

## LEGGED ROBOTS

# Attention-based map encoding for learning generalized legged locomotion

Junzhe He<sup>1\*</sup>, Chong Zhang<sup>1\*</sup>, Fabian Jenelten<sup>1</sup>, Ruben Grandia<sup>2</sup>, Moritz Bächer<sup>2</sup>, Marco Hutter<sup>1</sup>

Dynamic locomotion of legged robots is a critical yet challenging topic in expanding the operational range of mobile robots. It requires precise planning when possible footholds are sparse, robustness against uncertainties and disturbances, and generalizability across diverse terrains. Although traditional model-based controllers excel at planning on complex terrains, they struggle with real-world uncertainties. Learning-based controllers offer robustness to such uncertainties but often lack precision on terrains with sparse steppable areas. Hybrid methods achieve enhanced robustness on sparse terrains by combining both methods but are computationally demanding and constrained by the inherent limitations of model-based planners. To achieve generalized legged locomotion on diverse terrains while preserving the robustness of learning-based controllers, this paper proposes an attention-based map encoding conditioned on robot proprioception, which is trained as part of the controller using reinforcement learning. We show that the network learns to focus on steppable areas for future footholds when the robot dynamically navigates diverse and challenging terrains. We synthesized behaviors that exhibited robustness against uncertainties while enabling precise and agile traversal of sparse terrains. In addition, our method offers a way to interpret the topographical perception of a neural network. We have trained two controllers for a 12-degrees-of-freedom quadrupedal robot and a 23-degrees-of-freedom humanoid robot and tested the resulting controllers in the real world under various challenging indoor and outdoor scenarios, including ones unseen during training.

## INTRODUCTION

Humans and legged animals inhabit almost every corner of our planet and are adept at traversing various terrains in the wild. Similarly, legged robots hold vast potential for navigating complex natural landscapes that are typically inaccessible to other wheeled or tracked mobile robots. Yet, navigating challenging terrains demands precise planning when possible footholds are sparse [e.g., on construction debris (1)] and robustness in the presence of uncertainties and disturbances (2, 3). Furthermore, the locomotion controller must be able to generalize across diverse terrains.

Toward this goal, various methods have been explored, including learning-based ones, such as deep reinforcement learning (DRL); model-based ones, such as model predictive control (MPC); and hybrid approaches that combine both. However, achieving generalized legged locomotion across diverse terrains with both precision and robustness remains an open problem. Here, we offer an attention-based learning framework to train robust and generalized controllers that can precisely navigate on various terrains.

DRL has emerged as a powerful tool for enabling robust and agile legged locomotion on challenging terrains. By training an actuator model through supervised learning and appropriately randomizing the training environment, Hwangbo *et al.* (4) successfully transferred the dynamic motions learned through trial and error in simulation to the real world. Lee *et al.* (2) and Siekmann *et al.* (5) developed robust learning-based controllers for blind traversal of rough terrains and stairs by quadrupedal and bipedal robots, respectively. Regarding perceptive locomotion, Rudin *et al.* (6) trained holistic controllers that directly mapped the observations to joint-level outputs on uneven terrains by massively parallel DRL. Miki *et al.* (3) further advanced robust perceptive quadrupedal locomotion in the wild by

learning to filter out unreliable perception using a history of proprioceptive data. In addition, DRL approaches have enabled legged robots to perform dynamic parkour maneuvers across various structured terrains (7–13). However, these approaches struggle on sparse terrains because it is hard for the DRL algorithm to find valid footholds and learn from them. To tackle this problem, recent DRL-based work has achieved locomotion on sparse terrains through curriculum tuning to guide the training process with human intuition (14), but it only overfit a small range of terrains with a single multilayer perceptron (MLP) policy network and could not generalize. Another work enabled a bipedal robot to walk on fake stepping stones (QR-code tags on the flat ground) through MLP-based foothold feasibility prediction (15) but has not yet demonstrated locomotion on real terrains or generalization across various terrains.

Alternative to reinforcement learning (RL), model-based methods have been explored for decades to produce versatile movements and smooth trajectories that adhere to the kinematic and dynamic constraints of the robot. Some works (16, 17) have enabled stable locomotion on rough terrains, such as stairs, steps, inclines, etc. To achieve agility on uneven terrains or with various payloads, more adaptable frameworks such as MPC are used for locomotion tasks (18–22) by solving optimal control problems over a long horizon. More recent endeavors have demonstrated stronger traversal performance on the real robot, achieving real-time versatile legged locomotion on uneven and sparse terrains by predicting accurate footholds (23–25). Despite this progress, model-based controllers usually solve deterministic optimal control and thus struggle with measurement uncertainty, which may lead to degraded performance under noisy perception, imprecise motor actuation, etc. Most of the model-based methods use simplified dynamic and kinematic models, which may lead to failures because of model mismatches during deployment. Recent works have shown that solving optimization with full dynamics on the fly is possible (26, 27), but robustly handling measurement uncertainty still remains an open problem for these methods.

<sup>1</sup>Robotic Systems Lab, ETH Zurich, 8092 Zurich, Switzerland. <sup>2</sup>Disney Research Zurich, Stampfenbachstrasse 48, 8006 Zurich, Switzerland.

\*Corresponding author. Email: junzhe@ethz.ch (J.H.); chozhang@ethz.ch (C.Z.)

To combine the advantages of both model- and learning-based controllers, hybrid methods have been proposed. Kwon *et al.* (28) proposed to train a dynamic model for model-based controllers. Some works (29–31) warm-started nonlinear solvers from learned initialization to speed up convergence. Other works used DRL to generate footholds that are then tracked by model-based controllers (32, 33). These approaches leveraged learning to improve or bootstrap MPC while still running an MPC as the main controller and hence remained fragile to the uncertainty faced in real-world deployments. In contrast, to combine generalization capability and precise foothold predictions from a model-based planner and the robustness from a learning-based controller, Deep Tracking Control (DTC) (1) trained a DRL controller to track the reference state trajectories generated by model-based optimization (25), effectively overcoming model mismatch, slippage, and deformable terrains. Nonetheless, compared with holistic approaches, which directly map observations to actions, these hybrid controllers suffer from the complexity of the hierarchical structure. For example, during training, DTC required running the model-based planner on a central processing unit (CPU) and the training pipeline on a graphics processing unit (GPU) and took 14 days of training for full convergence. Furthermore, DTC required running the MPC controller during deployment, which is computationally demanding and relies on the performance of the model-based planner—the height map is fed into the model-based controller, which might generate infeasible guidance under degraded perception.

In summary, a holistic learning-based approach that can achieve precise, robust, and generalized locomotion on sparse terrains is missing. In addition, despite all of the advancements in legged locomotion under different task specifications, no approach has yet demonstrated the successful sim-to-real transfer of perceptive controllers to achieve dynamic locomotion on sparse terrains for both quadrupedal and bipedal robots.

Here, we propose to train an attention-based map encoding for generalized legged locomotion, which consists of two levels—a convolutional neural network (CNN) that embeds point-wise local terrain features in a robot-centric height map sampled from elevation mapping (34) and a multihead attention (MHA) (35) module that queries point-wise map features—and combines them with proprioceptive observations. MHA is a neural representation mechanism used in transformers (35), a deep learning architecture that enables highly effective handling of sequential and multimodal information. Recent works on legged locomotion learning have also introduced the transformer architecture to embed multimodal observations (36) or capture time-sequential features (37), achieving locomotion on unstructured terrains such as grasslands, slopes, and flat ground. Instead of using a complete transformer model, we investigated how MHA can enable generalized legged locomotion across diverse terrains by focusing on the areas that matter most. In our framework, a low-level CNN learned effective extraction of local features for various terrains, and the MHA module selected useful terrain points on the basis of the attention conditioned on proprioception, leading to a compact and generalizable representation of high-dimensional observations. We observed that the selected points indicate the future footholds in an emergent way without any explicit supervision on footholds.

To effectively synthesize our locomotion policies, we designed a two-stage training pipeline. In the first stage, we trained the controller on base terrains (defined in the “Terrains” section in the

Supplementary Materials) with perfect perception. This stage is crucial for initializing the map encoding learning and leads to basic locomotion skills under ideal conditions. In the second stage, we introduced more challenging terrains (defined in the “Terrains” section in the Supplementary Materials) with perception noise and drift to enhance the generalization to real-world conditions.

We demonstrate that our proposed control framework has resulted in substantial advancements over the state of the art (SOTA) in precise and generalized dynamic locomotion on a broad range of terrains (1) while achieving robustness against uncertainties and model mismatch (Movie 1). It also produces an interpretable representation of map scans that can be graphically visualized. The same framework was deployed on both a quadrupedal robot ANYmal-D (38) and a humanoid robot Fourier GR-1 (39), as shown in Fig. 1.

## RESULTS

### Precise and generalized locomotion

We tested our learned policies on both the GR-1 and ANYmal-D robots in simulation (Fig. 2) and extensive real-world experiments (Figs. 1 and 3) across a diverse range of terrains, including ones never encountered during the training phase. Figure 2A illustrates the controller’s performance on GR-1, trained from scratch during stage 1, where it exhibited precise foot placements across various terrain types, such as grid stones, pallets, beams, and gaps. Despite being exclusively trained on base terrain types partially shown in Fig. 2A, the controller demonstrated its ability to generalize to unseen terrains, including pentagon stones, single-column stones, narrow pallets, and consecutive gaps, shown in Fig. 2B. This highlights our controller’s capability to transfer the learned behaviors to terrains unseen during training.

A similar pattern of generalization was observed in the ANYmal-D controller. The controller was able to adapt and perform effectively on unseen terrains, as shown in Fig. 2D, despite being trained only on base terrains shown in Fig. 2E (note that the kinematics of ANYmal-D differs from that of GR-1, requiring adjustments in the terrain selection during training; more terrain details can be found in the “Terrains” section in the Supplementary Materials). This demonstrates the adaptability of our approach to managing different embodiments with varying kinodynamics, enabling generalization capabilities despite hardware variations.

To further enhance the policies’ precision on a wider range of terrains and robustness in real-world experiments, we fine-tuned them in stage 2 by incorporating additional terrain types and introducing disturbances and uncertainties. This fine-tuning process improved the controller’s ability to navigate more complex and unpredictable environments, reinforcing its stability and adaptability under real-world conditions. We demonstrate the resulting controllers in Fig. 2C, where both GR-1 and ANYmal-D achieved a 100% success rate in traversing a challenging obstacle parkour (with disturbances and uncertainties) designed by Grandia *et al.* (24) that had not been encountered during training. We enabled dynamic humanoid locomotion across such mixed sparse terrains with an online controller. The same controllers also achieved high robustness in real-world deployments, demonstrated in Figs. 1 and 3, unifying the properties of SOTA model-based and learning-based approaches. We deployed the fine-tuned policies for both ANYmal-D and GR-1 zero-shot on the real hardware. To validate the controller’s precision and robustness, we tested on various sparse terrains, including gaps,



**Movie 1. Attention-based map encoding enables generalized legged locomotion.**

stepping stones, and beams, as depicted in Fig. 3. Our results represent advances in robotic locomotion, demonstrating that our holistic DRL-based controller can achieve precise and generalized performance across a wide range of challenging terrains. Before our work, no other holistic DRL-based controllers achieved such a level of generalization.

### Agility and recovery reflexes by whole-body coordination

Our learned policies demonstrated advanced agility and robustness on the real robots, as shown in Fig. 4. Learning whole-body motion control has aided ANYmal-D and GR-1 in actively using the knees (Fig. 4A) and arms (Fig. 4C), respectively, for enhanced agility. The amplitude and frequency of GR-1's arm swing not only follow the gait but also depend on the terrain, as can be observed by comparing Fig. 4C to Fig. 4D. The emergent recovery behaviors can save the robots from slippage (Fig. 4, B and E) or stabilize themselves on shaky supports (Fig. 4D). When the current footholds and velocity commands made the next step hard to land (Fig. 4F), we observed that GR-1 performed a single-leg switch hop on one stepping stone, lifting the previous contact foot into the air to successfully reach the subsequent stone. All of these reflex behaviors, which enhance the system stability, are difficult to obtain from model-based methods, because they typically rely on contact state machines and handcrafted heuristics (24, 25, 40, 41).

### Versatile velocity tracking

Our learned policies can track velocity commands in a versatile way, with detailed evaluations and benchmarks shown in the “Simulation-based evaluations” section. This enabled us to maneuver ANYmal-D

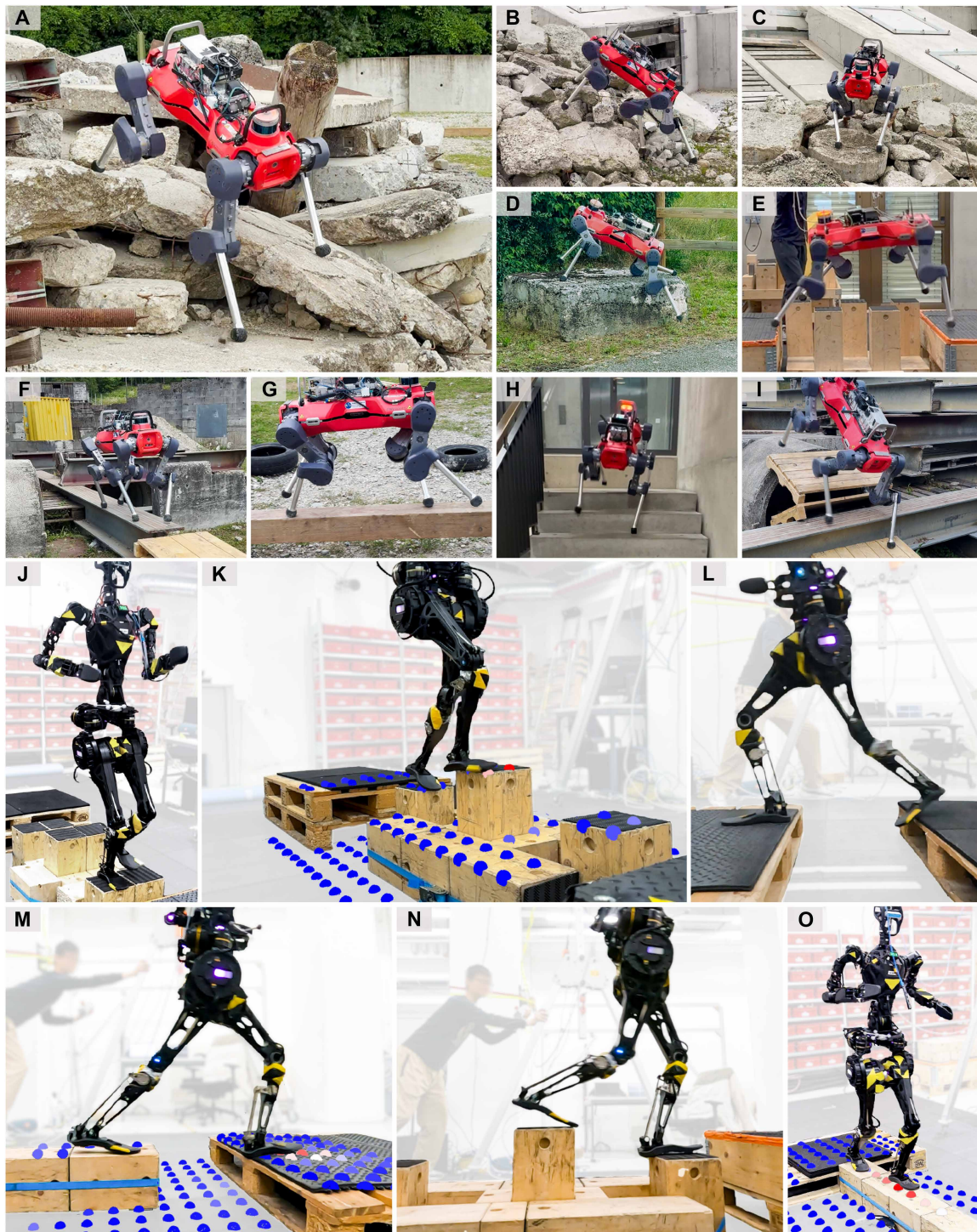
on debris. In Fig. 5A, the robot overcame sparse terrains with movable supports and showcased omnidirectional versatility. We can also command GR-1 with different velocities on challenging terrains, such as one row of uneven stepping stones or shaky balance beams, leading to different gait patterns and whole-body behaviors, as depicted in Fig. 5 (B and C). Such versatility can expand the operational range of our robots in complex terrains and confined spaces.

### Simulation-based evaluations

#### Benchmark with DTC and baseline RL controller

In the first evaluation experiment, we evaluated the performance of three different controllers that can dynamically navigate on sparse terrains with an ANYmal-D—the proposed method, DTC (1), and baseline-rl (14), which is a recent learning-based controller that demonstrates successful hardware experiments on stepping stones and beams. To do so, we used three comparing criteria similar to (1): the velocity tracking performances; the success, failure, and stuck rates on terrains that the compared controllers are trained on; and the success rates on individual terrains. All controllers were deployed in the same simulated environment with observation noises and drifts sampled with the same random seed.

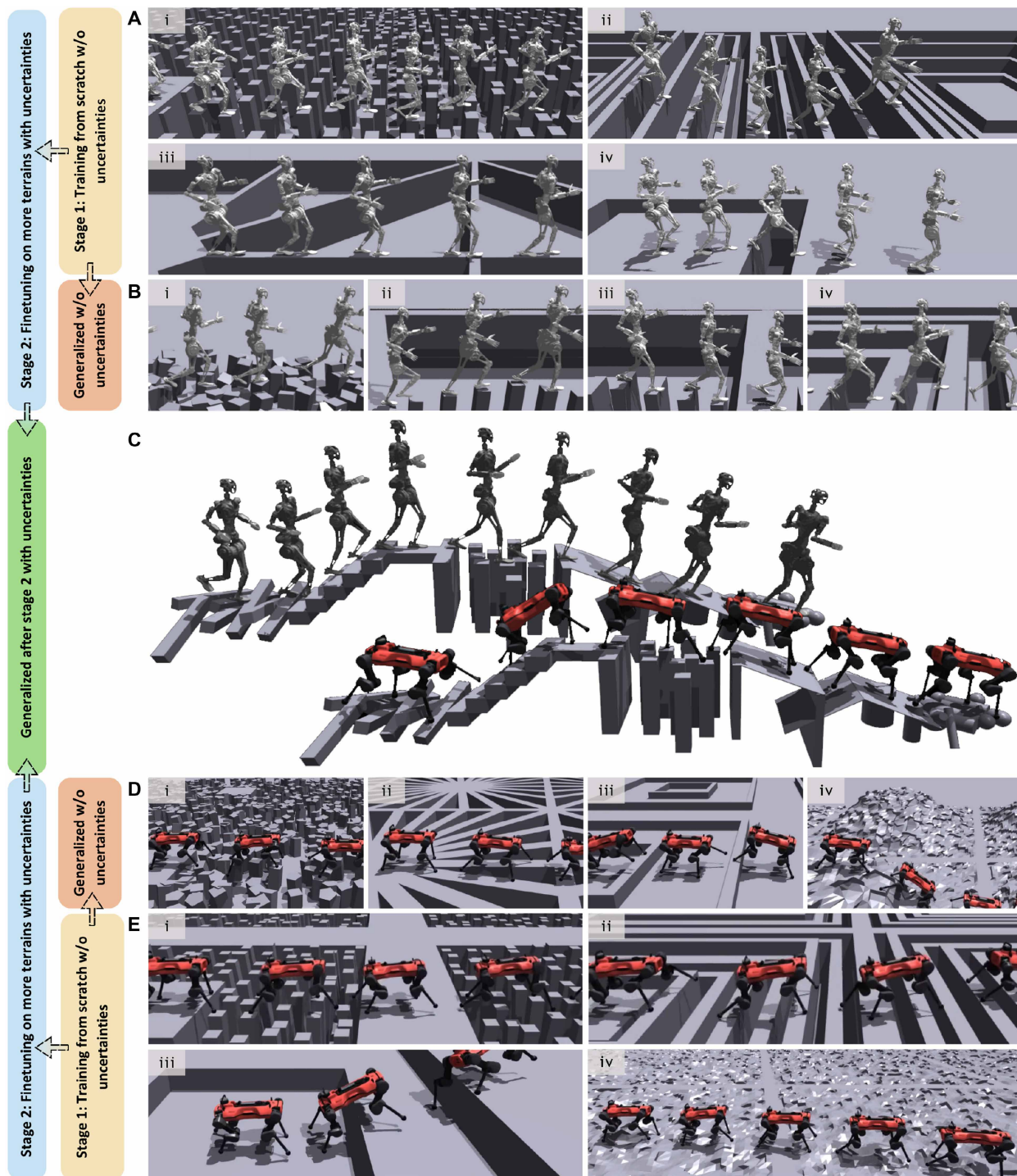
Because baseline-rl is designed under a goal-reaching setup and there is currently no learning-based velocity-tracking controller tailored for sparse terrains that can parallel DTC's performance on hardware in the literature, we only compared the proposed method and DTC for velocity tracking, as shown in Fig. 6A(i). The agents with different controllers were deployed with constant forward velocity commands on selected sparse terrains traversable by both controllers. The tracking error was only computed for surviving



**Fig. 1. Learning-interpretable, generalizable, agile, and robust legged locomotion on diverse terrains.** Our controllers enabled ANYmal-D (A to I) and GR-1 (J to O) to dynamically traverse diverse challenging terrains. Highly interpretable point-wise map encodings are graphically visualized in (K), (M), and (O), where more intensely red colors represent higher attention weights, indicating the next foothold.

agents. Our approach demonstrated substantially lower tracking errors except for gaps with small velocity commands, where the agents hesitate to proceed for both controllers. Noticeably, DTC exhibited high tracking errors with large velocity commands. This is mainly

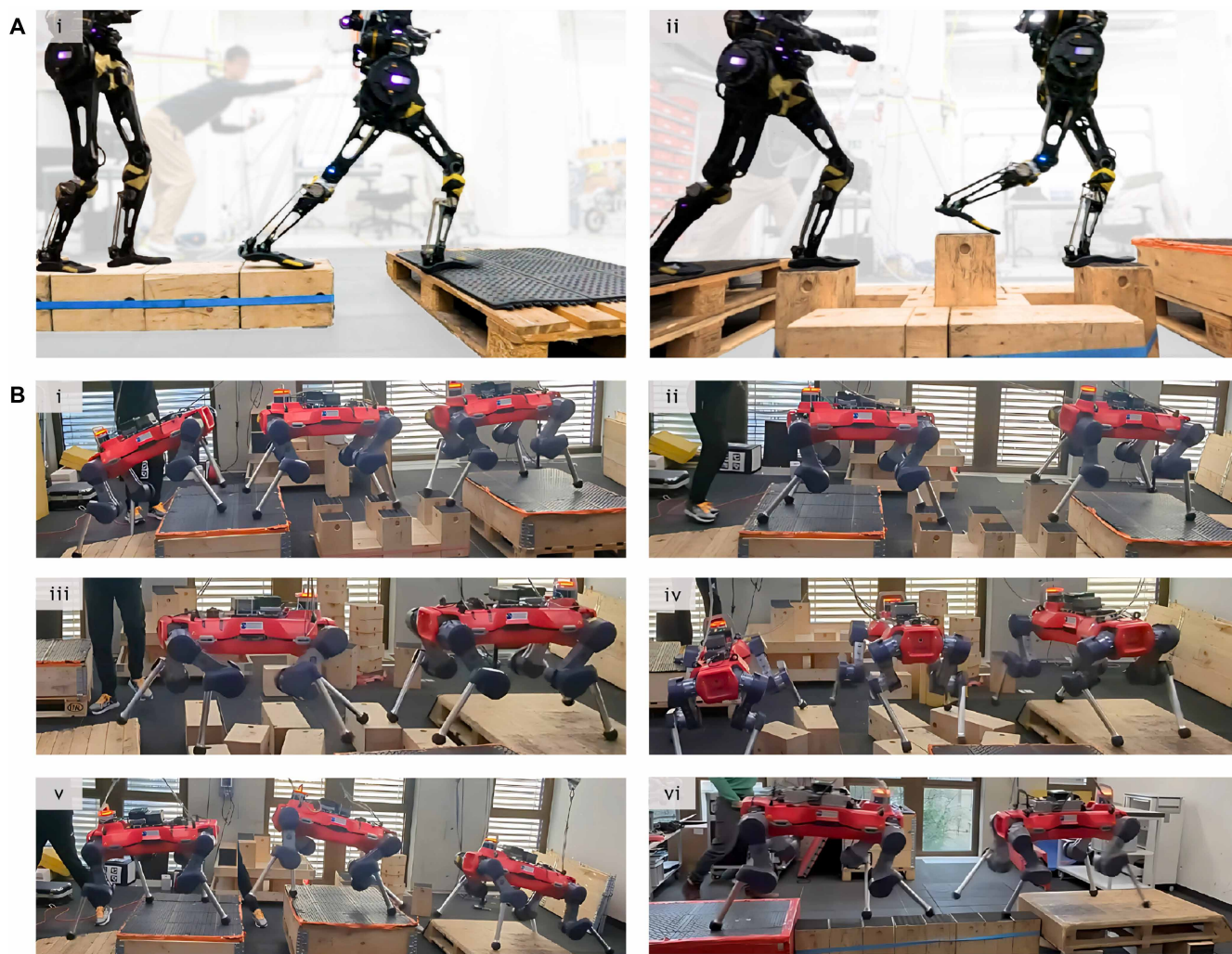
because DTC uses constant gait frequency, resulting in further away and thus harder-to-track footholds when velocity commands are higher. Our method, on the other hand, can adjust gait frequencies on different terrains with various velocity commands. Figure 6A(ii)



**Fig. 2. Precise and generalized locomotion in simulation.** (A) Selection of base terrains for GR1 stage 1 training, including grid stones (i), pallets (ii), beams (iii), and gaps (iv). (B) Selection of fine-tuning terrains for GR1 stage 2 training, including pentagon stones (i), single-column stones (ii), narrow pallets (iii), and consecutive gaps (iv). (C) GR1 and ANYmal-D on obstacle parkour (24). (D) Selection of fine-tuning terrains for ANYmal-D stage 2 training, including pentagon stones (i), beams (ii), rings (iii), and rough hills (iv). (E) Selection of base terrains for ANYmal-D stage 1 training, including grid stones (i), pallets (ii), pits (iii), and rough ground (iv).

shows the success, failure, and stuck rates on a complete set of terrains that the three controllers were trained on, where our method shows 26.5 and 77.3% higher success rates compared with DTC and baseline RL, respectively. The reason for the low overall success rates

of DTC and baseline RL can be explained in Fig. 6A(iii). DTC has lower than 20% success rates on grid stones (20 cm-by-20 cm randomly placed square stones), small grid stones (12 cm-by-12 cm randomly placed square stones), and narrow beams (15-cm width),



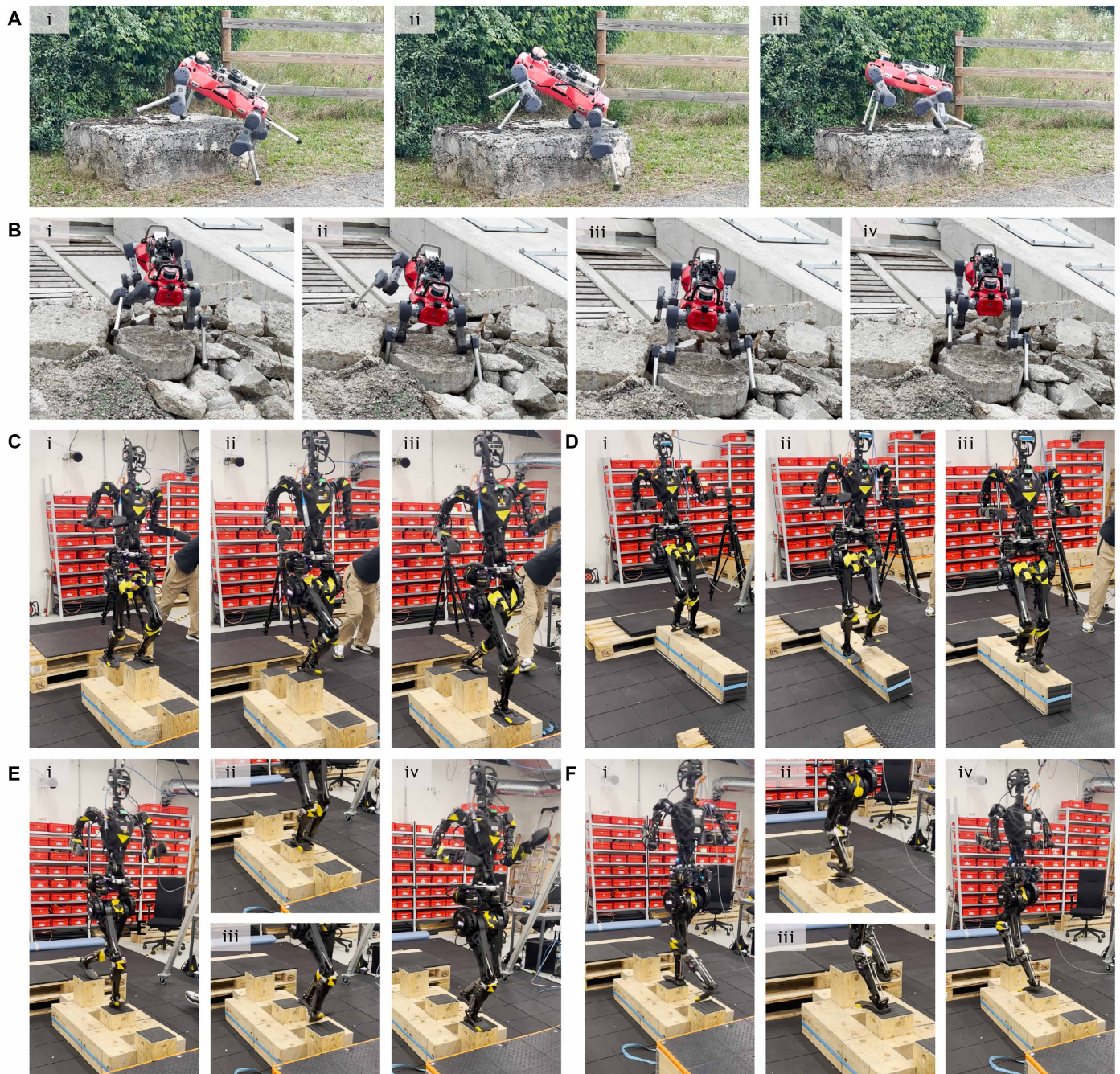
**Fig. 3. Precise and generalized locomotion on hardware.** (A) Selection of sparse terrains for GR1, including beam and gap (i) and single-column stones with height differences (ii). (B) Selection of sparse terrains for ANYmal-D, including stepping stones (i), stepping stones with height differences (ii), randomly placed stepping stones forward (iii), randomly placed stepping stones sideways (iv), boxes and gaps (v), and a 19-cm-wide beam (vi).

which is mainly because the high-level controller provides infeasible foothold guidance on these terrains because the support surface is smaller than the threshold that the model-based planner can accept as a steppable area. The baseline-rl controller is overfitted to grid stones and cannot generalize to unseen terrains, explaining its low success rates on other terrains.

To better visualize the generalization ability of our method versus baseline-rl, movie S1 demonstrates our controller (stage 1) and the baseline-rl controller on three terrains: the grid stone terrain used for training, the grid stone terrain with random beams, and the pentagon stone terrain shown in Fig. 2D(i). We show that our controller operated successfully on the grid stone terrain with random beams and the pentagon stone terrain (unseen during training) where baseline-rl failed. Although baseline-rl performed well on the training terrain, it struggled to generalize to unseen terrains even with very slight variations. These results indicate that generalization on sparse terrains is challenging for previous RL-based solutions and that the observed generalization can be attributed to our policy architecture design rather than the RL process alone.

### Ablation study on two-stage training

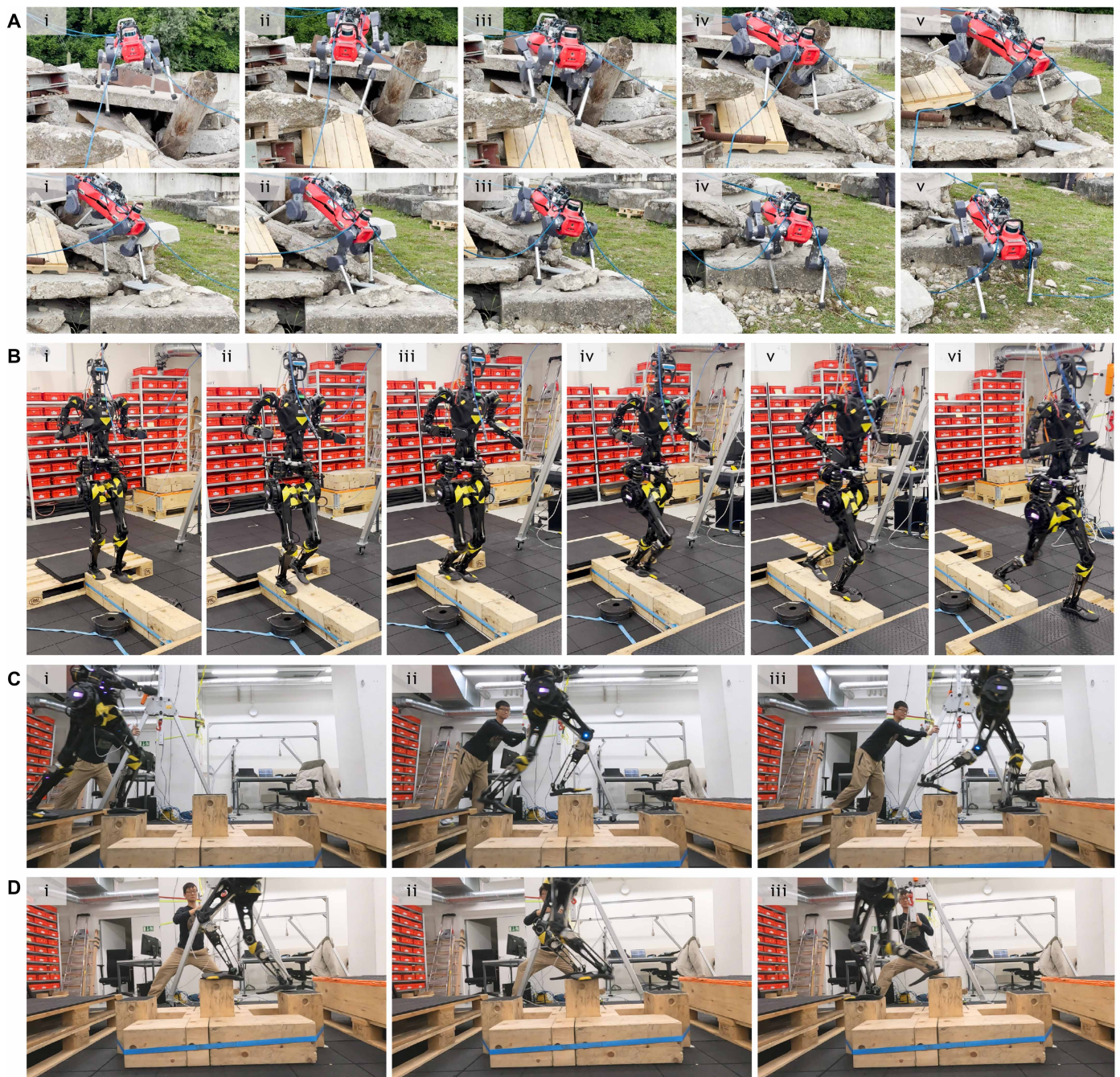
To justify the proposed two-stage training pipeline and the necessity of initializing the map encoding learning in stage 1 and introducing more terrains and uncertainties only in stage 2, we closely investigated and compared the training and deployment performance of three controllers. The first controller was trained with the proposed pipeline: first trained on base terrains with ground truth observations (observations with no noise or drift, only available in simulation) and then fine-tuned on all terrains (base terrains and fine-tuning terrains) with sensory drifts and noises. The second controller (C2) was trained on all terrains from scratch with ground truth observations, whereas the third controller (C3) was trained on base terrains from scratch with sensory drifts and noises. We used the terrain level (6) to compare their training performance. The terrains have 10 difficulty levels, indexed from 0 to 9. All robots were randomly assigned a terrain type and a level. The robot was upgraded to the next level if it walked out of the borders of its assigned terrain and downgraded to a lower level otherwise. Robots solving the highest level were then reset to a randomly selected level, which averaged to



**Fig. 4. Agility and recovery reflexes by whole-body coordination.** (A) ANYmal-D using its knee to climb up a large rock while rotating the torso. (B) ANYmal-D recovered from penetration of feet into the shaky debris as a result of slippage by knee support. (C) GR-1 traversing one row of 19-cm-wide uneven stepping stones with natural arm swing to aid the agile motions. (D) GR-1 stabilizing itself on a shaky, unfixed 19-cm-wide balance beam. (E) GR-1 encountered a slippage while traversing one row of uneven stepping stones and reacted with a fast step forward. (F) GR-1 encountered an inappropriate foothold as a result of left-biasing velocity commands while traversing a row of uneven stepping stones. With insufficient space for the left foot to land after the right foot's placement, GR-1 switched the foot contact in the air and successfully reached the subsequent stone with the right foot.

level 6 if all robots had solved the most difficult terrain. Perfect training will show a terrain level curve that first overshoots level 6 (robots attempting to upgrade to higher terrain levels) and then converges to level 6 (on the basis of the stationary distribution when all robots solve their terrains). Figure 6B(i) indicates that our proposed method

resulted in a convergent terrain level showing that most of the robots solve the most difficult level. However, C2 could not reach the same terrain level, meaning that the agents failed to upgrade to higher terrain levels. C3 could not converge back to 6, indicating that the agents could not solve the most difficult levels. We then deployed



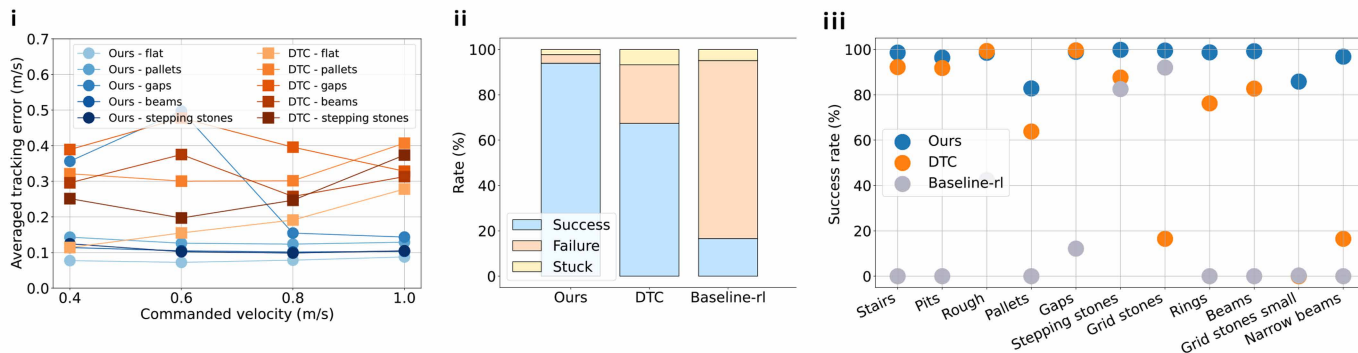
**Fig. 5. Versatile velocity tracking.** Our learned controllers demonstrated versatile velocity tracking capabilities on ANYmal-D and GR-1. (A) ANYmal-D maneuvering on the debris, overcoming sparse terrain with movable supports, and showcasing omnidirectional versatility. (B) GR-1 accelerating on the shaky balance beam when the velocity command changed from 0.7 to 1.5 m/s, taking longer strides. (C) With a 1.5 m/s forward velocity command, GR-1 had one step per stepping stone. (D) With a 0.7 m/s forward velocity command, GR-1 had two steps on each stepping stone.

the controllers and compared the success rates on different terrains, as depicted in Fig. 6B(ii). Our method showed substantially higher success rates than C2 and C3 on almost all terrains. C2 showed the worst performance even with perfect perception during testing. Although C3 showed satisfactory success rates on stairs, pits, and rough terrain, it performed worse on sparse terrains (grid stones, beams, etc.).

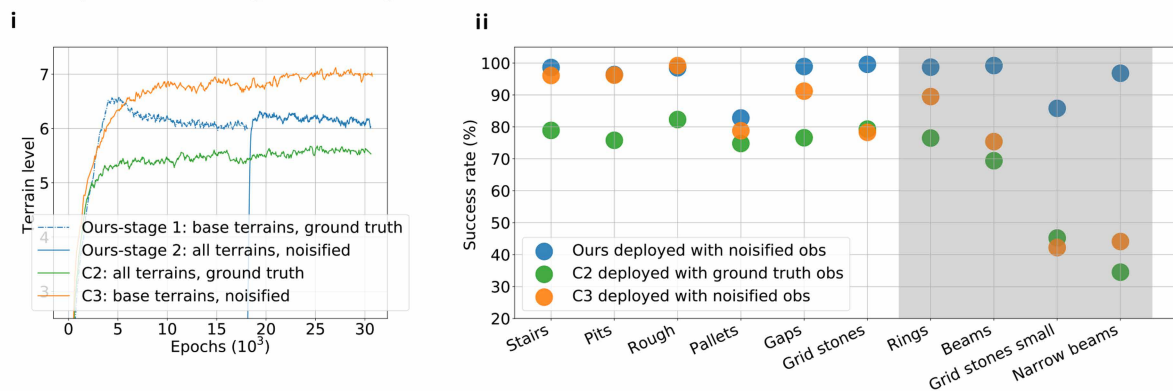
#### **Ablation study on network structure**

Our attention-based map encoding consists of two levels: a low-level CNN that embeds local terrain features and an MHA module that queries point-wise local features and combines them with proprioceptive observations. To show the efficacy of using an MHA module the way we propose, we compared our method with a transformer encoder similar to (36). To show the necessity of point-wise attention,

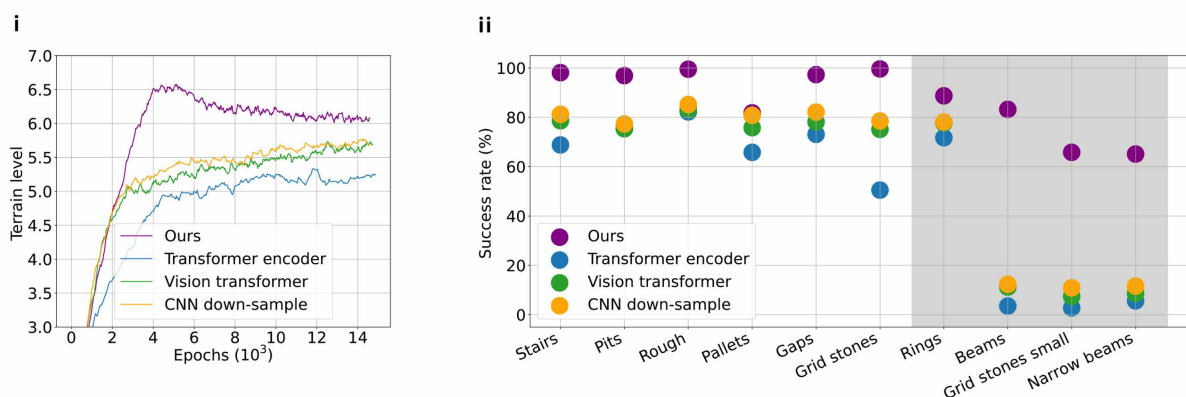
**A Benchmark with DTC and baseline RL**



**B Ablation study on two-stage training**



**C Ablation study on network structure**



**Fig. 6. Simulation-based evaluations.** We only evaluated the performance of our approach and benchmarked it with other methods on ANYmal-D. **(A)** Benchmark with DTC and baseline-rl. (i) Our method shows overall lower velocity tracking errors for different forward velocity commands on the selected terrains. (ii) Our method shows a substantially higher success rate and lower stuck and failure rates (a trial is counted as “successful” if the robot could walk out of the border of the terrain within a complete episode, “failed” if undesirable contacts happen, and “stuck” otherwise) on a combination of all training terrains. (iii) Our method demonstrates higher overall success rates on individual terrains. **(B)** Ablation study on two-stage training. (i) Terrain level training curves for the proposed two-stage training (ours), training from scratch on all terrains (base + fine-tuning terrains) with ground truth observations (C2), and training from scratch on the base terrains with sensory drift and noise (C3). Ours shows the best convergent behavior. (ii) Our method shows higher overall success rates on individual terrains, where the white background indicates the base terrains and the gray background fine-tuning terrains. **(C)** Ablation study on the network structure. (i) Terrain level training curves for different methods on base terrains. Ours shows the best convergent behavior. (ii) Our method shows higher overall success rates on individual terrains.

we down-sampled the map inputs with another CNN instead of only extracting local features without down-sampling. As another comparison, we used a vision transformer encoder (42) to process the map scans. The detailed structures can be found in the “Network structures” section in the Supplementary Materials. Similar to the “Ablation study on two-stage training” section, we compared the

training and deployment performance of the transformer encoder, the CNN down-sampling, and the vision transformer with our method through the terrain level and success rates, as shown in Fig. 6C(i),(ii), respectively. From the comparison, our method demonstrates higher convergent terrain levels during training and success rates during deployment. Our method shows substantially

higher success rates on unseen terrains, indicating better generalization than the other potential network structures.

Moreover, we found that our network architecture can facilitate exploration in RL training, enabling convergence without any forms of high-level guidance [e.g., the references in (1)] or exploration rewards [e.g., the ones in (14)]. To illustrate this, we compared the training performance on grid stones between our policy and an MLP policy under identical reward functions, as demonstrated in movie S2. Our policy efficiently learned feasible walking motions, whereas the MLP policy struggled with the training terrain. This indicates that policy architecture plays a crucial role in facilitating exploration.

### Interpretable attention-based map encoding

To further elucidate the interpretability provided by the proposed attention-based map encoding, we provide a detailed visualization of the attention weights from the MHA module across different types of terrains. The MHA module pays more attention to steppable areas based on the proprioceptive information of the robots and allocates higher attention weights to corresponding map scans. These visualizations highlight how the model prioritizes specific regions in the environment, guiding the controller to navigate complex discontinuous terrains. Figure 7A showcases the attention weights of the fine-tuned controller from stage 2 on a mixed terrain, which combines elements from various terrain types. The attention weights are concentrated around the next steppable region, indicating that the MHA module acts as a guidance to direct the controller toward feasible footholds, ensuring stability. Figure 7B depicts the attention weights when the robots are commanded to move in different directions on various terrains encountered during stage 1 training, including grid stones [Fig. 7B(i) forward, Fig. 7B(v) sideways, and Fig. 7B(ix) turning], pallets [Fig. 7B(ii) forward, Fig. 7B(vi) sideways, and Fig. 7B(x) turning], single beams [Fig. 7B(iii) forward, Fig. 7B(vii) sideways, and Fig. 7B(xi) turning], and gaps [Fig. 7B(iv) forward, Fig. 7B(viii) sideways, and Fig. 7B(xii) turning]. These terrains, characterized by their varying structural features, challenge the model to distribute attention effectively to critical regions to support footholds on the basis of the robots' kinematics and dynamics. Noticeably, the paid attention can refuse to follow infeasible velocity commands. For example, the agent is commanded to turn, but the attention still stays on the beam rather than on the ground to follow the command, as demonstrated in Fig. 7B(xi). In addition, Fig. 7C illustrates how the stage 1 controller generalizes to unseen terrains, such as pentagon stones [Fig. 7C(i)], narrow pallets [Fig. 7C(ii)], single-column stones [Fig. 7C(iii)], and consecutive gaps [Fig. 7C(iv)], highlighting the generalization of the attention mechanism in adapting to previously unseen environments. These visualizations reflect the model's ability to dynamically adjust its focus on the basis of the robot's proprioceptive information and command directions, promoting efficient and safe locomotion across diverse and unpredictable terrains. The consistent pattern of attention allocation supports our claim that the MHA module enhances both interpretability and generalization capabilities in complex locomotion tasks.

### DISCUSSION

In this work, we achieved generalized, agile, and robust legged locomotion using RL, mapping proprioceptive and exteroceptive observations directly to joint-level actions with a neural network. The key to achieving this is an attention-based map encoding module that

processes the map on the basis of proprioception and provides a representation that learns to focus on potential future footholds. The subsequent policy then translates this generalized representation into whole-body motions, enabling precise movements on sparse terrains. In addition, we developed a two-stage training pipeline to further enhance the generalization capability of the controller. The resulting controllers enabled both quadrupedal and humanoid robots to traverse diverse terrains and demonstrate highly interpretable neural encoding of terrain perception.

We enabled generalized legged locomotion while preserving the robustness of learning-based controllers. Previously, such generalization only appeared on quadrupedal locomotion controllers with model-based planning (1, 24). Hence, on top of the SOTA performance on quadrupedal and humanoid robots, this work demonstrates the possibility of achieving comparable generalization using DRL while overcoming the limits of model-based approaches when dealing with uncertainties and model errors.

Our policy network design mirrors the modular functions in model-based methods: The map encoding module implicitly selects future footholds like contact planners by attention, and the subsequent policy acts like a whole-body controller tracking the planned contacts. Yet, by integrating the entire controller from observations to joint-level actions, we can leverage data-driven learning to tackle the issues of model-based methods: computational burden, model mismatch, and violation of assumptions. We can also tune the controller as a whole rather than tuning separate modules, reducing system complexity.

This work also has limitations. First, training the policies can take several days, and parameter tuning becomes inefficient because of the training costs. Second, we used the 2.5-dimensional (2.5D) height map representation, which may be inapplicable to certain scenarios, such as confined spaces (43). Last, we focused on locomotion in this work; although the arms need to be used for manipulation as well, we have not studied how to balance the needs for both locomotion and manipulation with the arms.

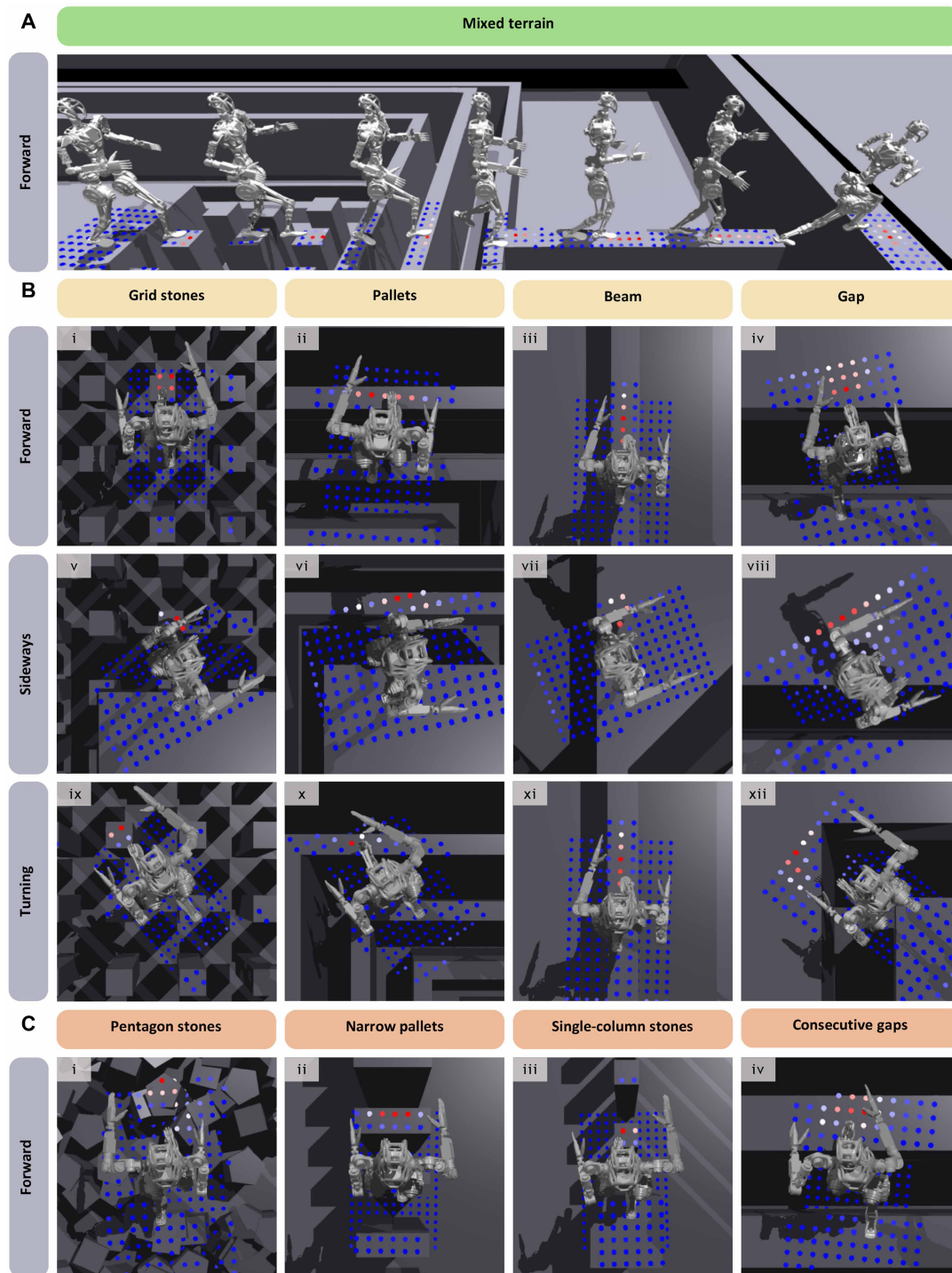
In the future, we will explore how we can improve training efficiency and develop effective 3D representations that can generalize to a broader range of scenarios. We also expect that the attention mechanism can be extended to locomanipulation tasks, including opening doors, moving obstacles, climbing with the help of hands, etc.

## MATERIALS AND METHODS

### Motivation

A general abstract model can be derived to describe the commonalities among model-based, learning-based, and hybrid methods discussed in the Introduction as pictured in Fig. 8A. In summary, an encoder observes proprioceptive and exteroceptive information and outputs a latent representation that is then fed into subsequent policy modules to generate actions.

For model-based controllers, the encoder can be interpreted as a high-level planner that predicts the base trajectories, joint angles, and footholds within its prediction horizon; the subsequent policy module can be interpreted as a whole-body controller that generates corresponding actions to track the solutions from the planner. Similarly, for learning-based controllers, neural networks are used to fuse the proprioceptive and exteroceptive information, and the subsequent layers generate actions. For hybrid methods, a network as the encoder to generate footholds was used in (32), with a model-based

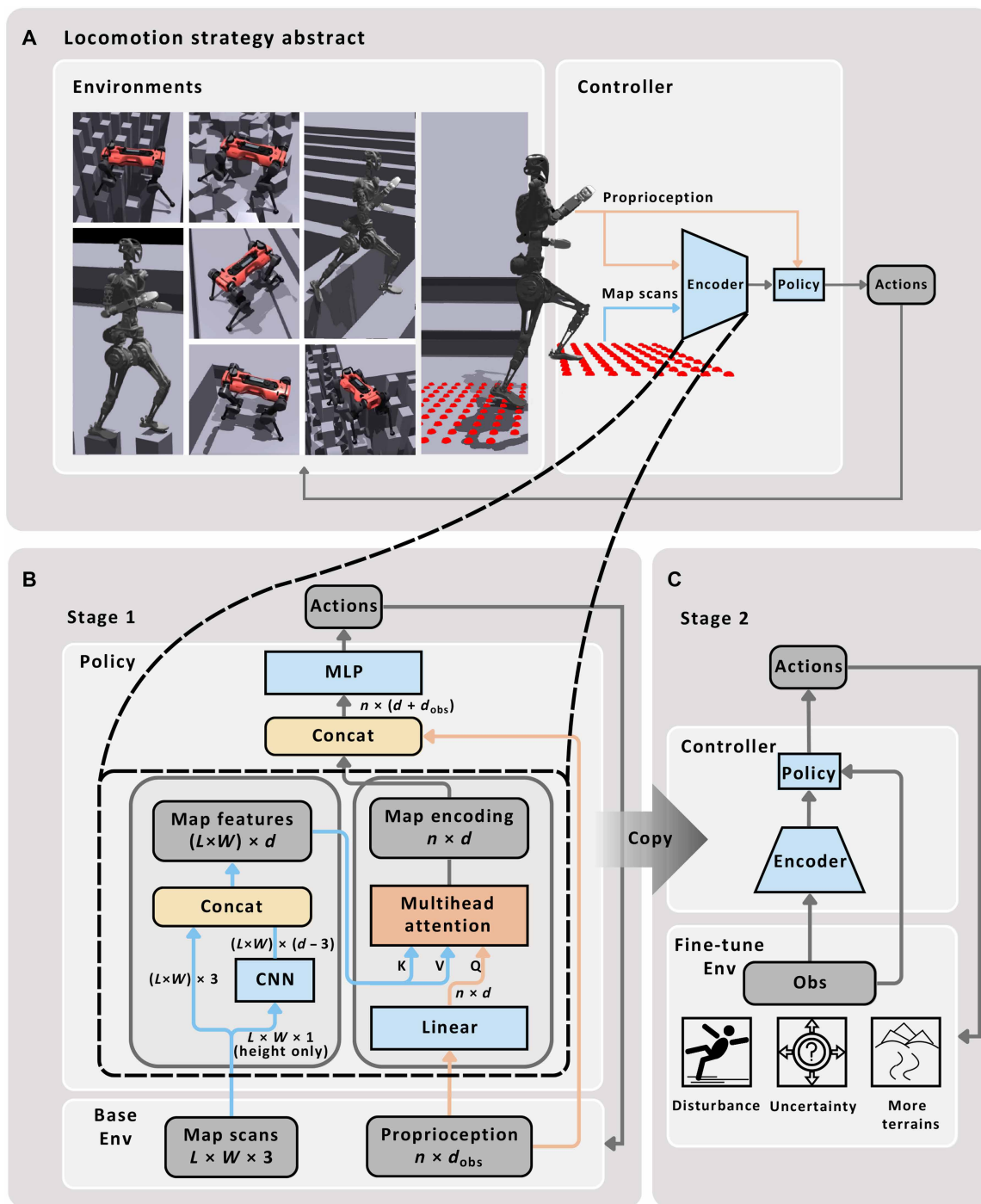


**Fig. 7. Attention weight visualization.** We visualized the height scans with associated attention weights on each scan, where higher-intensity red colors indicate higher attention. (A) Attention weights on a mixture of terrains for a stage 2 controller. (B) Attention weights on individual base terrains for a stage 1 controller with forward, sideways, and turning velocity commands. (C) Attention weights on individual fine-tuning terrains for a stage 1 controller with forward velocity commands.

whole-body controller serving as the subsequent module to track these footholds; TAMOLS (Terrain-Aware Motion Optimization for Legged Systems) was leveraged in (1) as the encoder to generate base trajectories and footholds and used an MLP policy to track

TAMOLS solutions. Despite the variety in these locomotion strategies, they share intrinsic similarities in their control pipelines.

Building on this general abstract framework, we can analyze the factors contributing to differences in robustness, precision, and



**Fig. 8. Proposed control pipeline.** (A) Locomotion strategy abstract. (B) Stage 1 training. (C) Stage 2 fine-tuning with disturbance, uncertainty, and additional harder terrains.

generalizability among the control strategies mentioned above. Model-based controllers have exhibited great precision, with model-based planners serving as the observation encoder. Model-based planners leverage trajectory optimization to predict footholds and base motions from map scans and robot states while adhering to the kinematic and dynamic constraints. However, they tend to be less robust because of their inability to handle model mismatches and uncertainties effectively. In contrast, learning-based methods demonstrate

greater robustness, because their encoders use neural networks, which are better at managing uncertainties when trained with proper domain randomization. Nevertheless, they are generally less precise, because it is difficult to impose hard constraints on neural networks and less generalizable to previously unseen terrains, given their susceptibility to overfitting. With these considerations, we focus on the encoder as a key component in developing a control strategy that balances precision, robustness, and generalizability across

terrains. In this work, we introduce an attention-based neural locomotion controller architecture that uses an MHA module to encode exteroception conditioned on proprioception and thus is capable of predicting precise foot placements on the basis of current robot states, as illustrated in Fig. 8B.

### MHA-based map encoding

MHA (35) is a neural representation mechanism that has revolutionized deep learning. At its core, MHA enables the model to focus on the important parts of the inputs and enriches the representation by generating multiple parallel attention outputs, colloquially known as “heads.” MHA therefore enhances interpretability by identifying important inputs and improves generalization by mitigating the effects of variations in less important inputs. These features align with our requirements for interpretability and generalization in locomotion controllers that can work across diverse terrains.

Technically, the inputs of MHA in our work include a query and a collection of key-value pairs. The output is determined by calculating a weighted sum of the values, with each value’s weight being determined by the compatibility of the query with the corresponding key. This process can be interpreted as paying more attention to the values, for which the corresponding keys are the most relevant to the query inputs. This interpretation can be naturally transferred to our case, where we expect the MHA module to focus more on the optimal steppable areas, informed by current proprioceptive data and commands. To this end, we use the proprioception embedding as the query and point-wise local features as the key-value pairs to get a map encoding conditioned on the proprioception. This structure is more expressive or, put another way, more suitable to model intricate relationships between the proprioceptive and exteroceptive information than MLPs because the attention weights paid to the map scans are state dependent. The outputs of MHA are then translated to joint-level actions by the subsequent policy module.

As demonstrated in Fig. 8B, the map scans ( $L \times W \times 3$ ,  $L$  points long,  $W$  points wide, and 3D coordinates in the robot frame for each point) make the exteroceptive observations. The  $z$  values of the points are first processed by a CNN that consists of two layers with zero padding to keep the original dimensionality and a kernel size of 5 to extract the local features for each point. The first layer has 16 hidden units, and the second has  $d - 3$  hidden units, where  $d$  is the dimension of the MHA module (64 in our case). Then, we concatenate the output of CNN [ $L \times W \times (d - 3)$ ] with the 3D coordinates to get point-wise local features of shape  $LW \times d$ . These local features extract the neighborhood features around each map point, enabling the subsequent MHA module to pay point-wise attention. Meanwhile, the proprioception ( $1 \times d_{\text{obs}}$ ) goes first through a linear layer that outputs a proprioception embedding ( $1 \times d$ ). Then, with the proprioception embedding as the query ( $Q$ , with  $n = 1$  as MHA’s query length) and the local features as the keys and values ( $K$  and  $V$ , respectively), the MHA module outputs the map encoding, with each head processing  $d/h$  dimensions of the inputs, where  $h$  is the number of heads.

### Two-stage training pipeline

To learn a robust and generalizable map encoding, we designed a two-stage training pipeline that progressively refines the controller’s capabilities. In the first stage, the controller was trained on base terrains with perfect perception. This stage warms up the map encoding learning and allows the controller to acquire locomotion skills using ground truth sensing.

In the second stage, we introduced more complex terrains with disturbances and uncertainties. These terrains simulated more realistic environments where perception may be imperfect, shaping the learned motions to be more adaptable and resilient. By exposing the robot to a broader range of terrains with added disturbances, this stage improved the controller’s generalization across diverse terrains and also enhanced its robustness against real-world uncertainties, which ensured that the final learned map encoding can operate effectively in unseen environments.

### Observation space

The policy network observes proprioception information and map scans in the robot-centric base frame (torso link for ANYmal-D and pelvis link for GR-1 are considered as the base), including the base linear velocity  $\mathbf{v}_b$ , angular velocity  $\boldsymbol{\omega}_b$ , gravity vector  $\mathbf{g}_b$ , joint positions  $\mathbf{q}_j$  and velocities  $\dot{\mathbf{q}}_j$ , previous actions  $\mathbf{a}_{t-1}$  inferred by the control policy, and vector map scans surrounding the robot’s base (a total of  $L \times W \times 3$  points). The critic network observes the same information but without noise. The symbols are defined in Table 1.

We did not use privileged observations for our policy, unlike prior works (2, 3) that combined them with teacher-student distillation to develop controllers robust to severely compromised perception. Our work primarily addresses traversal over diverse terrains, where robots are not expected to traverse stepping stones or balance beams under degraded perception. That said, our approach can be extended with privileged observations and teacher-student distillation if extreme scenarios are to be considered.

### Reward functions

The reward is a weighted sum of 14 terms for ANYmal-D and 16 terms for GR-1, detailed in Table 2. We divided our reward settings into three categories, including task, regularization, and style rewards.

The task rewards consist of tracking commanded linear and angular velocity while avoiding failures on terrains, where  $n_{\text{collision}}$  is defined as the number of collisions between specific body parts with the terrain and  $n_{\text{termination}}$  is defined as the number of early terminations. We terminated the episode when there was a collision between the torso and the terrain or bad orientation of the torso.

The regularization rewards were designed to smooth the actions and avoid drastic motions, overtorque, and overextensions. We used soft constraints with soft limits in rewards for joint limit violations to avoid hindering exploration, whereas other safe RL methods (44–46) may achieve similar effects with hard constraints.

To generate more natural movements, we also introduced style rewards, where we penalized unwanted velocities, stomping, foot slippage, jumping gaits, and tilted body parts. To confine the movements of arms for GR-1, we also introduced a reward to penalize deviation from default joint positions above a specific threshold for arm joints. After we obtained a baseline controller using the reward functions above, we introduced a stand-still penalty that penalizes joint motions while the robots are in the stance phase in stage 2. This is because we observed shaking and slight tapping on the hardware during the standing phase. The reward is shown with marks in Table 2.

### Training environment

We used a custom version of proximal policy optimization (PPO) (47) for training and a two-stage training pipeline, where we first trained a controller with ground truth observations for both the actor and critic and then fine-tuned the controller with noises and disturbances for

**Table 1. Nomenclature.**

$\mathbf{v}_i^*$	Commanded base linear velocity of rigid body $i$
$\mathbf{v}_i$	Linear velocity of rigid body $i$
$\boldsymbol{\omega}_i^*$	Commanded base angular velocity of rigid body $i$
$\boldsymbol{\omega}_i$	Angular velocity of rigid body $i$
$\mathbf{a}_j$	Desired $j$ th joint position inferred by the controller
$\mathbf{q}_{0,j}$	Default $j$ th joint position
$\mathbf{q}_j$	Joint positions
$\dot{\mathbf{q}}_j$	Joint velocities
$\ddot{\mathbf{q}}_j$	Joint accelerations
$\boldsymbol{\tau}_j$	Joint torques
$\mathbf{q}_{\text{lim},j}$	Joint position limits
$\dot{\mathbf{q}}_{\text{lim},j}$	Joint velocity limits
$\boldsymbol{\tau}_{\text{lim},j}$	Torque limits
$n_{\text{termination}}$	Number of terminations
$n_{\text{collision},i}$	Number of collisions of rigid body $i$
$n_{\text{zero\_contact}}$	Number of zero contact
$\mathbf{F}_f$	Net contact force on foot $f$
$C_f$	Contact state of foot $f$
$\mathbf{v}_f$	Velocity of foot $f$
$\mathbf{g}_i$	Gravity vector of rigid body $i$
$I_*$	Index of rigid body named *

the actor and ground truth information for the critic. The actor and critic use the same network structure as shown in Fig. 8B, where they share the same encoder module but use different subsequent MLPs. The actor MLP maps the MHA output to joint actions, whereas the critic MLP maps to a value. The hyperparameters are detailed in the Supplementary Materials (“PPO parameters” section).

### Domain randomization

We introduce noise to observations. At each simulation step, a customized noise is sampled from a uniform distribution and added to each observation term, except for the previous actions and velocity commands. The map scans also have random drifts sampled from a normal distribution for each terrain at the beginning of training. The perturbations are mainly designed to improve the robustness against sensor drifts during deployment. We also introduced artificial pushes by resetting the twist of the robots in simulation. To improve the controller’s robustness against payload and friction variations, we also randomized the torso mass and the friction coefficient of each contact foot.

### Terrains and curriculum

We used different terrain settings for each training stage, detailed in the Supplementary Materials (“Terrains” section). For each training stage, we used a curriculum introduced in (6). At the start of the training, all robots were randomly assigned a terrain type and a difficulty level (10 in total). During training, the robots that managed to walk out of their assigned terrains got upgraded to the next level and downgraded to a lower level otherwise. The difficulty level of each terrain was tuned heuristically such that the supporting surface

is big enough for the map scan resolution and the difficulty is not greater than the maximum robot capability. For example, for a 10-cm map resolution, the stepping stones should be larger than 10 cm; the gap width should not exceed the length of the quadruped robot/the maximum possible feet distance of the humanoid robot while assuming a walking gait (at least one foot on the ground).

### Training

The encoder network parameters are the same for ANYmal-D and Fourier GR-1, except that the dimensions of map scans for GR-1 (17 by 11) are smaller than those of ANYmal-D (26 by 16). We used  $d = 64$  for the MHA dimension,  $n = 1$  for the target sequence length, and  $h = 16$  for the number of heads. The policy network maps the map encoding and proprioceptive observation to actions, parameterized by an MLP. The action space comprises target joint positions tracked by a low-level PD controller (4). The joint actuators exhibit substantial delays: The ANYmal-D robot uses series elastic actuators, and the GR-1 robot uses motors with gear ratios of 51, 80, or 100. For both robots, joint delays are a genuine concern, and we identified the actuator dynamics before modeling their torque outputs in simulation. We used an actuator network in (4) for ANYmal-D and a dc motor model with inertia, friction, and damping identified for GR-1. We then trained the ANYmal controller through massive parallelization with 4096 robots for 18,000 epochs in stage 1 and 3600 epochs in stage 2. With 24 s of training time per epoch, the total training time is 6 days on an Nvidia Tesla A100-40GB GPU, which is a reduction of roughly 60% of training time compared with DTC. For the GR-1 controller, we trained for 15,000 epochs in stage 1 and 3200 epochs in stage 2, with 14 s per epoch, resulting in a total training time of 3.5 days on an Nvidia RTX 4090 GPU.

**Table 2. Definition of reward terms for ANYmal-D and GR-1** The  $x$  axis points forward, and the  $z$  axis points downward. The body, joint, and feet indices can be found in the Supplementary Materials. N/A, not applicable; –, term not used.

Reward terms	Functions	ANYmal-D		GR-1	
		Weights	Indices	Weights	Indices
<i>Task</i>					
Linear velocity tracking	$\exp(-\ \mathbf{v}_{xy,i}^* - \mathbf{v}_{xy,i}\ ^2)$	5.0	$i = [i_{\text{torso}}]$	5.0	$i = [i_{\text{torso}}]$
Angular velocity tracking	$\exp(-\ \boldsymbol{\omega}_{z,i}^* - \boldsymbol{\omega}_{z,i}\ ^2)$	3.0	$i = [i_{\text{torso}}]$	3.0	$i = [i_{\text{torso}}]$
Termination penalty	$-n_{\text{termination}}$	200	N/A	200	N/A
Collision penalty	$-n_{\text{collision},i}$	1	$i = [i_{\text{shank}}]$	–	–
<i>Regularization</i>					
Action rate	$-\ \mathbf{a}_{j,t} - \mathbf{a}_{j,t-1}\ ^2$	$5.0 \times 10^{-3}$	$j = [0:12]$	$5.0 \times 10^{-3}$	$j = [0:23]$
Joint acceleration penalty	$-\ \ddot{\mathbf{q}}_j\ ^2$	$2.5 \times 10^{-7}$	$j = [0:12]$	$1 \times 10^{-6}$	$j = [15:18, 19:22]$
Joint torque penalty	$-\ \boldsymbol{\tau}_j\ ^2$	$2.0 \times 10^{-5}$	$j = [0:12]$	$1 \times 10^{-4}$ $5 \times 10^{-5}$	$j = [15, 19]$ $j = [2,3,8,9]$
Joint position limits	$-\max( \mathbf{q}_j  - 0.9\mathbf{q}_{\text{lim},j}, 0)$	1.0	$j = [0:12]$	10	$j = [0:23]$
Joint velocity limits	$-\max( \dot{\mathbf{q}}_j  - 0.9\dot{\mathbf{q}}_{\text{lim},j}, 0)$	1.0	$j = [0:12]$	0.1	$j = [0:23]$
Joint torque limits	$-\max( \boldsymbol{\tau}_j  - 0.8\boldsymbol{\tau}_{\text{lim},j}, 0)$	0.2	$j = [0:12]$	$2 \times 10^{-3}$	$j = [0:23]$
<i>Style</i>					
Linear velocity penalty	$-\mathbf{v}_{z,i}^2$	1.0	$i = [i_{\text{torso}}]$	–	–
Angular velocity penalty	$-\ \boldsymbol{\omega}_{xy,i}\ ^2$	$5.0 \times 10^{-2}$	$i = [i_{\text{torso}}]$	$5.0 \times 10^{-2}$	$i = [i_{\text{torso}}]$
Contact force penalty	$-\max(\ \mathbf{F}_f\  - 700, 0)$	$2.5 \times 10^{-5}$	$f = [0:4]$	–	–
Foot slippage penalty	$-c_f^* \ \mathbf{v}_f\ $	0.5	$f = [0:4]$	1.0	$f = [0, 1]$
Joint deviation penalty	$\max(\ \mathbf{q}_j - \mathbf{q}_{0,j}\ ^2 - 0.25, 0.0)$	–	–	0.5	$j = [15:23]$
No fly	$-n_{\text{zero\_contact}}$	–	–	5.0	$f = [0, 1]$
Straight body	$-\ \mathbf{g}_j\ ^2$	–	–	3.0	$i = [i_{\text{torso}}, i_{\text{pelvis}}, i_{\text{feet}}]$
Standing joint positions penalty*	$-\ \mathbf{q}_j^* - \mathbf{q}_j\ $	0.1	$j = [0:12]$	–	–
Standing joint velocity penalty*	$-\ \dot{\mathbf{q}}_j^* - \dot{\mathbf{q}}_j\ $	0.5	$j = [0:12]$	0.2	$j = [0:23]$

\*The rewards marked are only activated during stage 2 fine-tuning.

## Deployment

We used ANYmal-D and GR-1 for our experiments and deployed the policies at a frequency of 50 Hz for both hardware platforms. For ANYmal-D, the control policy inference was done on a single Intel core-i7 8850H CPU, and the elevation mapping (34) ran on an on-board Nvidia Jetson. For GR-1, the policy inference was done on the Intel core-i7 13700h CPU, and the map scans were sampled through ray-casting on a predesigned terrain mesh given the robot's pose captured by the Qualysis Motion Capture system (48).

## Statistical analysis

We computed the averaged tracking errors shown in Fig. 6A(i), success rates shown in Fig. 6 [A(ii),(iii), B(ii), and C(ii)], and failure rates and stuck rates shown in Fig. 6A(ii) with data collected from 4096 environments in the Isaac Gym simulation. The tracking errors were sampled at 50 Hz, and the success rates, failure rates, and stuck rates were collected at the end of each entire episode. The terrain level shown in Fig. 6 [B(i) and C(i)] was computed by averaging the terrain difficulty levels (10 in total) for 4096 environments at the end of each entire episode during training.

## Supplementary Materials

### The PDF file includes:

Materials and Methods  
Figs. S1 to S3  
Table S1  
Legends for movies S1 and S2

### Other Supplementary Material for this manuscript includes the following:

Movies S1 and S2

## REFERENCES AND NOTES

1. F. Jenelten, J. He, F. Farshidian, M. Hutter, DTC: Deep Tracking Control. *Sci. Robot.* **9**, eadh5401 (2024).
2. J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, M. Hutter, Learning quadrupedal locomotion over challenging terrain. *Sci. Robot.* **5**, eabc5986 (2020).
3. T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, M. Hutter, Learning robust perceptive locomotion for quadrupedal robots in the wild. *Sci. Robot.* **7**, eabk2822 (2022).
4. J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, M. Hutter, Learning agile and dynamic motor skills for legged robots. *Sci. Robot.* **4**, eaau5872 (2019).
5. J. Siekmann, K. Green, J. Warila, A. Fern, J. W. Hurst, "Blind bipedal stair traversal via sim-to-real reinforcement learning" in *Robotics: Science and Systems XVII*, D. A. Shell, M. Toussaint, M. A. Hsieh, Eds. (RSS Foundation, 2021); 10.15607/RSS.2021.XVII.061.

6. N. Rudin, D. Hoeller, P. Reist, M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning" in *Proceedings of the 5th Conference on Robot Learning*, A. Faust, D. Hsu, G. Neumann, Eds., vol. 164 of *Proceedings of Machine Learning Research* (PMLR, 2022), pp. 91–100.
7. N. Rudin, D. Hoeller, M. Bjelonic, 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)* (IEEE, 2022), pp. 2497–2503.
8. D. Hoeller, N. Rudin, D. Sako, M. Hutter, ANYmal parkour: Learning agile navigation for quadrupedal robots. *Sci. Robot.* **9**, eadi7566 (2024).
9. Z. Zhuang, Z. Fu, J. Wang, C. Atkeson, S. Schwertfeger, C. Finn, H. Zhao, "Robot parkour learning" in *Proceedings of the 7th Conference on Robot Learning (CoRL)*, J. Tan, M. Toussaint, K. Darvish, Eds., vol. 229 of *Proceedings of Machine Learning Research* (PMLR, 2023), pp. 73–92.
10. X. Cheng, K. Shi, A. Agarwal, D. Pathak, "Extreme parkour with legged robots" in *2024 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2024), pp. 11443–11450.
11. Z. Zhuang, S. Yao, H. Zhao, "Humanoid parkour learning" in *Proceedings of the 8th Conference on Robot Learning*, P. Agrawal, O. Kroemer, W. Burgard, Eds., vol. 270 of *Proceedings of Machine Learning Research* (PMLR, 2024), pp. 1975–1991.
12. C. Zhang, W. Xiao, T. He, G. Shi, "WoCoCo: Learning whole-body humanoid control with sequential contacts" in *Proceedings of the 8th Conference on Robot Learning*, P. Agrawal, O. Kroemer, W. Burgard, Eds., vol. 270 of *Proceedings of Machine Learning Research* (PMLR, 2024), pp. 455–472.
13. H. Duan, B. Pandit, M. S. Gadde, B. Van Marum, J. Dao, C. Kim, 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.
14. C. Zhang, N. Rudin, D. Hoeller, M. Hutter, "Learning agile locomotion on risky terrains" in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2024).
15. H. Duan, A. Malik, M. S. Gadde, J. Dao, A. Fern, J. W. Hurst, "Learning dynamic bipedal walking across stepping stones" in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2022), pp. 6746–6752.
16. A. Shkolnik, M. Levashov, I. R. Manchester, R. Tedrake, Bounding on rough terrain with the LittleDog robot. *Int. J. Robot. Res.* **30**, 192–215 (2011).
17. P. Fankhauser, M. Bjelonic, C. Dario Bellicoso, T. Miki, 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; 10.1109/ICRA.2018.8460731.
18. D. Dimitrov, A. Sherikov, P.-B. Wieber, "A sparse model predictive control formulation for walking motion generation" in *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2011), pp. 2292–2299; 10.1109/IROS.2011.6095035.
19. M. Neunert, C. de Cousaz, F. Furrer, M. Kamel, F. Farshidian, R. Y. Siegwart, J. Buchli, "Fast nonlinear model predictive control for unified trajectory optimization and tracking" in *2016 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2016), pp. 1398–1404.
20. F. Farshidian, E. Jelavic, A. Satapathy, M. Gifftthaler, J. Buchli, "Real-time motion planning of legged robots: A model predictive control approach" in *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)* (IEEE, 2017), pp. 577–584.
21. J. Di Carlo, P. M. Wensing, B. Katz, G. Bledt, S. Kim, "Dynamic locomotion in the MIT Cheetah 3 through convex model-predictive control" in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2018), pp. 1–9; 10.1109/IROS.2018.8594448.
22. R. Grandia, F. Farshidian, R. Ranftl, M. Hutter, "Feedback MPC for torque-controlled legged robots" in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2019), pp. 4730–4737; 10.1109/IROS40897.2019.8968251.
23. C. Mastalli, W. Merkt, G. Xin, J. Shim, M. Mistry, I. Havoutis, S. Vijayakumar, Agile maneuvers in legged robots: A predictive control approach. arXiv:2203.07554 [cs.RO] (2022).
24. R. Grandia, F. Jenelten, S. Yang, F. Farshidian, M. Hutter, Perceptive locomotion through nonlinear model-predictive control. *IEEE Trans. Robot.* **39**, 3402–3421 (2023).
25. F. Jenelten, R. Grandia, F. Farshidian, M. Hutter, TAMOLS: Terrain-Aware Motion Optimization for Legged Systems. *IEEE Trans. Robot.* **38**, 3395–3413 (2022).
26. G. Kim, D. Kang, J.-H. Kim, S. Hong, H.-W. Park, Contact-implicit model predictive control: Controlling diverse quadruped motions without pre-planned contact modes or trajectories. *Int. J. Rob. Res.* **44**, 486–510 (2025).
27. V. Dh'Edin, A. K. Chinnakkonda Ravi, A. Jordana, H. Zhu, A. Meduri, L. Righetti, B. Schölkopf, M. Khadiv, "Diffusion-based learning of contact plans for agile locomotion" in *2024 IEEE-RAS 23rd International Conference on Humanoid Robots (Humanoids)* (IEEE, 2024), pp. 637–644; 10.1109/Humanoids58906.2024.10769875.
28. T. Kwon, Y. Lee, M. Van De Panne, Fast and flexible multilegged locomotion using learned centroidal dynamics. *ACM Trans. Graph.* **39**, 46:1–46:17 (2020).
29. O. Melon, M. Geisert, D. Surovik, I. Havoutis, M. Fallon, "Reliable trajectories for dynamic quadrupeds using analytical costs and learned initializations" in *2020 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2020), pp. 1410–1416; 10.1109/ICRA40945.2020.9196562.
30. D. Surovik, O. Melon, M. Geisert, M. Fallon, I. Havoutis, "Learning an expert skill-space for replanning dynamic quadruped locomotion over obstacles" in *Proceedings of the 2020 Conference on Robot Learning*, J. Kober, F. Ramos, C. Tomlin, Eds., vol. 155 of *Proceedings of Machine Learning Research* (PMLR, 2021), pp. 1509–1518.
31. O. Melon, R. Orsolino, D. Surovik, M. Geisert, I. Havoutis, M. Fallon, "Receding-horizon perceptive trajectory optimization for dynamic legged locomotion with learned initialization" in *2021 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2021), pp. 9805–9811; 10.1109/ICRA48506.2021.9560794.
32. S. Gangapurwala, M. Geisert, R. Orsolino, M. Fallon, I. Havoutis, RLOC: Terrain-aware legged locomotion using reinforcement learning and optimal control. *IEEE Trans. Robot.* **38**, 2908–2927 (2022).
33. Z. Xie, X. Da, B. Babich, A. Garg, M. v. de Panne, "GLiDE: Generalizable quadrupedal locomotion in diverse environments with a centroidal model" in *Algorithmic Foundations of Robotics XV*, S. M. LaValle, J. M. O'Kane, M. Otte, D. Sadigh, P. Tokekar, Eds. (Springer International Publishing, 2023), pp. 523–539.
34. T. Miki, L. Wellhausen, R. Grandia, F. Jenelten, T. Homberger, M. Hutter, "Elevation mapping for locomotion and navigation using GPU" in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2022), pp. 2273–2280.
35. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, I. Polosukhin, "Attention is all you need" in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett, Eds. (Curran Associates, 2017), vol. 30; [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
36. R. Yang, M. Zhang, N. Hansen, H. Xu, X. Wang, "Learning vision-guided quadrupedal locomotion end-to-end with cross-modal transformers" in *The 10th International Conference on Learning Representations (ICLR)* (ICLR, 2022); <https://openreview.net/forum?id=nhnJ3oo6AB>.
37. I. Radosavovic, T. Xiao, B. Zhang, T. Darrell, J. Malik, K. Sreenath, Real-world humanoid locomotion with reinforcement learning. *Sci. Robot.* **9**, eadi9579 (2024).
38. M. Hutter, C. Gehring, D. Jud, A. Lauber, C. D. Bellicoso, V. Tsounis, J. Hwangbo, K. Bodie, P. Fankhauser, M. Bloesch, R. Diethelm, S. Bachmann, A. Melzer, M. Hoepfner, "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; 10.1109/IROS.2016.7758092.
39. Fourier Intelligence, Fourier GR1 (2024); <https://www.ftai.com/products-gr1>.
40. M. Chignoli, D. Kim, E. Stanger-Jones, S. Kim, "The MIT humanoid robot: Design, motion planning, and control for acrobatic behaviors" in *2020 IEEE-RAS 20th International Conference on Humanoid Robots (Humanoids)* (IEEE, 2021), pp. 1–8.
41. S. Fahmi, V. Barasul, D. Esteban, O. Villarreal, C. Semini, VITAL: Vision-based terrain-aware locomotion for legged robots. *IEEE Trans. Robot.* **39**, 885–904 (2023).
42. A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale" in *The 9th International Conference on Learning Representations (ICLR)* (ICLR, 2021); <https://openreview.net/forum?id=YicbFdNTTy>.
43. T. Miki, J. Lee, L. Wellhausen, M. Hutter, "Learning to walk in confined spaces using 3D representation" in *2024 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, 2024), pp. 8649–8656; 10.1109/ICRA57147.2024.10610271.
44. E. Chane-Sane, P.-A. Leziart, T. Flayols, O. Stasse, P. Souères, N. Mansard, "CaT: Constraints as terminations for legged locomotion reinforcement learning" in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2024), pp. 13303–13310; 10.1109/IROS58592.2024.10802334.
45. J. Lee, L. Schroth, V. Klemm, M. Bjelonic, A. Reske, M. Hutter, "Exploring constrained reinforcement learning algorithms for quadrupedal locomotion" in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE, 2024), pp. 11132–11138; 10.1109/IROS58592.2024.10801341.
46. Y. Kim, H. Oh, J. Lee, J. Choi, G. Ji, M. Jung, D. Youm, J. Hwangbo, Not only rewards but also constraints: Applications on legged robot locomotion. *IEEE Trans. Robot.* **40**, 2984–3003 (2024).
47. J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, Proximal policy optimization algorithms. arXiv:1707.06347 [cs.LG] (2017).
48. Qualisys; <https://www.qualisys.com/>.
49. J. He, C. Zhang, F. Jenelten, R. Grandia, M. Bächer, M. Hutter, Data for "Attention-based map encoding for learning generalized legged locomotion," Zenodo (2024); 10.5281/zenodo.14499786.

**Acknowledgments:** We thank N. Rudin and V. Koltun for helpful discussion. **Funding:** This work was funded in part by the NCCR Automation, EU Project 10112321 and 852044, the

ETH Mobility Initiative, Fourier Intelligence, Apple Inc., and armasuisse Wissenschaft und Technologie W+T. Any views, opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and should not be interpreted as reflecting the views, policies, or positions, either expressed or implied, of Apple Inc. **Author contributions:** J.H. and C.Z. formulated the main ideas, trained the controllers, and conducted most of the experiments. F.J. and R.G. helped with experiments and provided insights into system development. M.B. and M.H. acquired the funding for the project. All authors helped to write, improve, and refine the paper. **Competing interests:** The authors declare that they have

a patent application pending. **Data and materials availability:** All data needed to support the conclusions of this manuscript are included in the main text or Supplementary Materials. Data to reproduce our plots are available in (49).

Submitted 16 December 2024

Accepted 29 July 2025

Published 27 August 2025

10.1126/scirobotics.adv3604

## Attention-based map encoding for learning generalized legged locomotion

Junzhe He, Chong Zhang, Fabian Jenelten, Ruben Grandia, Moritz Bächer, and Marco Hutter

*Sci. Robot.* **10** (105), eadv3604. DOI: 10.1126/scirobotics.adv3604

### View the article online

<https://www.science.org/doi/10.1126/scirobotics.adv3604>

### Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

---

*Science Robotics* (ISSN 2470-9476) is published by the American Association for the Advancement of Science, 1200 New York Avenue NW, Washington, DC 20005. The title *Science Robotics* is a registered trademark of AAAS.

Copyright © 2025 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works