initially follows the prompt but ultimately fails to reach the correct target position. To further enhance performance, improving generalizability is key, and one potential direction is larger-scale training on more realistic simulations. Additionally, image-based vision-language models require a significant number of tokens, which is computationally intensive and limits the number of history frames processed. This challenge could be addressed with recent advancements in long-context LLMs, allowing for more efficient handling of longer sequences. We leave these as future works.