# AME-2: Agile and Generalized Legged Locomotion via Attention-Based Neural Map Encoding

Chong Zhang[123*], Victor Klemm[1], Fan Yang[1], Marco Hutter[1]

Website: sites.google.com/leggedrobotics.com/ame-2

*Abstract*—Achieving agile and generalized legged locomotion across terrains requires tight integration of perception and control, especially under occlusions and sparse footholds. Existing methods have demonstrated agility on parkour courses but often rely on end-to-end sensorimotor models with limited generalization and interpretability. By contrast, methods targeting generalized locomotion typically exhibit limited agility and struggle with visual occlusions. We introduce AME-2, a unified reinforcement learning (RL) framework for agile and generalized locomotion that incorporates a novel attention-based map encoder in the control policy. This encoder extracts local and global mapping features and uses attention mechanisms to focus on salient regions, producing an interpretable and generalized embedding for RL-based control. We further propose a learning-based mapping pipeline that provides fast, uncertainty-aware terrain representations robust to noise and occlusions, serving as policy inputs. It uses neural networks to convert depth observations into local elevations with uncertainties, and fuses them with odometry. The pipeline also integrates with parallel simulation so that we can train controllers with online mapping, aiding sim-to-real transfer. We validate AME-2 with the proposed mapping pipeline on a quadruped and a biped robot, and the resulting controllers demonstrate strong agility and generalization to unseen terrains in simulation and in real-world experiments.

*Index Terms*—Legged Robots, Humanoid and Bipedal Locomotion, Motion Control, Reinforcement Learning.

## I. INTRODUCTION

AGILE and generalized locomotion is essential for legged robots to operate reliably in diverse real-world environments. It requires tightly coupled real-time perception and control, robustness to sensor noise and occlusions, agility through whole-body control, and precise behavior on terrains with sparse footholds. Designing such perceptive locomotion systems that are simultaneously agile and able to generalize across diverse terrains remains a critical challenge.

Classical pipelines address perceptive locomotion by combining model-based control with explicit mapping and state estimation [1]–[8]. Typically, they maintain a state estimator [9], build elevation maps from sensor data [10], [11], and use model-based planning and control methods to command the robot based on this mapping information. These approaches have proven effective for deliberate walking and carefully planned maneuvers on moderately structured terrain, where accurate state estimation and dense, slowly changing maps can be maintained. However, their reliance on precise models and often deterministic optimal control makes them sensitive to errors in estimation and mapping, especially under visual occlusions and violations of modeling assumptions [4]. Elevation maps are often updated at lower rates than the control loop and heuristically filtered [2], [4], [11], which often requires per-terrain tuning to ensure safety and makes it difficult to obtain a general solution for diverse terrains that demand agile motions. Moreover, the optimization and planning components in these pipelines tend to be computationally heavy, which can limit agility on real robots. These limitations motivate learning-based approaches such as RL that can better exploit raw sensor data, operate under uncertainty, and produce more dynamic behaviors.

Learning-based methods, in particular RL, have recently shown strong potential for legged locomotion [12]. One line of work combines RL with classical state estimation and mapping pipelines, using elevation maps or related representations as policy inputs [13]–[26]. These methods can improve robustness and generalization compared to purely model-based controllers, but they also inherit the computational cost and failure modes of the underlying estimation and mapping systems [27], [28]. In simulation, noise is typically injected into both the state estimates and the mapping observations. A common strategy is to assume a fully observed egocentric map during training and then tune the real-world mapping stack to approximate this assumption [13], [15], [21], which simplifies training but can make performance sensitive when occlusions occur or when agile whole-body contacts violate estimator assumptions. Another strategy is to heavily randomize state-estimation and mapping noise so that the controller becomes robust and conservative under mapping uncertainties in dense terrains [14], yet such policies tend to struggle on terrains that require very high agility or precise foothold placement.

More recently, neural network–based mapping has been explored, where networks are trained in simulation to reconstruct maps that can be used as policy inputs [29]–[34]. While this yields efficient mapping suitable for real-time control, the learned mappings are usually fitted to particular terrain distributions and sensor configurations, offer limited generalization to unseen terrains, and do not explicitly model uncertainties such as occlusions.

Another line of work directly maps raw exteroceptive perception and proprioception to actions [35]–[47]. These sensorimotor policies have demonstrated highly agile behaviors on challenging courses [36], [38], [42], but typically exhibit limited generalization beyond the training environments and offer little interpretability, since terrain reasoning is implicitly encoded in the policy. More recent sensorimotor generalist models [48] distill a single policy from multiple expert

[1]Robotic Systems Lab, ETH Zurich, Switzerland
[2]Secure, Reliable, and Intelligent Systems Lab, ETH Zurich, Switzerland
[3]ETH AI Center, Switzerland
*Corresponding: chong.zhang@ai.ethz.ch

Fig. 1. Our method enables agile and generalized legged locomotion across diverse terrains with onboard sensing and computation.

controllers to achieve strong performance, yet such methods still show limited generalization to unseen terrains and often require finetuning on new terrains, as evaluated in our benchmarks (Sec. VII). Overall, existing RL-based approaches tend to trade off agility, generalization, mapping efficiency, and interpretability rather than addressing these requirements within a unified framework.

In this work, we present AME-2, a unified RL framework for perceptive legged locomotion with an **A**ttention-based **M**ap **E**ncoding architecture trained jointly with the controller. The AME-2 encoder, built upon the design in [15], first extracts local features (pixel-wise representation of terrain details) and global features (capturing global terrain context) from the elevation map, then uses the global features together with proprioception to assign attention weights to the local features. The resulting weighted local features are concatenated with the global features and proprioception to form a terrain-aware representation for policy learning. This design allows the policy to downweigh local regions that are less relevant for the current task, improving generalization to new terrains. Because the attention is conditioned on global terrain context, the policy can also learn distinct attention and motion patterns across different terrains. To enable learning agile controllers with this representation, we adopt a goal-reaching locomotion formulation from prior works [20], [26], [32], [40], [49], [50] and adapt it so that the same reward functions and training settings can be used for different robots. Together, the AME-2 encoder and unified training formulation enable training a single generalist policy with terrain-aware skills, which in our experiments exhibits strong generalization to unseen terrains while maintaining high agility.

To enable real-world deployment of agile and generalized locomotion controllers, we develop a learning-based elevation mapping pipeline to remove the dependency on classical mapping stacks while maintaining generalization. The pipeline projects depth images into local grids and uses a lightweight neural network trained via Bayesian learning [51] to predict local elevations with per-cell uncertainty estimates. These local maps are fused into a consistent global frame with odometry, providing a fast and uncertainty-aware representation that accounts for occlusions and sensor noise. We then query egocentric elevations and associated uncertainties as the map inputs for the controller. This design is inspired by [52], which fuses per-frame map predictions instead of training temporal networks for off-road navigation, thereby reducing data requirements and mitigating overfitting. Furthermore, we synthesize random terrain data, including many terrains that are not feasible for the robots, to broaden the training distribution for the local map predictor and improve generalization.

The same mapping pipeline runs both in parallel simulation and on the real robots, enabling online mapping during training rather than relying on idealized or hand-tuned maps. For training efficiency, we adopt a teacher–student scheme [48], [53]: a teacher policy is first trained with ground-truth elevation maps, and a student policy with the same AME-2 architecture is then trained using the proposed mapping under teacher supervision alongside RL. The resulting student policy is directly deployable, operates with the lightweight, uncertainty-aware mapping in the loop, and retains the generalization and agility of the teacher.

We evaluate AME-2 on an ANYmal-D quadruped [54] and a LimX TRON1 biped [55] both in simulation and on real hardware. As shown in Fig. 1, the resulting controllers exhibit strong agility and robust generalization to a wide range of terrains. The proposed mapping pipeline produces high-quality elevation maps in real time, supporting these agile motions with onboard perception and computation.

The main contributions of this work are summarized as

follows:

1) We propose **AME-2**, a unified RL framework with an attention-based map encoder to achieve agile and generalized legged locomotion.
2) We develop a **lightweight, uncertainty-aware elevation mapping** pipeline that generalizes to diverse terrains, explicitly models occlusions and noise, and helps bridge the sim-to-real gap.
3) We use a **teacher-student scheme** that yields deployable controllers using the learned mapping while retaining agility, generalization, and interpretability.
4) We demonstrate a **state-of-the-art combination of agility and generalization** with a quadruped and a biped robot across diverse challenging terrains under the same training setup.

## II. RELATED WORKS

### A. Perceptive Locomotion with Explicit Mapping

Early explorations of perceptive locomotion use model-based methods with either prebuilt maps or classical online mapping. Most of them first perform foothold or trajectory planning and then apply a model-based tracking controller [1]–[3], [5]–[8], while some directly optimize a reactive controller over feasible footholds [4]. Although these methods offer strong guarantees and can generalize well under their modeling assumptions, their real-world performance is limited by model mismatch, state-estimation uncertainties, mapping errors, and computational burden.

To address these limitations, RL emerges as an alternative for perceptive locomotion [12]. One line of work keeps a strong model-based component and uses RL mainly as an add-on [21], [22], [56]. DTC [21] augments a model-based planner with an RL tracking controller to improve robustness while maintaining generalization, but the planner remains sensitive to mapping uncertainties and constrains agility through its modeling assumptions and computational cost. RLOC [22] instead uses a model-based controller to track an RL planner, yet the overall performance is capped by the limitations of the controller.

Other approaches remove the model-based modules entirely while still relying on explicit maps. When given prebuilt maps, RL controllers can exhibit highly agile motions [26]. With classical online mapping, they outperform model-based counterparts in terms of agility and robustness [13]–[19], [23], [24], but their agility remains limited by the update rate and quality of the mapping. In addition, many such controllers suffer from generalization issues, as they tend to overfit to the training terrains [21]. A recent work [15] shows that RL can achieve generalized locomotion across diverse terrains, but its agility is still constrained by the online classical mapping stack and the learning framework.

More recently, RL-based methods have begun to use learned maps, where neural networks reconstruct egocentric terrains for the controller in real time [29]–[34]. These approaches demonstrate strong agility and successful real-world deployment, but the learned maps often fit the training distribution and show limited generalization in unseen, unstructured

environments. In contrast, our work combines a fully RL-based controller with a learned mapping module designed for generalization. As a result, our policies achieve agility comparable to previous state of the art while enabling generalized locomotion across diverse terrains.

### B. Perceptive Locomotion with Raw Sensor Data

Perceptive locomotion with raw sensor data seeks to bypass the limitations of explicit mapping. Instead of constructing and maintaining a map, some works directly train policies from camera images [35]–[39], [41]–[47] or lidar point clouds [40]. Early approaches in this direction still separate perception and control: a visual module outputs footholds or trajectories, which are then tracked by model-based controllers [57]. Between explicit mapping and fully end-to-end policies, other works learn intermediate representations from raw sensors (such as ray distances [49] and neural volumetric memory [58]) that remain cheap to compute while enabling fast control.

Many recent methods adopt a monolithic design, where a single neural network maps raw sensor streams and proprioception directly to joint commands. Within this family, some distill deployable student policies from privileged teacher policies — often using maps as privileged inputs — to achieve efficient and stable learning [35], [36], [38], [42], [43], [46], [48], while others learn from scratch with the help of representation learning techniques such as world models [44] or privileged-information reconstruction [41], [45]. The current state of the art in agile generalist locomotion [48] trains multiple mapping-based teacher controllers and distills them into a single generalist student policy that operates directly on raw depth images. This student policy combines the skills of the teachers and can be finetuned efficiently on new terrains. However, its zero-shot generalization to previously unseen environments is limited.

In contrast, our work also trains neural networks from raw sensor data, but uses them exclusively within the mapping process rather than in a monolithic perception-to-action policy. Because our learned mapping is lightweight and well suited to massive parallelization in simulation, we can run the full mapping loop during training and then deploy exactly the same mapping at test time, preserving both speed and consistency. Without training a monolithic policy that directly maps raw sensor data to joint actions, we still obtain an agile controller while also achieving strong generalization across diverse terrains.

### C. Learning-Based Mapping

Learning-based mapping aims to retain the structure and interpretability of explicit maps while overcoming the computational and modeling limitations of classical mapping. Classical mapping pipelines combine geometric projection, hand-tuned filtering, and probabilistic fusion to produce high-quality maps, and can handle a wide range of complex terrains [10], [11], [23], [24]. However, these methods are often too slow for agile locomotion and require heuristic, terrain-specific filtering to cope with occlusions and sensor noise.
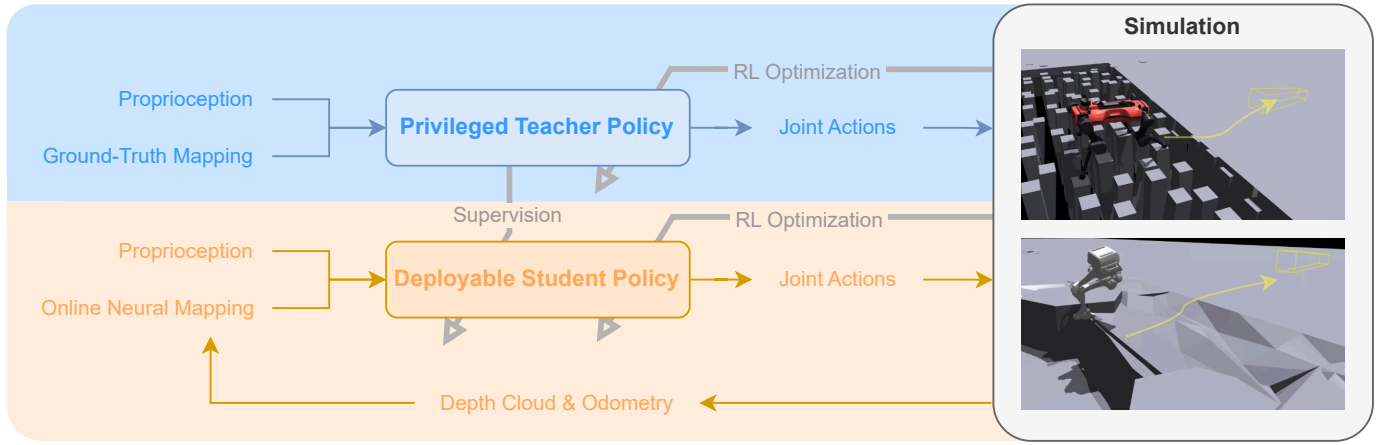
Fig. 2. An overview of our system. We use RL to train a teacher policy with ground-truth mapping in simulation, and a student policy under the teacher's supervision with our proposed neural mapping. The policies output joint-level actions that actuate the robots to reach position and heading goals, as illustrated by the yellow wireframes.

To enable faster control, several works directly train terrain reconstruction modules in or after the policy training loop [29]–[34]. These learned mappings can infer local geometry quickly enough to support agile controllers, but are usually trained on a narrow distribution of environments, showing limited generalization beyond the training terrains.

Recent off-road navigation approaches take a different route to obtain uncertainty-aware neural mapping with improved accuracy and data efficiency. In particular, [52] leverages neural processes [59] to model terrain elevation: a network predicts per-frame local elevation maps together with associated uncertainties, which are then fused over time using odometry. By predicting local maps with uncertainty instead of learning a fully end-to-end sensor-to-map model, this approach reduces the data requirements for accurate estimation and improves generalization in complex off-road environments. However, the computational costs of the models in [52] are too high (each environment takes > 13 GB GPU memory) and cannot be deployed with thousands of parallel simulation environments.

Inspired by these ideas, we design a learning-based elevation mapping module tailored for agile legged locomotion. We project depth-camera point clouds into local elevation grids and train lightweight, robust networks to predict local elevation and its uncertainty from these inputs, which are then fused over time using odometry. This design suppresses sensor noise and naturally models occlusions as regions of high uncertainty. We further improve generalization through extensive randomization in a custom data pipeline. The resulting module is lightweight enough to run in massively parallel simulation and fast enough to support highly agile motions in deployment, while providing an uncertainty-aware terrain representation for our controllers.

## III. SYSTEM OVERVIEW

### A. Formulation

As illustrated in Fig. 2, we use RL to train terrain-aware locomotion policies that reach position and heading goals. We formulate the problem as a partially observable Markov decision process (POMDP) and optimize the policies using PPO [60] in parallel simulation [61]. By maximizing goal-reaching rewards (detailed in Sec. IV), the policies learn agile locomotion skills that enable them to traverse diverse terrains. We first train a privileged teacher policy and then transfer the skills to a deployable student policy, which is also detailed in Sec. IV.

### B. Controller Inputs and Outputs

Our controllers run at 50 Hz. The policy actions $a$ are joint PD targets [62] tracked at 400 Hz on real hardware.

Proprioception observations include base linear velocity $v_b$ (only accessible to the teacher), base angular velocity $\omega_b$, projected gravity $g_b$, joint positions $q$, joint velocities $\dot{q}$, previous actions, and goal commands $c$.

Ground-truth mapping represents each point in the ego-centric elevation grid by its 3D coordinates $(x, y, z)$. Our neural mapping augments this with an uncertainty channel and produces a 4D representation $(x, y, z, u)$ for each point, where $u$ is the uncertainty metric.

## IV. LOCOMOTION CONTROL

### A. Policy Architecture and AME-2 Encoder

The general policy architecture is illustrated in Fig. 3, with the AME-2 encoder serving as the core feature extractor for the mapping input. We have a proprioception encoder to embed proprioceptive observations, and the AME-2 encoder to embed the map observations based on the proprioception embedding. We feed both the proprioception embedding and the map embedding through a multilayer perceptron (MLP) to output the actions.

The AME-2 encoder first extracts local map features with a convolutional neural network (CNN) and computes a positional embedding for each point with an MLP. These are then fused by another MLP to obtain pointwise local features. Next, an additional MLP followed by max pooling over the pointwise features produces global features that capture the overall terrain context. We combine these global features with a proprioceptive embedding through an MLP to obtain a query vector, which is used in a multi-head attention (MHA)
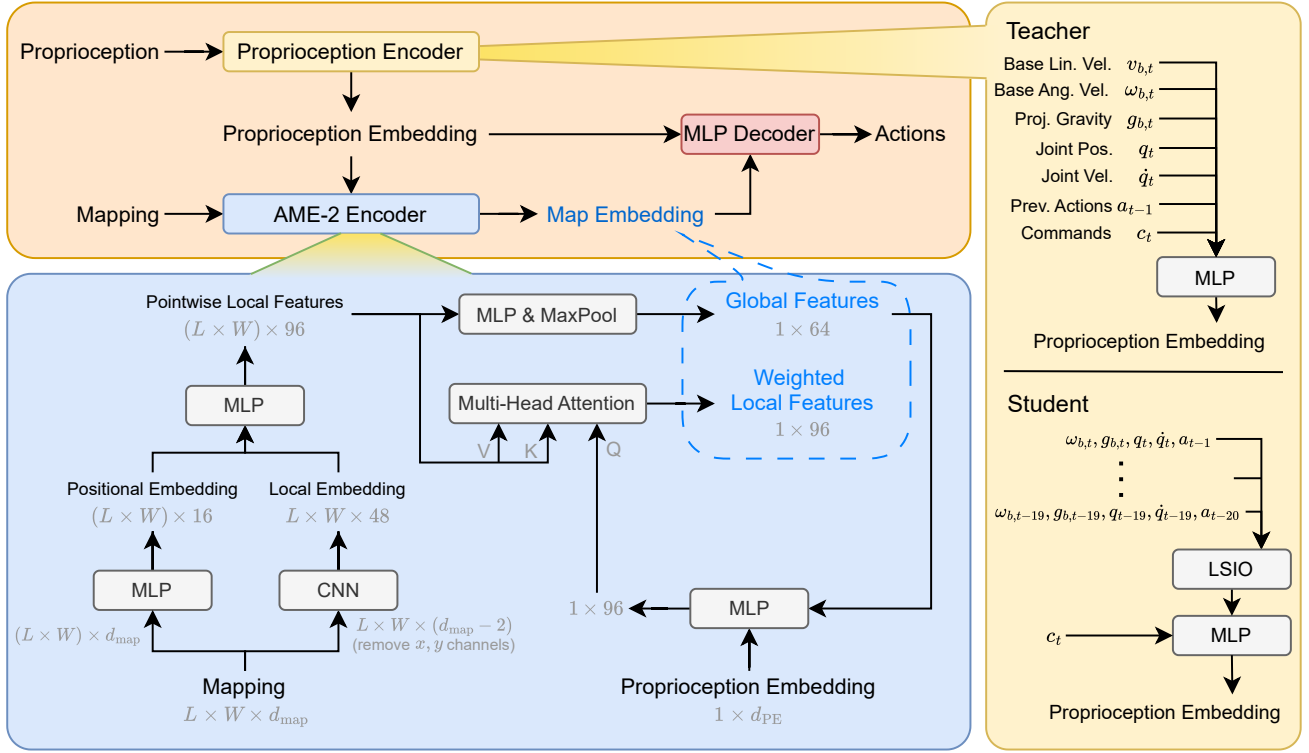
Fig. 3. An illustration of our AME-2 policy architecture. **Left top:** The general abstract of our policies. Proprioceptive observations are encoded into a proprioception embedding, which is used together with the mapping as the inputs of the AME-2 encoder to produce the map embedding. The proprioception embedding and the map embedding are then concatenated and fed into an MLP to generate actions. **Left bottom:** The AME-2 encoder design. We extract both local features and global features from the map, and then use the global features and the proprioception embedding to produce the attention-weighted local features. $L$ and $W$ are the length and width of the map, $d_{map}$ is the dimension of map representation (3 for the teacher, 4 for the student), $d_{PE}$ is the dimension of proprioception embedding. **Right:** The proprioception encoder designs. We have different designs for teacher and student policies to facilitate sim-to-real transfer, while both designs can fit into our overall architecture.

module [63] with the pointwise local features serving as keys and values. This yields a weighted local feature embedding that focuses on important terrain regions based on the current proprioceptive state and global context. Finally, the global features and weighted local features are concatenated to form the map embedding fed into the policy's action decoder.

This design builds on the attention-based map encoder in previous work [15], which we refer to as AME-1 for clarity. Unlike the AME-1 encoder, our AME-2 encoder additionally computes global features and uses them to weigh the local features. This enables more generalized locomotion over complex terrains, where motion patterns should vary with the terrain. We demonstrate the resulting performance gap between AME-1 and AME-2 encoders in Sec. VI and Sec. VII.

The choice of proprioception encoder in our architecture is flexible, as long as it produces an effective proprioception embedding. Hence, we use different designs for the teacher and the student. For the teacher, since ground-truth proprioceptive observations are available, we use a plain MLP over them to obtain the proprioception embedding. For the student, because of various environmental uncertainties and following lessons from previous works [64]–[66], we stack the proprioceptive observations (except base linear velocities and commands) from the past 20 steps and use Long-Short I/O (LSIO) [66] to obtain a temporal embedding of both robot states and environment dynamics. This temporal embedding and the commands are then jointly fed into an MLP to produce

the student's proprioception embedding.

### B. Asymmetric Actor Critic

We adopt asymmetric actor–critic training [67]. For the critic, we do not use the same attention-based design as in the actor, because the critic does not need to generalize beyond the training terrains and optimizing an MHA module with $(L \times W)$ local feature inputs is costly. Instead, we use the mixture-of-experts (MoE) design from [68], which is powerful for function fitting yet much more computationally efficient to optimize. Note that, although we can train a generalized teacher with our policy architecture and an MoE critic, an MoE actor does not yield a generalized teacher, as shown in Sec. VII.

On top of the noiseless proprioceptive observations used by the teacher actor and the ground-truth map, we additionally provide the contact state of each link to the critic. We do not provide such contact information to the teacher, since we can already train a robust, generalized teacher without it, and a larger information gap between teacher and student can hurt student training [69]. We also apply left-right symmetry augmentation from [32] to the critic to improve sample efficiency and motion style, but not to the actor to avoid additional computational cost of actor optimization. The critic design is shared for the teacher and the student.

TABLE I
REWARDS FOR CONTROLLER TRAINING.

| Reward Term | Expression | Weight ($\times d\tau^*$) | Notes |
|---|---|---|---|
| *Task Rewards* | | | |
| Position Tracking | Eq. (1) | 100 | |
| Heading Tracking | Eq. (3) | 50 | |
| Moving to Goal | Eq. (4) | 5 | |
| Standing at Goal | Eq. (5) | 5 | |
| *Regularization and Penalties* | | | |
| Early Termination | 1 if early termination triggered | $-10/d\tau$ | |
| Undesired Events | 1 for each undesired event | $-1$ | |
| Base Roll Rate | $[\omega_b]_x^2$ | $-0.1$ | |
| Joint Regularization | $\|\dot{q}\|^2 + 0.01\|\tau\|^2 + 0.001\|\ddot{q}\|^2$ | $-0.001$ | $\tau$ denotes the joint torques |
| Action Smoothness | $\|a_t - a_{t-1}\|^2$ | $-0.01$ | |
| Link Contact Forces | $\|\max(F_{\text{con}} - G, \mathbf{0})\|^2$ | $-0.00001$ | $F_{\text{con}}$ denotes the contact forces for each link and $G$ is the robot weight |
| Link Acceleration | $\sum_l \|\dot{v}_l\|$ | $-0.001$ | $v_l$ is the velocity of link $l$; summing over all links |
| *Simulation Fidelity* | | | summing over all joints |
| Joint Position Limits | $\sum_j \max(0, q_j - 0.95q_j^{\max}, 0.95q_j^{\min} - q_j)$ | $-1000$ | $q_j^{\max}$ and $q_j^{\min}$ are position limits for joint $j$ |
| Joint Velocity Limits | $\sum_j \max(0, |\dot{q}_j| - 0.9\dot{q}_j^{\max})$ | $-1$ | $\dot{q}_j^{\max}$ is the velocity limit for joint $j$ |
| Joint Torque Limits | $\sum_j \max(0, |\tau_j| - 0.8\tau_j^{\max})$ | $-1$ | $\tau_j^{\max}$ is the torque limit for joint $j$ |

$^*d\tau = 0.02$ sec is the policy interval.

## C. Teacher-Student RL

We use Teacher–Student RL [14], [48], [53], [70] to facilitate sim-to-real transfer. Although directly training a student policy with our mapping pipeline is also possible, simulation can run at only about half the speed compared to using ground-truth mapping and requires significantly more GPU memory, making it computationally inefficient for us.

Our student training objective linearly combines the RL losses from PPO, the action distillation losses from [53], and a representation loss given by the mean squared error between the teacher and student map embeddings. In practice, we also disable the PPO surrogate loss during the first few iterations while using a large learning rate. These design choices enable stable, efficient, and well-aligned student training, which we will ablate in Sec. VII.

## D. Environments

*1) Rewards:* We use three types of rewards shared across both quadruped and biped training environments, as listed in Table I: task rewards which incentivize goal-reaching behaviors, regularization rewards which improve stability and safety, and simulation-fidelity rewards, which penalize near-limit joint states that can cause unrealistic simulation.

There are four task reward terms, respectively for goal position tracking, goal heading tracking, moving towards goal, and standing at goal. Following [49], we define the position tracking reward as

$$r_{\text{position\_tracking}} = \frac{1}{1 + 0.25d_{xy}^2} \cdot t_{\text{mask}}(4), \quad (1)$$

where $d_{xy}$ is the horizontal distance from the robot to the goal position, and $t_{\text{mask}}(\cdot)$ is a time-based mask function defined as

$$t_{\text{mask}}(T) = \frac{1}{T} \cdot \mathbf{1}(t_{\text{left}} < T), \quad (2)$$

with $t_{\text{left}}$ denoting the remaining time of the current episode. Further, we define the heading tracking reward as

$$r_{\text{heading\_tracking}} = \frac{1}{1 + d_{\text{yaw}}^2} \cdot t_{\text{mask}}(2) \cdot \mathbf{1}(d_{xy} < 0.5), \quad (3)$$

Intuitively, these two tracking rewards encourage the robot to be at the goal position with the desired heading at the end of the episode, without constraining how it reaches the goal, thereby allowing complex locomotion skills to emerge for traversing terrains.

To further facilitate exploration, we introduce a moving-to-goal reward:

$$r_{\text{move}} = \mathbf{1}\Big(d_{xy} < 0.5 \vee \\ \big(\cos\theta_{v_b,\text{goal}} > 0.5 \wedge v_{\min} \leq \|[v_b]_{xy}\| \leq v_{\max}\big)\Big), \quad (4)$$

where $\theta_{v_b,\text{goal}}$ is the angle between the base velocity and the vector from the base to the goal position, $v_{\min} = 0.3\,\text{m/s}$ is the lower bound of the horizontal base velocity $[v_b]_{xy}$ to be considered as moving, and $v_{\max} = 2\,\text{m/s}$ is an upper bound chosen according to hardware and SLAM constraints. In other words, $r_{\text{move}} = 1$ if the robot is already close to the goal or if it is moving roughly towards the goal, and $r_{\text{move}} = 0$ otherwise.

To enforce a stable standing posture after reaching the goal, we define the standing reward:

$$r_{\text{stand}} = \mathbf{1}(d_{xy} < 0.5 \wedge d_{\text{yaw}} < 0.5) \\ \cdot \exp\left(-\frac{d_{\text{foot}} + d_{\text{g}} + d_{\text{q}} + d_{xy}}{4}\right), \quad (5)$$

where $d_{\text{foot}}$ is the number of feet not in contact divided by the total number of feet, $d_{\text{g}} = 1 - [g_b]_z^2$ measures the base tilt relative to gravity, and $d_{\text{q}}$ is the mean deviation of joint positions from the standing reference. This reward ensures the robot maintains a static, upright configuration at the goal.

The regularization and simulation fidelity rewards are presented in Table I. We set all weights to integer powers of 10, and make each term's contribution $1 \sim 2$ orders of magnitude smaller than the task rewards. This enables successful learning without extensive weight tuning.

In the regularization rewards, we penalize each occurrence of these undesired events:

- Spinning too fast: Yaw rate $|[\omega_b]_z| > 2.0\,\text{rad/s}$, which can trigger drifts in the odometry.
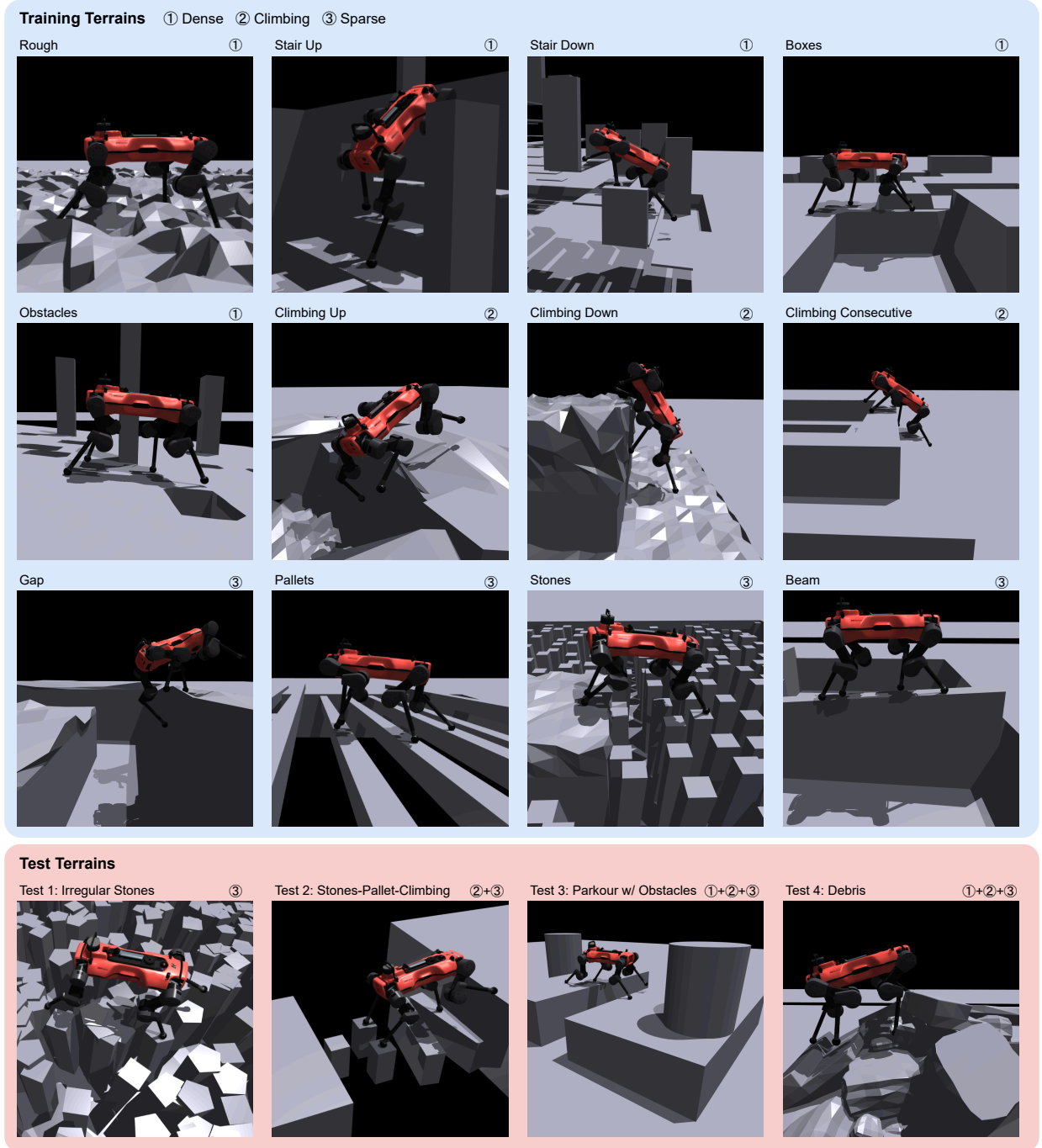
Fig. 4. Training and Test Terrains. We train our locomotion policies on primitive terrains of three categories: ① dense, ② climbing, and ③ sparse. To test generalization, we evaluate the policies on four terrains that are either entirely unseen or combinations of training and/or unseen terrains.

- Leaping on flat terrain: All feet are off the ground when the elevation difference is $< 30$ cm.
- Non-foot contacts: We penalize both general non-foot contacts and non-foot contact switches (from no contact to contact). We separate these to allow stable interactions when necessary (e.g., climbing) while discouraging unnecessary non-foot contacts.
- Stumbling: Any link having a horizontal contact force larger than the vertical force.
- Slippage: Any link moving while in contact.
- Self-collision: Collisions between robot links.

Unlike previous works [19], [23], [33], we do not explicitly reward or penalize foot contact positions in sparse terrains; instead, we formulate rewards for all robot links to enable emergent whole-body contacts. Those foothold position rewards are also difficult to define in a terrain-agnostic way: for example, quadruped climbing benefits from active knee contacts and near-edge foot placements, as shown in Fig. 10, whereas on sparse terrains near-edge foot placements are undesirable.

*2) Termination:* To facilitate training, we set the following early termination conditions:

- Bad orientation: The projected gravity vector satisfies $|[g_b]_x| > 0.985$, $|[g_b]_y| > 0.7$, or $[g_b]_z > 0.0$ (robot flipped).
- Base collision: The base link experiences a contact force larger than the robot's total weight.
- High thigh acceleration: The acceleration of any thigh link exceeds a specific threshold. We design this to reduce impact during jumping, as stiff landings can damage hardware or reduce lifespan. Interestingly, we use typical data from dog (60 m/s$^2$ [71]) and human (100 m/s$^2$ [72]) biokinetics literature, and they directly work to emerge impact absorption or avoidance behaviors, which is shown in Sec. VI.
- Stagnation: Movement in the past 5 s is less than 0.5 m while the robot is still $> 1$ m away from the goal position.

Triggering any of these conditions terminates the episode immediately and incurs early termination penalties.

*3) Terrains and Curriculum:* We train our controllers on primitive terrains and evaluate generalization on complex test terrains. Both are shown in Fig. 4, while quantitative results are presented in Sec. VII. We use a terrain curriculum that scales difficulty from easy to hard (detailed in Appendix A), similar to [26]. To stabilize learning across different challenging terrains, we estimate the robot's success rate using an exponential moving average. If the robot reaches the goal and its estimated success rate exceeds 0.5, we promote the environment to the next difficulty level. Conversely, the level is demoted if the robot remains $> 4$ m away from the goal (which is the average starting distance). Once the robot passes the highest level, the environment resets to a random difficulty.

Additionally, we use a perception noise curriculum and an initial heading curriculum. In the first 20% iterations of teacher training, we linearly increase mapping noise from zero to the maximum level, and simultaneously expand the initial heading from facing the goal to a random yaw in $[-\pi, \pi]$.

*4) Domain Randomization:* We employ domain randomization [73] during training to enhance robustness and facilitate sim-to-real transfer. We apply the following randomization setup:

- Robot Dynamics: We uniformly randomize the payload, friction coefficients, and actuation delays. For the TRON1 biped, we additionally randomize PD gains and motor armatures, following [15].
- Observation Noise: We apply uniform noises to the policy observations. We also degrade the student's depth clouds by simulating missing points and sensor artifacts.
- Mapping: During student training, we randomly select a subset of environments to access complete maps, while others rely on partial online maps. This setup enables map reuse when the robot traverses the same terrain repeatedly, as illustrated in Fig. 5. We also corrupt the student's mapping by randomly removing points and assigning random height values with high uncertainties. Additionally, we simulate mapping drifts for both the teacher and the student.

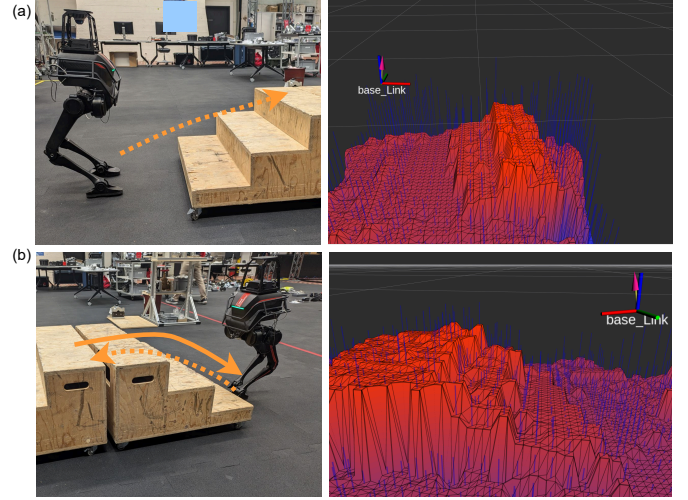Implementation details are provided in Appendix B.



Fig. 5. The controller can reuse the built map for the same terrain, enabled by our mapping randomization. (a) When the robot moves in the direction of the dotted orange arrow, it observes a partial map of the staircase (visualized on the right). (b) After going down along the solid orange arrow, it perceives a complete map of the staircase (visualized on the right) when attempting to move back along the dotted orange arrow.

*E. Training and Deployment*

*1) Training Setup:* We train our controllers in Isaac Gym [61] with the PPO implementation in RSL-RL [74]. The both teacher policies are trained with 80000 iterations, and the both student policies are trained with 40000 iterations (surrogate loss disabled during first 5000 iterations). The hyperparameters are listed in Appendix C. The training cost for the ANYmal-D policies is $\sim 60$ RTX-4090-days with 8 GPUs in parallel. The training cost for the TRON1 policies is $\sim 30$ RTX-4090-days with 4 GPUs in parallel.

The difference between both robots' training costs are due to the mapping sizes. For ANYmal-D, we use a map of size $36 \times 14$, with 8-cm resolution, centered at $x = 0.6$ m, y $= 0$ m in the base frame. For TRON1, we use a map of size $18 \times 13$, with 8-cm resolution, centered at $x = 0.32$ m, y $= 0$ m in the base frame. These numbers are designed based on the robot dimensions and terrain sizes.

*2) Deployment Setup:* We deploy the controllers using ONNX Runtime [75]. On the onboard Intel Core i7-8850H CPUs, the policy inference time is approximately 2 ms.

*3) Other Sim-to-Real Designs:* We model the actuator dynamics for sim-to-real transfer. In simulation, we use the actuator network [76] for ANYmal-D, and an identified DC-motor model for TRON1. We clip the applied torques using the joint torque-velocity constraints in [77].

To enable continuous deployment (infinite-horizon execution), we define the command observations differently from prior goal-reaching works [48]–[50]. Following prior practice, the critic receives the full command: the goal's relative position, the sine and cosine of the relative yaw, and the remaining episode time. The actor, however, receives a modified command representation: we clip the observed goal distance to a maximum of 2 m and remove the remaining time. Additionally,
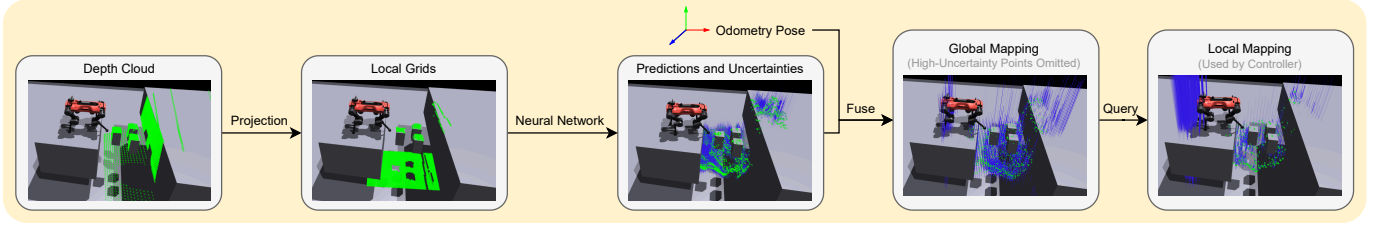
Fig. 6. Our proposed mapping pipeline. For each frame of depth clouds, we project it into local grids, and use a lightweight neural network to predict elevation estimations with uncertainties. The predictions are then fused into the global map via the odometry, and local maps can be queried from the global map to serve as the controller inputs. In visualization, the green points indicate the depth points or elevation estimations, and the blue lines indicate the uncertainties.
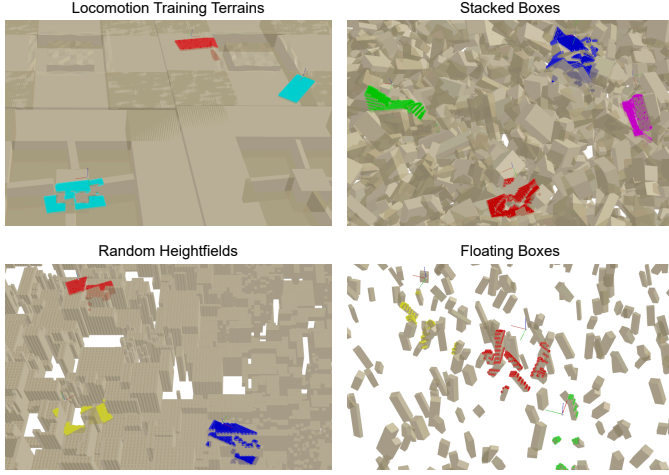


Fig. 7. Terrain meshes used for mapping model training. We use four meshes: the locomotion training terrains, randomly stacked boxes, random heightfields, and random floating boxes. The colored points visualized are the elevation scans from random poses.

if the actual goal distance exceeds $2\,\mathrm{m}$, we randomize the observed yaw command during training. This design decouples policy deployment from finite-horizon training and simplifies steering control.

## V. NEURAL MAPPING

### A. Mapping Pipeline

Our proposed neural mapping pipeline is illustrated in Fig. 6. It can not only achieve real-time computation on the hardware, but also run with thousands of parallel environments in simulation, thereby bridging the sim-to-real gap with identical mapping pipelines.

For each sensing frame, we project the point cloud onto a local 2D height grid. If multiple points fall into the same cell, we keep the maximum $z$-value, which is most relevant for locomotion, and we assign a fixed minimum value to cells with no points. The resulting local elevations, however, are often noisy and incomplete due to occlusions. To mitigate this, we use a lightweight CNN trained with Bayesian learning (detailed in Sec. V-B) to jointly predict base-relative elevations and their uncertainties (in the form of log-variance). The predicted uncertainties capture both measurement noise and occlusions, while the elevation predictions themselves can also suppress noise.

We then use the odometry poses to fuse these local predictions into a global grid map $\mathcal{M}$. It has two layers: the elevation

layer and the uncertainty layer (in the form of variance). The global map is initialized as flat ground at the robot's standing height (ground height relative to the base) with large uncertainties. At each new frame, given the local elevation estimations with uncertainties and the current base pose, we project the local predictions to the global grid cells. For any global grid point $(u, v)$ covered by the local grids, we denote the new estimation as $h_t$ and its uncertainty as $\sigma_t^2$, fusing them with existing values $h_{prior}$ and $\sigma_{prior}^2$.

We do not use standard Bayesian fusion because repeated observations of the same occluded or uncertain area should not have reduce uncertainties simply due to consistent predictions. Instead, we employ a *Probabilistic Winner-Take-All* strategy. First, we calculate an effective measurement variance $\hat{\sigma}_t^2$ that is lower-bounded by the prior to prevent over-confidence:

$$\hat{\sigma}_t^2 = \max(\sigma_t^2, 0.5 \cdot \sigma_{prior}^2). \tag{6}$$

An update is considered valid only if the effective measurement variance is not significantly larger than the prior ($\hat{\sigma}_t^2 < 1.5\sigma_{prior}^2$), or if the absolute uncertainty is low ($\hat{\sigma}_t^2 < 0.2^2$).

For valid updates, we determine the probability $p_{\mathrm{win}}$ of overwriting the map based on the relative precision:

$$p_{\mathrm{win}} = \frac{(\hat{\sigma}_t^2)^{-1}}{(\hat{\sigma}_t^2)^{-1} + (\sigma_{prior}^2)^{-1}}. \tag{7}$$

Finally, the map is updated stochastically. We sample $\xi \sim \mathcal{U}[0, 1]$ and let the new prediction take over the cell if the sample falls within the probability threshold:

$$(h_{new}, \sigma_{new}^2) \leftarrow \begin{cases} (h_t, \hat{\sigma}_t^2) & \text{if } \xi < p_{\mathrm{win}}, \\ (h_{prior}, \sigma_{prior}^2) & \text{otherwise.} \end{cases} \tag{8}$$

For the controller inputs, we then just query the grids around the robot pose in the latest global map.

This *Probabilistic Winner-Take-All* strategy offers the following benefits:

- The uncertainty of the same occluded point will not decrease through consistent predictions.
- Over-confident predictions, if not consistent, cannot take over the cell.
- The system can rapidly update the map in response to dynamic terrain changes when high-confidence measurements are available.
- The pipeline is easy to integrate with parallel simulation, and fast enough to run on the hardware.
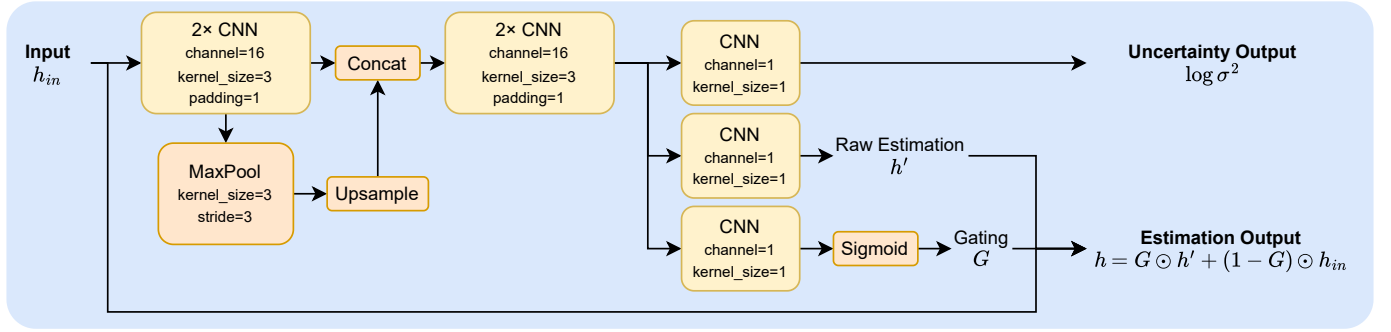
Fig. 8.  Mapping model architecture. We use a lightweight U-Net [78] model with a gated residual design for the estimation.

## B. Training

We train a CNN to predict per-frame elevation estimates and uncertainties with synthetic data and random terrains. The technical implementations are detailed below.

*1) Terrains and Data Sampling:* We use both the locomotion training terrain meshes and additional procedurally generated terrains to train the mapping model, as shown in Fig. 7. Instead of relying on physical simulation, we directly sample local elevation grids from random poses above these meshes using raytracing with Warp [79], enabling sampling at hundreds of thousands of frames per second on a single RTX 4090 GPU.

*2) Training Data Synthesis:* We apply the following augmentations to the sampled local elevation grids:

- additive uniform noise with random magnitude on each cell;
- random cropping of the map from the four borders;
- simulated occlusions with random sensor positions and field of view;
- random ranges to clip the elevations;
- random missing points and outliers at random ratios.

By doing so, we synthesize diverse, noisy, and partially observable local grids as the model inputs, and the original ground-truth elevations as the labels.

*3) Model and Optimization:* We train the model to reconstruct the ground-truth elevations with $\beta-$NLL loss ($\beta = 0.5$) from [80]:

$$L_{0.5} = \mathbb{E}_{X,Y}\left[\text{sg}\left[\hat{\sigma}(X)\right]\left(\frac{\log \hat{\sigma}^2(X)}{2} + \frac{(Y - \hat{\mu}(X))^2}{2\hat{\sigma}^2(X)}\right)\right],$$
(9)

where $X$ denotes the inputs, $Y$ denotes the ground-truth elevations, and $\hat{\mu}(X)$ and $\hat{\sigma}^2(X)$ denote the predicted estimation and variance, respectively. The operator $\text{sg}[\cdot]$ denotes the stop-gradient operation. Compared to the standard negative log-likelihood (NLL) loss used in classical Bayesian learning [51], this formulation reduces the tendency of the model to overestimate uncertainty on hard samples to trivially reduce the loss. It encourages the model to output high uncertainty when accurate predictions are not possible, and low uncertainty together with accurate predictions when they are, thereby capturing the noise and occlusions.

For batched optimization, terrain roughness can vary significantly across samples. As a result, flat terrains can dominate the batch loss and reduce the effective emphasis on challenging

cases. To mitigate this, we re-weigh samples in each batch during training using their total variation (TV) [81]:

$$\text{TV}(Y_b) = \frac{1}{HW}\left(\|\nabla_x Y_b\|_1 + \|\nabla_y Y_b\|_1\right),$$
$$w_b = \frac{\text{TV}(Y_b)}{\sum_{b'=1}^{B} \text{TV}(Y_{b'}) + \varepsilon}.$$
(10)

Here, $Y_b$ is the ground-truth elevations of the $b$-th sample in a batch, and $H$ and $W$ are the height and width. The weight $w_b$ is normalized across the batch, and $\varepsilon$ is a small positive constant. By doing so, we assign higher weights to samples with larger elevation variations.

We use a shallow U-Net [78] model with a gated residual design, as illustrated in Fig. 8. The CNNs output the uncertainty, a raw estimation, and a gating map. The final estimation is obtained by the gated combination of the raw estimation and the input, preserving accuracy in clearly observed areas while selectively overwriting noisy or occluded areas.

We train the models on 54 million frames for each robot, each model takes less than 1 hour to converge. For ANYmal-D, the local grids are of shape $51 \times 31$, with 4-cm resolution, centered at $x = 1.0$ m, $y = 0$ m in the base frame. For TRON1, the local grids are of shape $31 \times 31$, with 4-cm resolution, centered at $x = 0.6$ m, $y = 0$ m in the base frame.

## C. Simulation Integration and Deployment

When deployed in simulation with 1000 parallel ANYmal-D environments, the model's inference time is below $0.3$ ms, and the GPU memory consumption is about 3 GB (for TRON1, these numbers are roughly 60% of those for ANYmal-D). Storing $8, \text{m} \times 8, \text{m}$ global maps for 1000 environments requires about $0.3$ GB, and the global map size can be flexibly chosen since we recenter the map around the robot when it approaches the boundary. Obtaining depth clouds and handling other intermediate overheads takes an additional $\sim 1.2$ GB of GPU memory per 1000 ANYmal-D environments.

On real hardware, the mapping pipeline takes approximately $5$ ms per frame on the onboard CPU, with around $2.5$ ms spent on model inference using ONNX Runtime. This enables fast, low-latency mapping that can keep up with the depth camera's frame rate. For ANYmal-D, we merge the point clouds from two front facing cameras for processing. For TRON1, we use the only front facing camera.

Regarding odometry, we use CompSLAM [82] with Graph-MSF [83] on ANYmal-D to obtain a high-frequency, accurate
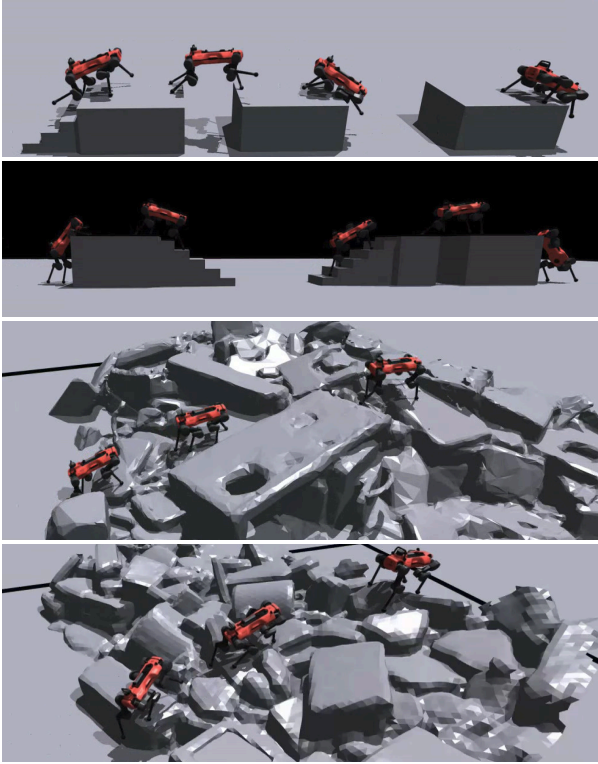
Fig. 9. Our ANYmal-D controller zero-shots the hardest parkour and rubble pile terrains reported in prior work [32], [48].

LiDAR–inertial odometry. On TRON1, we use DLIO [84] instead, as it better handles the platform's higher accelerations. Notably, although these odometry solutions can run reliably at high frequency, their velocity estimates are observed to be noisy and delayed. This is intuitive: for example, a 1 cm position drift over a 20 ms timestep can translate to a velocity error of 0.5 m/s, yet has little impact on map fusion. Therefore, we remove the linear velocity observations from the student policy.

## VI. RESULTS

### A. Agility Comparison with Prior Art

We achieve state-of-the-art agility on both ANYmal-D and TRON1, measured by the difficulty of the terrains they can climb up and down. On ANYmal-D, despite being trained only on primitive terrains, our policy zero-shots the hardest parkour and rubble pile terrains reported in prior work [32], [48], as shown in Fig. 9. On TRON1, our policy climbs up platforms up to $0.48$ m and climbs down up to $0.88$ m, while prior work on bipeds reports $0.5$ m platforms [17] using the Unitree H1 robot (with $4\times$ peak torque and $1.5\times$ base height compared to TRON1).

Across all terrains, our controllers on both robots reach peak forward velocities above $1.5$ m/s, demonstrating a combination of agility and terrain-aware locomotion. On top of this agility, our controllers also exhibit stronger generalization than existing methods, as shown below and benchmarked in Sec. VII.

### B. Real World Results

*1) Quadruped Parkour:* We demonstrate the agility and generalization of our trained ANYmal-D controller on a parkour course, as shown in Fig. 10. This course is not included in the training terrains of either the locomotion controller or the mapping model, yet our system can stably traverse it at speeds of up to 2 m/s, composing climbing and jumping maneuvers.

*2) Sparse Terrain Locomotion:* As shown in Fig. 11, our controllers enable both ANYmal-D and TRON1 to traverse a variety of sparse terrains, many of which are unseen during training, exhibiting generalization capabilities.

*3) Biped maneuvering:* To demonstrate the omnidirectional perceptive locomotion capabilities of our controller, we command TRON1 to traverse a 38-cm-high platform, a gap, a staircase, and rough terrain, as shown in Fig. 12. The robot executes smooth maneuvers in all directions.

*4) Robustness:* We demonstrate the robustness of our controllers on moving terrains. As shown in Fig. 13, TRON1 can traverse or balance on a platform cart with unlocked wheels and spring suspension that can easily tilt to one side with the robot weight. We also show that ANYmal-D can recover from unfixed tilted stepping stones using its knees, after an occasional trip causes a misstep that tilts the blocks.

*5) Mapping Results:* In Fig. 14, we show the mapping results after ANYmal-D traverses terrains that are unseen during training. These maps capture fine-grained terrain details, such as gaps and supports, that help the controller achieve agile and generalized locomotion. Among prior works, only [32] has demonstrated maps of comparable quality during agile motions. However, that approach cannot reliably infer unseen parts of obstacles in unstructured terrains [48], is not efficient enough to support concurrent locomotion training, and does not explicitly encode uncertainty under partial observations such as occlusions. For generalized locomotion, explicitly modeling uncertainty is important: occluded regions are not completed by learned priors but kept uncertain, and newly observed geometry can be integrated into the map based on the predicted uncertainty once it becomes visible.

### C. Emergent Behaviors

*1) Active Perception:* Trained with partial observations and uncertainty-aware mapping, our locomotion controllers exhibit emergent active perception behaviors. In Fig. 15, we exemplify this with ANYmal-D attempting to climb onto a high obstacle. In one attempt, the robot makes contact with the obstacle and cannot immediately climb up, due to its limited field of view and imperfect depth sensing. This interaction, however, exposes higher parts of the obstacle to the map. In the subsequent attempt, the controller uses this updated map and successfully climbs onto the obstacle. The maps are visualized to show how information gained from previous interactions is reused to make subsequent movements more informed and successful.

Notably, this addresses a stated limitation of prior state-of-the-art work [48], where a recurrent controller fails to climb onto a high box when starting close to it due to a lack of longer-term memory. In contrast, our mapping module explicitly maintains long-term spatial information, providing the
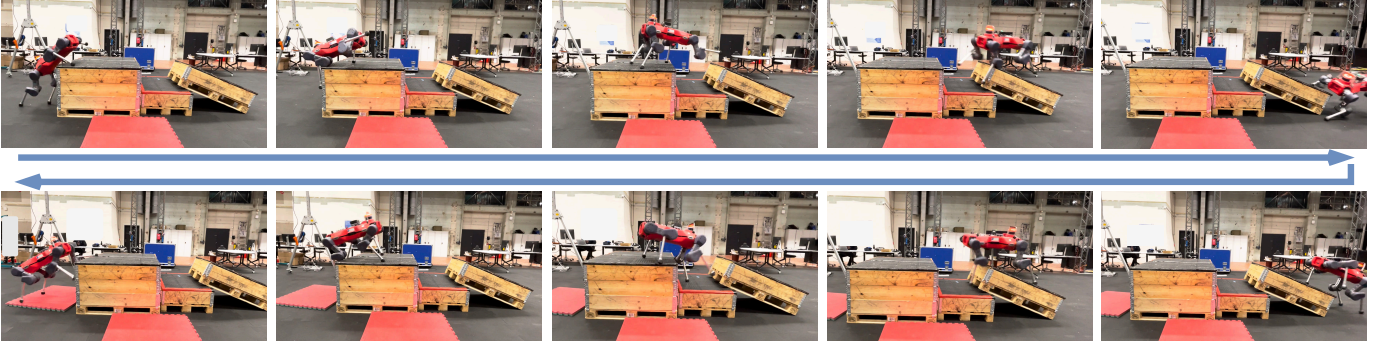
Fig. 10. Our system enables ANYmal-D to move back and forth on a parkour course, a terrain that is unseen during training of both the controller and the mapping model.
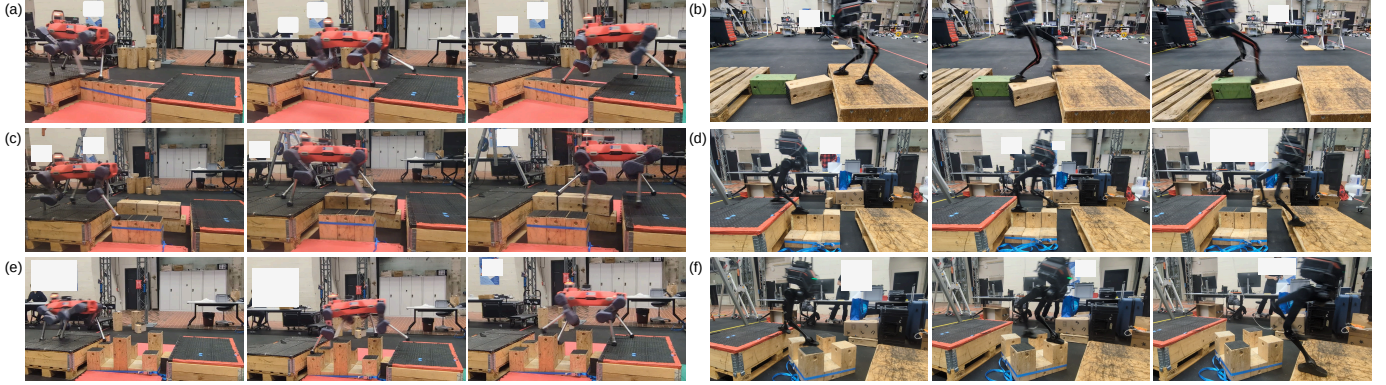


Fig. 11. Our systems enable ANYmal-D and TRON1 to traverse diverse sparse terrains. (a) ANYmal-D traverses a 19-cm wide balance beam, which is in the training terrains. (b) TRON1 traverses two unfixed floating 19-cm wide blocks, which make an unseen curved beam terrain. (c) ANYmal-D traverses two unfixed beams with gaps, which is unseen during training. (d) TRON1 traverses a beam followed by a gap, which is an unseen combination during training. (e) ANYmal-D traverses two rows of 19-cm wide stepping stones with 10-cm height differences, which is unseen during training. (f) TRON1 traverses diamond-layout stepping stones, which is unseen during training.



Fig. 12. TRON1 maneuvers over a 38-cm-high platform, a gap, a staircase, and rough terrain, showcasing its omnidirectional perceptive locomotion capabilities. Orange arrows indicate the robot's approximate trajectories.

controller with an effective representation of the environment over time. Moreover, consistent with [48], we also observe the robot looking upward while climbing onto platforms, gaining additional terrain information.

*2) Loco-Navigation:* Since our locomotion controllers are trained to reach goals, they also exhibit local navigation capabilities. As shown in Fig. 16, the controllers can align the robot with the terrain before traversal and avoid collisions with obstacles. These behaviors can simplify remote human

operation by reducing the need for high-frequency, navigation-level commands (such as velocity commands [13]).

*3) Whole-Body Contact:* As previously shown in Fig. 10 and Fig. 13, our ANYmal-D controller actively use knee contacts to help stabilize the robot and traverse challenging terrains. On one hand, such whole-body contact behaviors are beneficial for agility and robustness; on the other hand, because most existing legged robots are designed primarily for foot contacts, these behaviors may stress the hardware. That said, we believe that whole-body contact will become increasingly important for achieving whole-body dexterity in complex environments [64], [85], [86].

*4) Impact Reduction:* As mentioned in Sec. IV-D2, we use human and dog motion data to set thigh-acceleration thresholds for early termination. This induces impact-reduction behaviors, as shown in Fig. 17: when climbing down, ANYmal-D uses its knees to gently touch down, while TRON1 retracts the support leg to absorb the impact.

### D. Interpretable Feature Patterns

We visualize the feature patterns of our ANYmal-D student policy's AME-2 encoder in Fig. 18. For local features, we visualize their attention weights with a red–blue colormap. For global features, which are obtained via max pooling, we treat pooled points as selected (weight 1) or not (weight 0) and average this binary mask over the feature dimension.
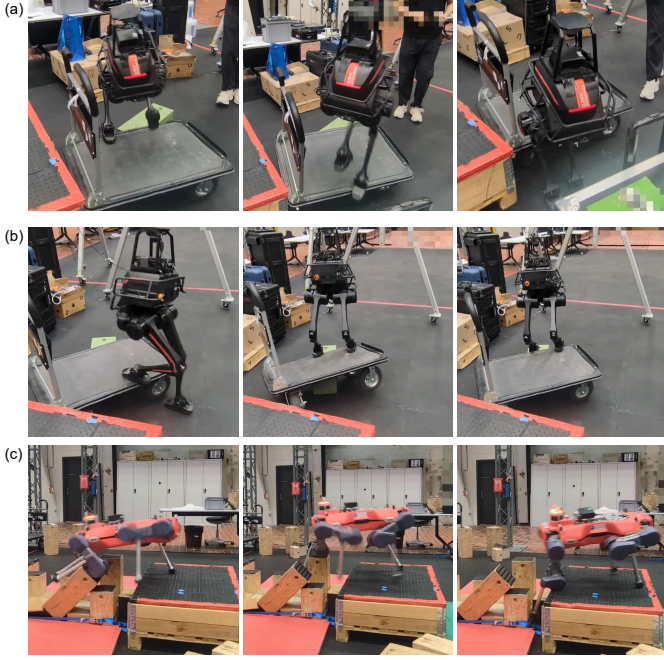
Fig. 13. Our controllers are robust to moving terrains. (a) TRON1 climbs over an unlocked tiltable platform cart. (b) TRON1 climbs onto and balances on an unlocked tiltable platform cart. (c) ANYmal-D recovers from unfixed tilted stepping stones using its knees.
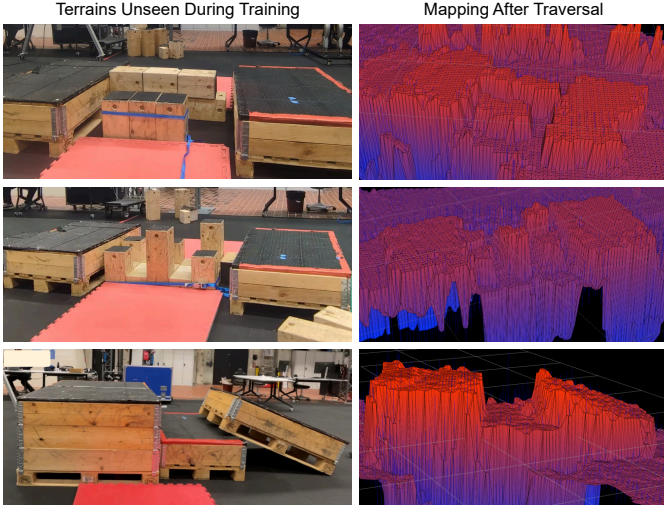


Fig. 14. Obtained maps after ANYmal-D's traversal of terrains that are unseen during training. Our mapping supports agile and generalized locomotion. In the visualizations, meshes represent the estimated elevations, and thin lines indicate the uncertainties. The meshes are colored with a red–blue colormap to enhance height contrast, excluding regions that are never covered by the cameras.

By inspecting these feature maps, we observe that local attention weights often emphasize fine-grained terrain details, while global features tend to concentrate on a sparse set of distinctive points across terrain types (e.g., obstacle boundaries, high platform surfaces, and beam centers). We also find that these patterns persist on test terrains, which helps explain our controller's generalization to unseen terrains despite being trained only on primitive terrains.
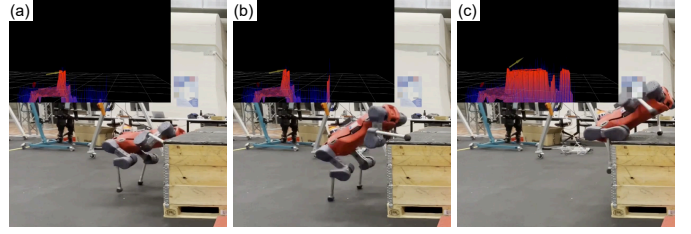


Fig. 15. Emergent active perception behavior. (a) In a climb-up attempt, ANYmal-D collides with the obstacle due to limited field of view and depth sensing quality. (b) Although unsuccessful, this trial reveals higher parts of the obstacle in the map, enabling the robot to succeed in the immediate retrial. (c) ANYmal-D then climbs onto the obstacle. The maps are visualized in the top-left corner of each subfigure, excluding regions that are never covered by the cameras.
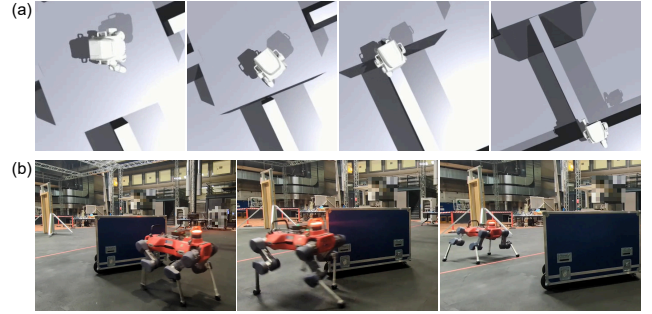


Fig. 16. Emergent local navigation capabilities. (a) TRON1 turns around and aligns itself with the beam to traverse the terrain. (b) ANYmal-D avoids an obstacle when reaching the goal.

## VII. ABLATIONS AND BENCHMARKING

In this section, we benchmark our design choices on test terrains (depicted in Fig. 4) that are unseen during training. All benchmarking experiments are conducted on ANYmal-D in simulation, where richer prior works and baselines are available.

### A. Teacher Architecture

We compare our teacher policy architecture against two baselines: AME-1 [15] and MoE [68]. AME-1 is the prior state of the art in generalized locomotion, outperforming policies based on MLPs, CNNs, multimodal transformers [37], and Vision Transformers [87]. The MoE architecture from [68] has demonstrated strong scalability by enabling a humanoid to track diverse human motions.

The results are summarized in Table II. All architecture designs scale well on the training terrains, highlighting the stability of our training framework. However, MoE shows very limited generalization to unseen test terrains. AME-1 generalizes well on sparse terrains (test 1) but struggles when different terrains are mixed (test 2-4), which we attribute to its encoder assigning attention to local features based only on proprioception, thus missing global context. In contrast, our AME-2 teacher generalizes reliably across all test terrains.

### B. Student Designs

Having a generalized teacher does not automatically yield a generalized deployable student. We therefore compare our AME-2 student against the following baselines:

TABLE II
SUCCESS RATES (%) OF DIFFERENT TEACHER ARCHITECTURES ON TRAINING AND TEST TERRAINS

| Teacher Architecture | Training Terrains | | | Test Terrains | | | | |
|---|---|---|---|---|---|---|---|---|
| | Dense | Climbing | Sparse | Test 1 | Test 2 | Test 3 | Test 4 | avg. Test |
| **AME-2 (Teacher)** | **97.9** | **96.0** | **95.2** | **96.8** | **93.8** | **99.7** | **90.6** | **95.2** |
| AME-1 [15] | **96.0** | **95.7** | **95.5** | **99.2** | 29.2 | 32.2 | 44.1 | 51.2 |
| MoE [68] | **95.0** | **96.7** | **90.5** | 52.9 | 55.3 | 41.0 | 30.8 | 45.0 |

Best results (within 5% of the maximum) are shown in **bold**.

TABLE III
SUCCESS RATES (%) OF DIFFERENT STUDENT DESIGNS ON TRAINING AND TEST TERRAINS

| Student Design | Training Terrains | | | Test Terrains | | | | |
|---|---|---|---|---|---|---|---|---|
| | Dense | Climbing | Sparse | Test 1 | Test 2 | Test 3 | Test 4 | avg. Test |
| **AME-2 (Student)** | **96.4** | **96.4** | **91.8** | **90.0** | 77.7 | 89.1 | **72.9** | **82.4** |
| Visual Recurrent Student [48] | **98.0** | **97.9** | 86.4 | 34.5 | 20.9 | **99.8** | 50.6 | 51.5 |
| AME-2 (Student) w/o RL | **95.8** | 87.9 | 84.6 | 46.9 | 59.1 | 90.6 | 46.0 | 60.7 |
| AME-2 (Student) w/o rep. loss | **95.8** | **94.8** | 86.5 | 81.9 | **82.6** | 69.5 | 60.4 | 73.6 |
| AME-2 (Teacher) | 97.9 | 96.0 | 95.2 | 96.8 | 93.8 | 99.7 | 90.6 | 95.2 |

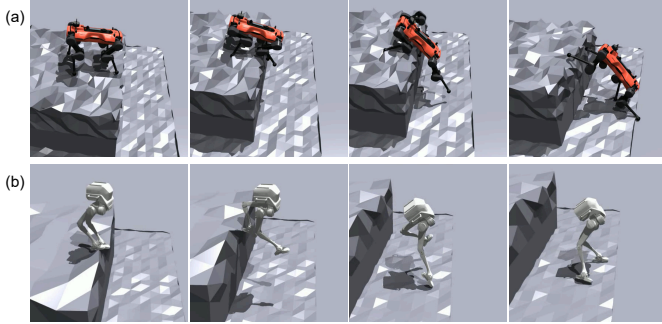Best student results (within 5% of the maximum) are shown in **bold**.

Fig. 17. Emergent impact reduction behaviors. (a) ANYmal-D uses its knee to support the body when climbing down, and gently touches the ground to reduce the impact. (b) TRON1 extends one leg while jumping down, and retracts it during landing before switching support to the other leg.

1) The end-to-end student design from [48], the previous state of the art in agile generalist locomotion, which uses recurrent networks with vision inputs.
2) An AME-2 student trained without RL losses, using only teacher action distillation and the representation loss.
3) An AME-2 student trained without the representation loss.

The results are summarized in Table III. On the training terrains, all student policies perform well, with our proposed design achieving slightly higher success rates on sparse terrains. On unseen test terrains, our student achieves the best overall success rates.

The end-to-end visual recurrent student from [48] performs slightly better than ours on *Test 3*, a parkour course with obstacles, but underperforms on all other test terrains that contain sparse regions. On *Test 3*, most of our student's failures occur when it attempts to climb the obstacles, which can be misinterpreted as higher platforms due to occlusions and the observed surface elevations. In contrast, using the full depth image, as in the visual recurrent design, may provide richer cues and a longer reaction horizon for obstacle avoidance. Despite this modest drop in generalization around obstacles, our student controller achieves much higher success rates on sparse and mixed terrains.

Regarding the learning setup, we find that direct teacher supervision (action distillation plus representation alignment) already yields good performance on the training terrains, but does not produce strong generalization to unseen terrains without RL. Incorporating RL is important for performance on the test terrains, and the representation loss gives a further improvement. Overall, our student policy architecture performs better on unseen terrains than the end-to-end visual recurrent policy, and our learning design, with mixed RL losses and teacher supervision losses (action distillation and representation alignment), further enhances the generalization.

### C. Mapping Designs

We benchmark our mapping pipeline against two baselines:

1) Our pipeline, with the mapping model trained only on locomotion training terrains.
2) A temporal recurrent model that directly predicts egocentric elevation maps from local grid observations and delta transforms.

The first baseline tests whether additional diverse meshes (both traversable and untraversable) are necessary for better generalization. The second contrasts our per-frame fusion approach with an end-to-end recurrent neural memory model: the recurrent model tends to overfit terrain patterns seen during training and produces lower-quality maps on unseen terrains, as we will show in the qualitative comparison in Fig. 19.

All pipelines are evaluated on data collected from the same rollout trajectories, with results reported in Table IV. We use the $L_{0.5}$ loss defined in Eq. 9 as the metric, which balances estimation accuracy and uncertainty quantification. The temporal recurrent model (detailed in Appendix D) uses a CNN encoder for local grids, an MLP encoder for delta transforms, and Spatially-Enhanced Recurrent Units (SRU), an architecture that enables 100 m-scale mapless visual navigation [88], as the memory module. To stabilize training and prevent the model from collapsing to large uncertainties and poor estimates, we have to train the recurrent model only on locomotion training terrains with the policy in the loop (as done in prior work [29],
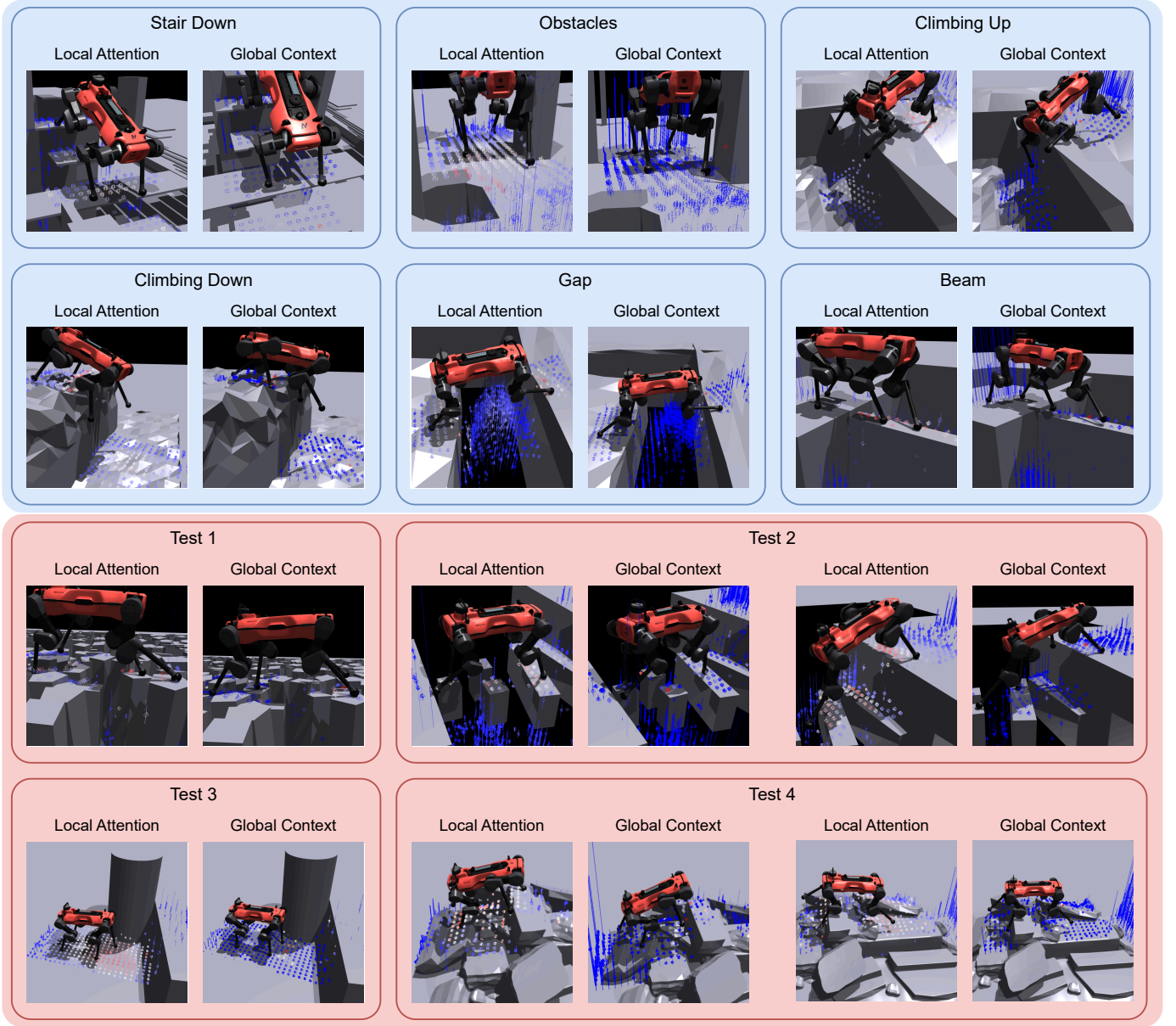
Fig. 18. Feature patterns of our ANYmal-D controller. We visualize the neural elevation map together with the corresponding local-attention and global-context weights. Points indicate elevation estimates, lines indicate uncertainties, and higher-intensity red colors correspond to higher weights. Local attention weights provide fine-grained patterns for locomotion control, while global context features primarily focus on sparse points that characterize the terrain type.

[33]) and add an auxiliary L2 loss on the elevation estimates in addition to $L_{0.5}$. This setup requires roughly $5\times$ more training time to converge than our mapping models.

It is worth noting that the recurrent model can optimize both the elevation estimates and the associated uncertainties in unseen regions based on the training distribution, which is advantageous under the $L_{0.5}$ loss, whereas our pipeline uses heuristic defaults. However, as illustrated in Fig. 19, our mapping pipeline produces higher-quality maps in practice, and can effectively model occlusions with high uncertainties.

### D. Robustness to Visual Noise

We validate the robustness of our system to noise and sensor degradation. We evaluate our student policy's success rates under missing points (removed from depth clouds), artifacts (replaced by random points in depth clouds), and with the front upper camera disabled. The tested setups are all above the randomization ranges during training.

The results are presented in Table V. Under missing points and artifacts, we find the policy performance does not drop, and can sometimes be slightly better. This can be explained by that, fewer effective points in the depth clouds result in more uncertain regions and make the policy behave more conservatively.

When the front upper camera is disabled, the system still performs well on most terrains without re-training, highlighting the flexibility of our method for different sensor configurations. However, on terrains that require active perception, such as the high obstacles to climb in *Test 2* and *Test 3*, the policy struggles, suggesting that these behaviors remain sensitive to the available sensor viewpoints.

TABLE IV
$L_{0.5}$ Loss of Different Neural Mapping Methods on Training and Test Terrains

| Mapping Pipeline | Training Terrains | | | Test Terrains | | | | |
|---|---|---|---|---|---|---|---|---|
| | Dense | Climbing | Sparse | Test 1 | Test 2 | Test 3 | Test 4 | avg. Test |
| **Ours** | -0.006 | -0.108 | 0.020 | **0.033** | 0.227 | **0.036** | -0.111 | **0.046** |
| Ours (Loco-Only) | 0.020 | -0.090 | 0.081 | 0.104 | 0.234 | 0.087 | -0.075 | 0.088 |
| Temporal Recurrent | **-0.162** | **-0.114** | **-0.101** | 0.316 | **0.135** | 0.066 | **-0.176** | 0.085 |

Best (lowest) results are shown in **bold**.

TABLE V
Success Rates (%) of the Student Policy under Noise and Sensor Degradation

| Setup | Training Terrains | | | Test Terrains | | | | |
|---|---|---|---|---|---|---|---|---|
| | Dense | Climbing | Sparse | Test 1 | Test 2 | Test 3 | Test 4 | avg. Test |
| Nominal | 96.4 | 96.4 | 91.8 | 90.0 | 77.7 | 89.1 | 72.9 | 82.4 |
| 20% Missing Points | 95.5 | 94.9 | 91.3 | 91.1 | 77.3 | 90.3 | 74.8 | 83.4 |
| 3% Artifacts | 96.4 | 97.6 | 93.2 | 95.2 | 72.8 | 92.5 | 80.3 | 85.2 |
| 1 Camera Disabled | 91.2 | 82.0 | 85.9 | 95.1 | 56.0 | 1.1 | 74.6 | 56.7 |

## VIII. Discussion

In this paper, we propose AME-2, a unified RL framework for agile and generalizable legged locomotion with attention-based neural map encoding. We use goal-reaching rewards to incentivize agile behaviors, and design the AME-2 encoder to achieve generalization across diverse challenging terrains. To facilitate sim-to-real deployment, we combine teacher-student learning with RL and introduce a lightweight neural mapping pipeline that explicitly models partial observability while running identically in simulation and on real robots. Our framework applies to both a quadruped and a biped without changing the rewards or training settings.

### A. Modular v.s. End-to-End

Classical locomotion systems [2], [4] use modular designs (perception, planning, and control) to make the problem tractable, and can generalize well when their underlying assumptions hold (e.g., low uncertainty, stable foot contacts, accurate mapping and dynamics modeling). Recent works instead reduce system complexity through end-to-end sensor-to-motor policy learning [38], [48], but have not yet demonstrated strong generalization to unseen terrains.

Our work tries to bridge these two. We use a mapping module and a controller module, and train the neural controller end-to-end to generate joint actions from observations. In the controller, our AME-2 encoder also learns planning-level representations (global context and local contact-relevant features) for the subsequent decoder, playing a similar role to the planning module in classical solutions. This design preserves a modular structure at the system level while avoiding explicit online model-based planning. In the mapping module, we keep a simple heuristic structure, but also run a lightweight neural network in the loop to improve robustness and generalization under noise, partial observability, and unseen terrain patterns.

### B. Goal Reaching v.s. Velocity Tracking

We use a goal-reaching formulation [20] instead of velocity tracking [15]. This encourages more agile behaviors and reduces reliance on high-frequency navigation commands. A drawback is that it offers less direct control over the robot's intermediate motion, which can be troublesome when semantic obstacles lie along the path but appear traversable from the geometry. We envision combining both command interfaces within a single controller in future work, e.g., via masked commands and distillation [89].

### C. Whole-Body Skills

We achieve whole-body locomotion skills in this work, but for a given terrain type the learned motions tend to follow similar contact patterns. For systems with higher DoFs, such as humanoids, the same terrain type at different scales may require different contact strategies. For example, a humanoid may step onto low obstacles with a single leg, use a two-leg jump for medium heights, and leverage both arms and legs for higher obstacles. It remains unclear whether our method, without additional references or priors such as [86], can automatically discover such diverse contact patterns to handle a broader range of terrains in higher-DoF systems.

### D. Towards Higher Success Rates on Unseen Terrains

We still observe more failures on unseen test terrains than on training terrains, as expected for learning-based systems. Notably, most failures occur during skill transitions. For example, on *Test 2* the robot must first decelerate on a sparse region and then execute a climbing maneuver while maintaining precise hind-leg footholds. While these specific cases could likely be addressed by finetuning [48], it remains unclear whether there is a scalable, principled approach for learning such challenging skill transitions with zero-shot generalization.

### E. Limitations and Future Work

This work has several limitations. First, we use a 2.5-D elevation map and do not address fully 3-D locomotion. Future work could incorporate multi-layer elevation maps [24], which are directly compatible with our policy architecture, or explore attention-based voxel representations. Second, our controllers are not designed for severely degraded perception, such as high grass or snow. A promising direction is to combine robust
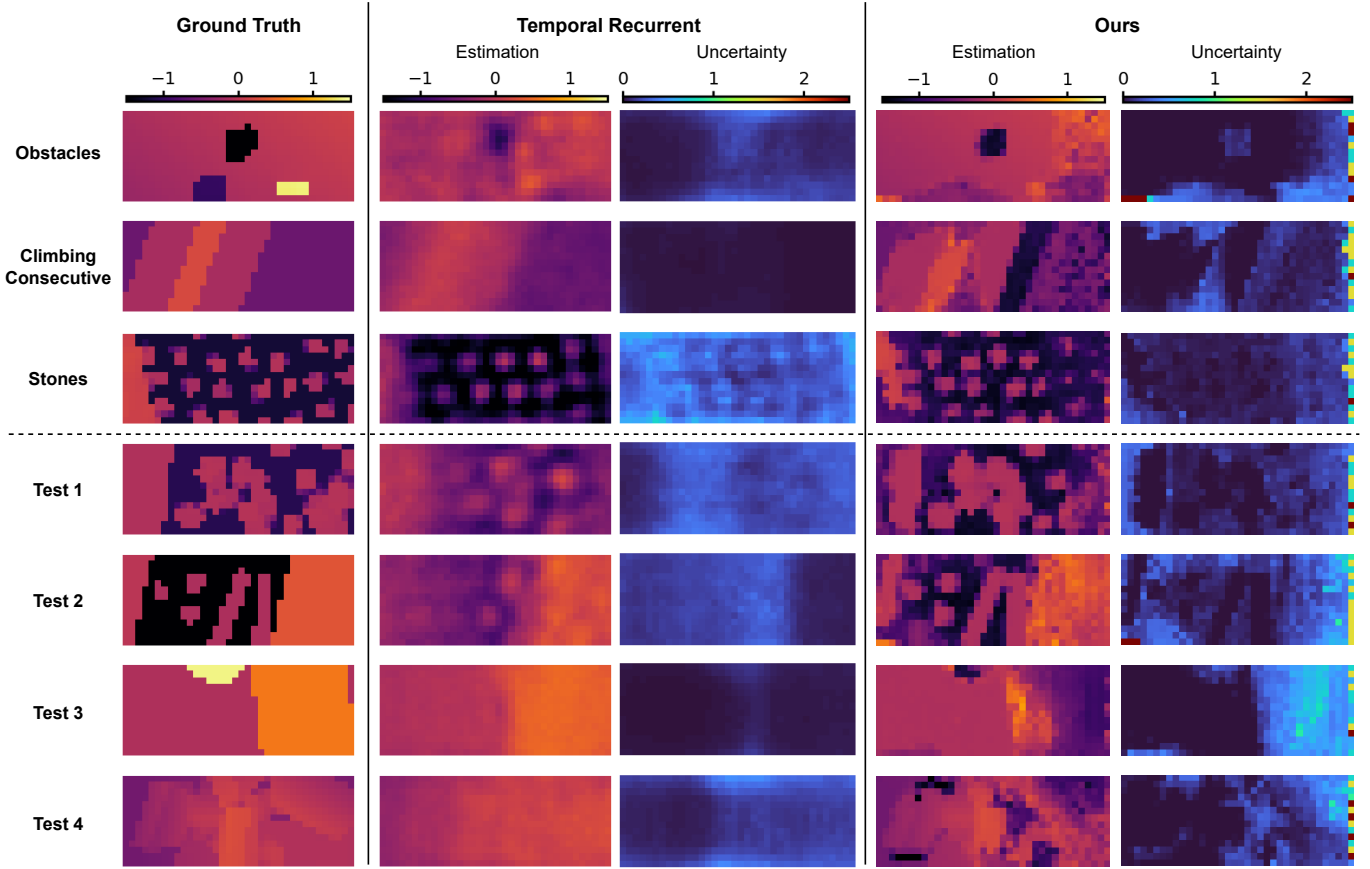
Fig. 19. Qualitative comparison between the temporal recurrent model and our mapping pipeline on training and test terrains. The visualized maps are egocentric, with the robot facing rightward. Our mapping pipeline produces accurate estimations when confident and assigns high uncertainty to occluded regions. In contrast, the temporal recurrent model generates less meaningful uncertainty maps, loses fine-grained details, and performs worse on unseen terrains.

controllers [14] with our agile and generalized controllers into a single scene-aware policy. Third, our mapping module can fail in highly dynamic environments with occlusions, and could be extended to explicitly reason about moving elements in the scene.

## APPENDIX A
## TERRAIN PARAMETERS

We use the following terrain proportions and parameters for training:

- **Rough (Dense, 5%)**: Heightfields generated from uniform noise, with terrain resolution smaller than the map resolution. Over the curriculum, the noise range increases from $\pm 0$ m to $\pm 0.2$ m for ANYmal-D, and from $\pm 0$ m to $\pm 0.15$ m for TRON1. Goals are sampled anywhere on the terrain.

- **Stair Down (Dense, 5%)**: Floating stairs descending with random widths and random walls on both sides. The slope is increased from $5°$ to $45°$ over the curriculum for both robots. Goals are sampled anywhere on the terrain.

- **Stair Up (Dense, 5%)**: Floating stairs ascending with random widths and random walls on both sides. The slope is increased from $5°$ to $45°$ over the curriculum for both robots. Goals are sampled anywhere on the terrain.

- **Boxes (Dense, 5%)**: Heightfields generated by adding and removing random boxes on flat ground. Over the curriculum, the maximum box height increases from $0.05$ m to $0.4$ m for ANYmal-D, and from $0.05$ m to $0.3$ m for TRON1. Goals are sampled anywhere on the terrain.

- **Obstacles (Dense, 5%)**: Random positive and negative obstacles placed on random slopes. Over the curriculum, the obstacle density increases from $0$ m$^{-2}$ to $0.5$ m$^{-2}$ for both robots. Goals are sampled anywhere on the terrain.

- **Climbing Up (Climbing, 20%)**: A pit to climb out of, with vertices near the edges randomly perturbed by up to 20% of the height. Over the curriculum, the pit height increases from $0.1$ m to $1.0$ m for ANYmal-D, and from $0.1$ m to $0.48$ m for TRON1. Goals are sampled at the

other end of the terrain.

- **Climbing Down (Climbing, 5%)**: A platform to climb down from, with vertices near the edges randomly perturbed by up to 20% of the height. Over the curriculum, the platform height increases from 0.2 m to 1.0 m for ANYmal-D, and from 0.2 m to 0.88 m for TRON1. Goals are sampled at the other end of the terrain.
- **Climbing Consecutive (Climbing, 5%)**: Two stacked layers of rings to climb out of. The robot needs to consecutively climb up twice and then climb down twice. Over the curriculum, the first ring height increases from 0.05 m to 0.5 m for ANYmal-D and from 0.05 m to 0.3 m for TRON1. The second ring height increases from 0.05 m to 0.4 m for ANYmal-D and from 0.05 m to 0.3 m for TRON1. Goals are sampled at the other end of the terrain.
- **Gap (Sparse, 5%)**: A gap to jump over, with vertices near the edges randomly perturbed by up to 10% of the distance, and a height difference between the two sides randomly sampled within $\pm30\%$ of the distance. Over the curriculum, the distance increases from 0.1 m to 1.1 m for ANYmal-D and from 0.1 m to 0.6 m for TRON1. Goals are sampled at the other end of the terrain.
- **Pallets (Sparse, 5%)**: Parallel beams with random orientations. Over the curriculum, the beam width decreases from 0.4 m to 0.16 m for both robots. The gap width increases from 0.08 m to 0.35 m for ANYmal-D and from 0.08 m to 0.2 m for TRON1. The inter-beam height difference increases from 0 m to 0.3 m for ANYmal-D and from 0 m to 0.2 m for TRON1. Goals are sampled at the other end of the terrain.
- **Stones (Sparse, 30%)**: We follow the "Stones-Everywhere" design in [26] and [15]. Goals are sampled anywhere on the terrain.
- **Beam (Sparse, 5%)**: A beam connecting two platforms, randomly inclined with roll, pitch, and yaw sampled within $\pm0.1$ rad, and with a platform height difference sampled within $\pm0.2$ m. Over the curriculum, the beam width decreases from 0.9 m to 0.18 m for both robots. Goals are sampled at the other end of the terrain.

For sparse terrains, we add physical floors to half of them and virtual floors (visible in the map without physical collision) to the other half. The floor height is randomized between $-1.5$ m and $-0.35$ m. This discourages the robot from walking on the floor when traversing shallow sparse terrains.

## APPENDIX B
### RANDOMIZATION PARAMETERS

For both robots, we randomize the payload within $[-5, 5]$ kg, actuation delays within $[0, 0.02]$ s, and friction coefficients within $[0.3, 1.0]$. For TRON1, we additionally apply random biases to PD gains and armatures within $\pm15\%$ of the raw values.

For policy observations, we add zero-mean uniform noise with the following maximum magnitudes: 0.1 m/s for base linear velocity, 0.2 rad/s for base angular velocity, 0.05 for

projected gravity, 0.01 rad for joint positions, 1.5 rad/s for joint velocities, and 0.05 m for map observations. We also add a random drift within $[-0.03, 0.03]$ m when querying map observations.

During student policy training, we make 15% of depth cloud points missing and 2% artifacts. We give 10% of the environments access to complete maps with variance 0.0025 m$^2$. We corrupt 1% points in the student map observations by assigning random values with a random variance larger than 1 m$^2$.

## APPENDIX C
### PPO PARAMETERS

We use the following PPO parameters in Table VI.

TABLE VI
PPO PARAMETERS

| Parameter | Value |
|---|---|
| Value Loss Coefficient | 1.0 |
| Value Loss Clip | 0.2 |
| Entropy Coefficient | Decay from 0.004 to 0.001 |
| Learning Epochs per Iteration | 4 |
| Mini Batches per Iteration | 3 |
| Simulation Steps per Iteration | 24 |
| Number of Environments | 4800 |
| Learning Rate | Adaptive |
| $\gamma$ | 0.99 |
| $\lambda$ | 0.95 |
| Desired KL Divergence | 0.01 |
| *Student only* | |
| Learning Rate When Surrogate Loss Disabled | 0.001 |
| Action Distillation Loss Coefficient | 0.02 |
| Representation Loss Coefficient | 0.2 |

## APPENDIX D
### TEMPORAL RECURRENT BASELINE MAPPING MODEL

The architecture of the recurrent model used as the mapping baseline is shown in Fig. 20. We provide the model with ground-truth delta transforms, which are sufficient for map reconstruction and serve as a clean replacement for the proprioceptive observations used in prior work [29], [33], where transforms are not directly available.
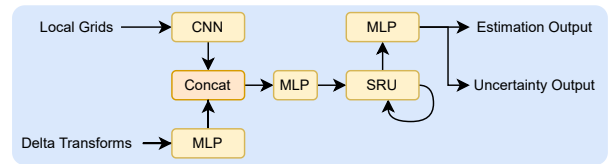


Fig. 20. The temporal recurrent model architecture to predict egocentric elevation mapping.

## REFERENCES

[1] P. Fankhauser, M. Bjelonic, C. D. Bellicoso, T. Miki, and M. Hutter, "Robust rough-terrain locomotion with a quadrupedal robot," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 5761–5768.

[2] F. Jenelten, R. Grandia, F. Farshidian, and M. Hutter, "Tamols: Terrain-aware motion optimization for legged systems," *IEEE Transactions on Robotics*, vol. 38, no. 6, pp. 3395–3413, 2022.

[3] S. Fahmi, V. Barasuol, D. Esteban, O. Villarreal, and C. Semini, "Vital: Vision-based terrain-aware locomotion for legged robots," *IEEE Transactions on Robotics*, vol. 39, no. 2, pp. 885–904, 2022.

[4] R. Grandia, F. Jenelten, S. Yang, F. Farshidian, and M. Hutter, "Perceptive locomotion through nonlinear model-predictive control," *IEEE Transactions on Robotics*, vol. 39, no. 5, pp. 3402–3421, 2023.

[5] R. Akizhanov, V. Dhédin, M. Khadiv, and I. Laptev, "Learning feasible transitions for efficient contact planning," *arXiv preprint arXiv:2407.11788*, 2024.

[6] R. J. Griffin, G. Wiedebach, S. McCrory, S. Bertrand, I. Lee, and J. Pratt, "Footstep planning for autonomous walking over rough terrain," in *2019 IEEE-RAS 19th International Conference on Humanoid Robots (Humanoids)*, 2019, pp. 9–16.

[7] C. Mastalli, W. Merkt, G. Xin, J. Shim, M. Mistry, I. Havoutis, and S. Vijayakumar, "Agile maneuvers in legged robots: a predictive control approach," *arXiv preprint arXiv:2203.07554*, 2022.

[8] F. Jenelten, T. Miki, A. E. Vijayan, M. Bjelonic, and M. Hutter, "Perceptive locomotion in rough terrain – online foothold optimization," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5370–5376, 2020.

[9] F. Jenelten, J. Hwangbo, F. Tresoldi, C. D. Bellicoso, and M. Hutter, "Dynamic locomotion on slippery ground," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 4170–4176, 2019.

[10] P. Fankhauser, M. Bloesch, and M. Hutter, "Probabilistic terrain mapping for mobile robots with uncertain localization," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3019–3026, 2018.

[11] T. Miki, L. Wellhausen, R. Grandia, F. Jenelten, T. Homberger, and M. Hutter, "Elevation mapping for locomotion and navigation using gpu," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 2273–2280.

[12] S. Ha, J. Lee, M. van de Panne, Z. Xie, W. Yu, and M. Khadiv, "Learning-based legged locomotion: State of the art and future perspectives," *The International Journal of Robotics Research*, vol. 44, no. 8, pp. 1396–1427, 2025.

[13] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Conference on robot learning*. PMLR, 2022, pp. 91–100.

[14] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science robotics*, vol. 7, no. 62, p. eabk2822, 2022.

[15] J. He, C. Zhang, F. Jenelten, R. Grandia, M. Bächer, and M. Hutter, "Attention-based map encoding for learning generalized legged locomotion," *Science Robotics*, vol. 10, no. 105, p. eadv3604, 2025.

[16] A. Allshire, H. Choi, J. Zhang, D. McAllister, A. Zhang, C. M. Kim, T. Darrell, P. Abbeel, J. Malik, and A. Kanazawa, "Visual imitation enables contextual humanoid control," *arXiv preprint arXiv:2505.03729*, 2025.

[17] J. Long, J. Ren, M. Shi, Z. Wang, T. Huang, P. Luo, and J. Pang, "Learning humanoid locomotion with perceptive internal model," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 9997–10 003.

[18] J. Ren, T. Huang, H. Wang, Z. Wang, Q. Ben, J. Long, Y. Yang, J. Pang, and P. Luo, "Vb-com: Learning vision-blind composite humanoid locomotion against deficient perception," *arXiv preprint arXiv:2502.14814*, 2025.

[19] H. Wang, Z. Wang, J. Ren, Q. Ben, T. Huang, W. Zhang, and J. Pang, "Beamdojo: Learning agile humanoid locomotion on sparse footholds," in *Robotics: Science and Systems (RSS)*, 2025.

[20] N. Rudin, D. Hoeller, M. Bjelonic, and M. Hutter, "Advanced skills by learning locomotion and local navigation end-to-end," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 2497–2503.

[21] F. Jenelten, J. He, F. Farshidian, and M. Hutter, "Dtc: Deep tracking control," *Science Robotics*, vol. 9, no. 86, p. eadh5401, 2024.

[22] S. Gangapurwala, M. Geisert, R. Orsolino, M. Fallon, and I. Havoutis, "Rloc: Terrain-aware legged locomotion using reinforcement learning and optimal control," *IEEE Transactions on Robotics*, vol. 38, no. 5, pp. 2908–2927, 2022.

[23] Y. Dong, J. Ma, L. Zhao, W. Li, and P. Lu, "Marg: Mastering risky gap terrains for legged robots with elevation mapping," *IEEE Transactions on Robotics*, 2025.

[24] Y. Chen, J. Ma, Z. Luo, Y. Han, Y. Dong, B. Xu, and P. Lu, "Learning autonomous and safe quadruped traversal of complex terrains using multi-layer elevation maps," *IEEE Robotics and Automation Letters*, 2025.

[25] T. Miki, J. Lee, L. Wellhausen, and M. Hutter, "Learning to walk in confined spaces using 3d representation," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 8649–8656.

[26] C. Zhang, N. Rudin, D. Hoeller, and M. Hutter, "Learning agile locomotion on risky terrains," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 11 864–11 871.

[27] F. Shi, C. Zhang, T. Miki, J. Lee, M. Hutter, and S. Coros, "Rethinking robustness assessment: Adversarial attacks on learning-based quadrupedal locomotion controllers," in *Robotics: Science and Systems XX, Delft, The Netherlands, July 15-19, 2024*, D. Kulic, G. Venture, K. E. Bekris, and E. Coronado, Eds., 2024.

[28] J. Lee, M. Bjelonic, A. Reske, L. Wellhausen, T. Miki, and M. Hutter, "Learning robust autonomous navigation and locomotion for wheeled-legged robots," *Science Robotics*, vol. 9, no. 89, p. eadi9641, 2024.

[29] J. Sun, G. Han, P. Sun, W. Zhao, J. Cao, J. Wang, Y. Guo, and Q. Zhang, "Dpl: Depth-only perceptive humanoid locomotion via realistic depth synthesis and cross-attention terrain reconstruction," *arXiv preprint arXiv:2510.07152*, 2025.

[30] Y. Yang, G. Shi, C. Lin, X. Meng, R. Scalise, M. G. Castro, W. Yu, T. Zhang, D. Zhao, J. Tan *et al.*, "Agile continuous jumping in discontinuous terrains," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 10 245–10 252.

[31] D. Hoeller, N. Rudin, C. Choy, A. Anandkumar, and M. Hutter, "Neural scene representation for locomotion on structured terrain," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 8667–8674, 2022.

[32] D. Hoeller, N. Rudin, D. Sako, and M. Hutter, "Anymal parkour: Learning agile navigation for quadrupedal robots," *Science Robotics*, vol. 9, no. 88, p. eadi7566, 2024.

[33] R. Yu, Q. Wang, Y. Wang, Z. Wang, J. Wu, and Q. Zhu, "Walking with terrain reconstruction: Learning to traverse risky sparse footholds," *arXiv preprint arXiv:2409.15692*, 2024.

[34] H. Duan, B. Pandit, M. S. Gadde, B. Van Marum, J. Dao, C. Kim, and A. Fern, "Learning vision-based bipedal locomotion for challenging terrain," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 56–62.

[35] A. Agarwal, A. Kumar, J. Malik, and D. Pathak, "Legged locomotion in challenging terrains using egocentric vision," in *Proceedings of The 6th Conference on Robot Learning*. PMLR, 2023.

[36] X. Cheng, K. Shi, A. Agarwal, and D. Pathak, "Extreme parkour with legged robots," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 11 443–11 450.

[37] R. Yang, M. Zhang, N. Hansen, H. Xu, and X. Wang, "Learning vision-guided quadrupedal locomotion end-to-end with cross-modal transformers," in *International Conference on Learning Representations*, 2022.

[38] Z. Zhuang, S. Yao, and H. Zhao, "Humanoid parkour learning," in *Conference on Robot Learning*. PMLR, 2025, pp. 1975–1991.

[39] S. Li, S. Luo, J. Wu, and Q. Zhu, "Move: Multi-skill omnidirectional legged locomotion with limited view in 3d environments," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025, pp. 7647–7653.

[40] Q. Ben, B. Xu, K. Li, F. Jia, W. Zhang, J. Wang, J. Wang, D. Lin, and J. Pang, "Gallant: Voxel grid-based humanoid locomotion and local-navigation across 3d constrained terrains," *arXiv preprint arXiv:2511.14625*, 2025.

[41] S. Luo, S. Li, R. Yu, Z. Wang, J. Wu, and Q. Zhu, "Pie: Parkour with implicit-explicit learning framework for legged robots," *IEEE Robotics and Automation Letters*, vol. 9, no. 11, pp. 9986–9993, 2024.

[42] E. Chane-Sane, J. Amigo, T. Flayols, L. Righetti, and N. Mansard, "Soloparkour: Constrained reinforcement learning for visual locomotion from privileged experience," in *Proceedings of The 8th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, vol. 270. PMLR, 2025, pp. 4268–4285.

[43] S. Kareer, N. Yokoyama, D. Batra, S. Ha, and J. Truong, "Vinl: Visual navigation and locomotion over obstacles," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 2018–2024.

[44] H. Lai, J. Cao, J. Xu, H. Wu, Y. Lin, T. Kong, Y. Yu, and W. Zhang, "World model-based perception for visual legged locomotion," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025, pp. 11 531–11 537.

[45] P. Li, H. Li, Y. Ma, L. Chang, X. Yang, R. Yu, Y. Zhang, Y. Cao, Q. Zhu, and G. Sartoretti, "Kivi: Kinesthetic-visuospatial integration for dynamic and safe egocentric legged locomotion," *arXiv preprint arXiv:2509.23650*, 2025.

[46] H. Su, H. Luo, S. Yang, K. Jiang, W. Zhang, and H. Chen, "Lipm-guided reinforcement learning for stable and perceptive locomotion in bipedal robots," in *2025 IEEE-RAS 24th International Conference on Humanoid Robots (Humanoids)*, 2025, pp. 1031–1038.

[47] D. Wang, X. Wang, X. Liu, J. Shi, Y. Zhao, C. Bai, and X. Li, "More: Mixture of residual experts for humanoid lifelike gaits learning on complex terrains," *arXiv preprint arXiv:2506.08840*, 2025.

[48] N. Rudin, J. He, J. Aurand, and M. Hutter, "Parkour in the wild: Learning a general and extensible agile locomotion policy using multi-expert distillation and rl fine-tuning," *arXiv preprint arXiv:2505.11164*, 2025.

[49] T. He, C. Zhang, W. Xiao, G. He, C. Liu, and G. Shi, "Agile but safe: Learning collision-free high-speed legged locomotion," in *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, July 2024.

[50] C. Zhang, J. Jin, J. Frey, N. Rudin, M. Mattamala, C. Cadena, and M. Hutter, "Resilient legged local navigation: Learning to traverse with compromised perception end-to-end," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 34–41.

[51] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" *Advances in neural information processing systems*, vol. 30, 2017.

[52] S. Jung, D. Gwak, B. Boots, and J. Hays, "Uncertainty-aware accurate elevation modeling for off-road navigation via neural processes," in *Conference on Robot Learning*. PMLR, 2025, pp. 2802–2822.

[53] H. Wang, Z. Shi, C. Zhu, Y. Qiao, C. Zhang, F. Yang, P. Ren, L. Lu, and D. Xuan, "Integrating learning-based manipulation and physics-based locomotion for whole-body badminton robot control," *arXiv preprint arXiv:2504.17771*, 2025.

[54] M. Hutter, C. Gehring, D. Jud, A. Lauber, C. D. Bellicoso, V. Tsounis, J. Hwangbo, K. Bodie, P. Fankhauser, M. Bloesch *et al.*, "Anymal-a highly mobile and dynamic quadrupedal robot," in *2016 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2016, pp. 38–44.

[55] LimX Dynamics, "TRON 1: The first multi-modal biped robot," https://www.limxdynamics.com/en/tron1, 2025, accessed: 2025-11-23.

[56] H. J. Lee, S. Hong, and S. Kim, "Integrating model-based footstep planning with model-free reinforcement learning for dynamic legged locomotion," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 11 248–11 255.

[57] W. Yu, D. Jain, A. Escontrela, A. Iscen, P. Xu, E. Coumans, S. Ha, J. Tan, and T. Zhang, "Visual-locomotion: Learning to walk on complex terrains with vision," in *5th Annual Conference on Robot Learning*, 2021.

[58] R. Yang, G. Yang, and X. Wang, "Neural volumetric memory for visual locomotion control," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 1430–1440.

[59] H. Kim, A. Mnih, J. Schwarz, M. Garnelo, A. Eslami, D. Rosenbaum, O. Vinyals, and Y. W. Teh, "Attentive neural processes," *arXiv preprint arXiv:1901.05761*, 2019.

[60] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[61] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa *et al.*, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.

[62] X. B. Peng and M. van de Panne, "Learning locomotion skills using deeprl: does the choice of action space matter?" in *Proceedings of the ACM SIGGRAPH / Eurographics Symposium on Computer Animation*, ser. SCA '17. New York, NY, USA: Association for Computing Machinery, 2017.

[63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[64] C. Zhang, W. Xiao, T. He, and G. Shi, "Wococo: Learning whole-body humanoid control with sequential contacts," in *Proceedings of The 8th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, vol. 270. PMLR, 2025, pp. 455–472.

[65] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. M. Kitani, C. Liu, and G. Shi, "Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning," in *Proceedings of The 8th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, vol. 270. PMLR, November 2025, pp. 1516–1540. [Online]. Available: https://proceedings.mlr.press/v270/he25b.html

[66] Z. Li, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath, "Reinforcement learning for versatile, dynamic, and robust bipedal locomotion control," *The International Journal of Robotics Research*, 2024.

[67] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning," in *Robotics: Science and Systems XIV, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA, June 26-30, 2018*, H. Kress-Gazit, S. S. Srinivasa, T. Howard, and N. Atanasov, Eds., 2018. [Online]. Available: http://www.roboticsproceedings.org/rss14/p08.html

[68] Z. Chen, M. Ji, X. Cheng, X. Peng, X. B. Peng, and X. Wang, "Gmt: General motion tracking for humanoid whole-body control," *arXiv preprint arXiv:2506.14770*, 2025.

[69] N. Messikommer, J. Xing, E. Aljalbout, and D. Scaramuzza, "Student-informed teacher training," in *7th Robot Learning Workshop: Towards Robots with Human-Level Abilities*, 2025. [Online]. Available: https://openreview.net/forum?id=q06EAFOwDI

[70] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science robotics*, vol. 5, no. 47, p. eabc5986, 2020.

[71] T. Pfau, A. G. de Rivaz, S. Brighton, and R. Weller, "Kinetics of jump landing in agility dogs," *The Veterinary Journal*, vol. 190, no. 2, pp. 278–283, 2011.

[72] S. McErlain-Naylor, M. King, and S. Allen, "Surface acceleration transmission during drop landings in humans," *Journal of Biomechanics*, vol. 118, p. 110269, 2021.

[73] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2017, pp. 23–30.

[74] C. Schwarke, M. Mittal, N. Rudin, D. Hoeller, and M. Hutter, "Rsl-rl: A learning library for robotics research," *arXiv preprint arXiv:2509.10771*, 2025.

[75] O. R. developers, "Onnx runtime," https://onnxruntime.ai/, 2025.

[76] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, "Learning agile and dynamic motor skills for legged robots," *Science Robotics*, vol. 4, no. 26, p. eaau5872, 2019.

[77] Y.-H. Shin, T.-G. Song, G. Ji, and H.-W. Park, "Actuator-constrained reinforcement learning for high-speed quadrupedal locomotion," *arXiv preprint arXiv:2312.17507*, 2023.

[78] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[79] M. Macklin, "Warp: A high-performance python framework for gpu simulation and graphics," https://github.com/nvidia/warp, March 2022, nVIDIA GPU Technology Conference (GTC).

[80] M. Seitzer, A. Tavakoli, D. Antic, and G. Martius, "On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks," in *International Conference on Learning Representations*, 2022.

[81] A. Chambolle, "An algorithm for total variation minimization and applications," *Journal of Mathematical imaging and vision*, vol. 20, no. 1, pp. 89–97, 2004.

[82] S. Khattak, T. Homberger, L. Bernreiter, J. Nubert, O. Andersson, R. Siegwart, K. Alexis, and M. Hutter, "Compslam: Complementary hierarchical multi-modal localization and mapping for robot autonomy in underground environments," *arXiv preprint arXiv:2505.06483*, 2025.

[83] J. Nubert, S. Khattak, and M. Hutter, "Graph-based multi-sensor fusion for consistent localization of autonomous construction robots," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 10 048–10 054.

[84] K. Chen, R. Nemiroff, and B. T. Lopez, "Direct lidar-inertial odometry: Lightweight lio with continuous-time motion correction," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 3983–3989.

[85] Z. Zhuang and H. Zhao, "Embrace contacts: humanoid shadowing with full body ground contacts," in *9th Annual Conference on Robot Learning*, 2025.

[86] L. Yang, X. Huang, Z. Wu, A. Kanazawa, P. Abbeel, C. Sferrazza, C. K. Liu, R. Duan, and G. Shi, "Omniretarget: Interaction-preserving data generation for humanoid whole-body loco-manipulation and scene interaction," *arXiv preprint arXiv:2509.26633*, 2025.

[87] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *ICLR*, 2021.

[88] F. Yang, P. Frivik, D. Hoeller, C. Wang, C. Cadena, and M. Hutter, "Spatially-enhanced recurrent memory for long-range mapless navigation via end-to-end reinforcement learning," *The International Journal of Robotics Research*, p. 02783649251401926, 2025.

[89] T. He, W. Xiao, T. Lin, Z. Luo, Z. Xu, Z. Jiang, J. Kautz, C. Liu, G. Shi, X. Wang, L. J. Fan, and Y. Zhu, "HOVER: versatile neural whole-body controller for humanoid robots," in *IEEE International Conference on Robotics and Automation, ICRA*. IEEE, 2025, pp. 9989–9996.