

Visual-Locomotion: Learning to Walk on Complex Terrains with Vision

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Vision is one of the most important perception modalities for legged
2 robots to safely and efficiently navigate uneven terrains, such as stairs and stepping
3 stones. However, training robots to effectively understand high-dimensional visual
4 input for locomotion is a challenging problem. In this work, we propose
5 a framework to train a vision-based locomotion controller which enables a
6 quadrupedal robot to traverse uneven environments. The key idea is to introduce a
7 hierarchical structure with a high-level vision policy and a low-level motion controller.
8 The high-level vision policy takes as inputs the perceived vision signals as
9 well as robot states and outputs the desired footholds and base movement of the
10 robot. These are then realized by the low level motion controller composed of a
11 position controller for swing legs and a MPC-based torque controller for stance
12 legs. We train the vision policy using Deep Reinforcement Learning and demon-
13 strate our approach on a variety of uneven environments such as randomly placed
14 stepping stones, quincuncial piles, stairs, and moving platforms. We also vali-
15 date our method on a real robot to walk over a series of gaps and climbing up a
16 platform.

17 **Keywords:** Legged Robot, Reinforcement Learning, Visual Locomotion

18 1 Introduction



Figure 1: Our approach trains a visual-locomotion policy for the Laikago robot in simulation (left) and transfers to the real world (right). We use two depth sensors on the robot, as marked in red.

19 Reproducing natural vision-based locomotion skills seen in animals on artificial creatures such as
20 legged robots has long been at the forefront of robotics research. Progress towards developing
21 robust visual locomotion controllers not only deepens our understanding of how humans and animals
22 perceive the environment and control our limbs, but also enables us to build autonomous machines
23 that can reliably traverse real-world environments. To tackle this problem, most existing methods
24 adopt a three-stage pipeline [1, 2, 3, 4, 5, 6, 7]: perception, motion planning, and control. In the
25 perception stage, raw sensor data such as RGB image, depth image, and/or LiDAR point clouds are
26 carefully fused with proprioceptive streams such as IMU, motor angles, and wheel odometry. For
27 mobile robots, a SLAM (simultaneous localization and mapping) or equivalent approach is often
28 required to produce a elevation terrain map centered around the robot [8]. The generated elevation
29 map is then fed to the downstream motion planning modules to select paths, motion style, and foot

30 placements (on legged robots). Finally, the planned robot pose or joint angles are tracked by a low-
31 level motion controller: model predictive control (MPC) [9] or whole body control methods (WBC)
32 [4] are popular choices for unstable, highly dynamical platforms such as legged robots.

33 While there have been great results from previous works that adopt the three-stage control pipeline,
34 it often leads to a system that is overly complex and requires significant manual effort to develop.
35 For instance, SLAM algorithms require careful parameter tuning to achieve a balance between lat-
36 ency and accuracy for mobile robots [10]. On the other hand, recent developments in reinforcement
37 learning (RL) based methods open an alternative path towards creating vision-based locomotion con-
38 trollers without relying on terrain reconstruction or extensive prior knowledge for foothold selection.
39 Researchers have proposed learning-based algorithms that teach robot arms to retrieve objects from
40 cluttered environments [11] and teach drones to avoid obstacles [12] using vision input directly.
41 They demonstrate the great potential in an end-to-end learning approach to achieve a low-latency
42 control pipeline and reduce the prior knowledge required to design a working controller.

43 Inspired by this recent progress, we propose a novel control architecture that enables legged robots
44 to successfully solve various visual-locomotion tasks. Specifically, we adopt the philosophy of
45 end-to-end learning and merge the perception and motion planning modules using a neural network
46 (NN). Our contribution is a learnable hierarchical system that contains two individual layers: a high-
47 level vision policy and a low-level motion controller. The high-level vision policy takes two depth
48 images and outputs the desired pose of the robot’s base and foothold placements for all swing legs,
49 thereby eliminating the need for a complex SLAM algorithm. The low-level locomotion controller
50 takes inputs from the high-level vision policy, and outputs the target motor positions and torques
51 to achieve the desired states. This hierarchical approach allows us to achieve dynamic locomotion
52 in challenging simulated environments including randomly placed stepping stones, staircases, and
53 moving platforms. We also demonstrate zero-shot sim-to-real transfer of visual locomotion policies
54 on the real hardware for walking over stepping stones and climbing up a platform.

55 2 Related Work

56 Recently, visual-locomotion researchers have developed various promising approaches that tackle
57 the problem using on-board sensors [1, 2, 3, 4, 5, 6, 7, 13, 14, 15, 16]. The main idea in these works
58 is to perform explicit terrain shape reconstruction (i.e. local SLAM). For example, Fankhauser and
59 Hutter [13] developed a “Grid Map” stack to construct local elevation maps around the robot and
60 Kim et al. [14] utilized an Intel RealSense D435 depth sensor and a T265 tracking camera to obtain
61 an elevation map of the robot’s surroundings. However, in our work we do not employ any explicit
62 mapping mechanism; Instead, our learning based visual policy consumes raw depth images from
63 onboard cameras for decision making. Because we eliminate the highly complex mapping process
64 which requires special expertise to tune, our system has a simplified and low latency data pipeline
65 with less chance for manual error accumulation.

66 Given a representation of the environment (terrain heights, gradients, etc) provided by a perception
67 module, the controller needs to plan a sequence of leg movements that can guide the robot through
68 the environment and track the pre-planned leg motions as closely as possible. Unlike our approach
69 where perception and motion planning are subsumed into a single neural network, previous works
70 utilize heuristic-based approaches or local optimization to plan adequate motions [5, 3, 16, 2, 7,
71 17, 18, 19]. For instance, Jenelten et al. developed a metric for scoring each point on an elevation
72 map based on its roughness, closeness to an edge, as well as slope degrees. They then employed an
73 online batch optimization to choose the best landing positions for the feet. In contrast, our learning
74 system does not need any of these additional metrics to produce successful feet targets, and greatly
75 simplifies the system architecture.

76 Foothold optimization is commonly performed for the current locomotion phase (i.e. within a single
77 swing step), however, recent works have proposed multi-step contact sequence optimization [3, 16]
78 to account for upcoming changes in the terrain. In our work, we train the perception and planning
79 module together using deep reinforcement learning, which endows the policy with implicit long-
80 horizon planning capabilities. Although neural networks were also employed in prior works [2, 7]
81 to decide optimal footholds, our system differs in that it also decides the desired body pose and
82 speed. By simultaneous adjusting body pose and footholds, our system can learn challenging visual

83 locomotion tasks such as the stepping stones, which is deemed as difficult to solve with constant
 84 body velocities [5],

85 In addition to methods that leverages the vision input for controlling the robot, recent works have
 86 also shown remarkable results for training robots to traverse some uneven terrain without relying
 87 on visual perception [20, 21]. By carefully designing the reward function and training procedure,
 88 these methods produces policies that can robustly handle moderate uneven terrains by adjusting
 89 leg movements. Despite the impressive results demonstrated in these works, vision inputs are still
 90 beneficial for efficiently and safely traversing general unstructured terrains. For some tasks, such as
 91 walking over gaps, visual guidance is required for any practical solutions.

92 3 Methodology

93 3.1 Overview

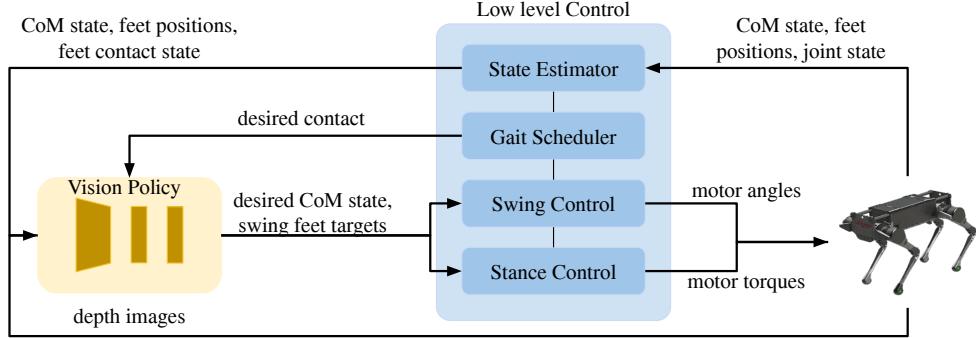


Figure 2: Overview of the visual-locomotion control architecture.

94 In this work, we adopt a hierarchical architecture for vision-locomotion control, as illustrated in
 95 Figure 2. The controller consists of two components: a high-level vision policy and a low-level
 96 motion controller. The high-level vision policy uses raw depth images from the onboard cameras
 97 and robot state to decide the desired center of mass (CoM) pose (and speed) and landing positions
 98 of the swing feet (Section 3.2). Then the low-level motion controller tracks the desired CoM state
 99 by combining a position-based swing leg controller and an model predictive control (MPC) based
 100 stance leg controller (Section 3.3). The vision policy runs at 20 Hz, while the motion controller
 101 runs at 250 Hz. We train our policy entirely in a physics simulation environment and transfer the
 102 trained policy to a real robot platform using techniques described in Section 3.4. Unless otherwise
 103 specified, all quantities are defined in the yaw aligned inertial frame (Appendix A.1).

104 3.2 Learning-Based Vision Policy

105 We formulate the visual-locomotion task as a Partially-observable Markov Decision Process
 106 (POMDP), $(O, S, A, R, P, p_0, \gamma)$, where O is the observation space, S is the state space, A is the
 107 action space, $R : S \times A \mapsto \mathbb{R}$ is the reward function, $P : S \times A \mapsto S$ is the dynamics equation, p_0
 108 is the initial state distribution and γ is the discount factor. Our goal is to find a policy $\pi : O \mapsto A$
 109 that maximizes the expected accumulated reward: $J(\pi) = \mathbb{E}_{\tau=(s_0, a_0, \dots, s_T)} \sum_{t=0}^T \gamma^t r(s_t, a_t)$.

110 **Observation and action spaces.** We design the observation space in our experiments as $O =$
 111 $(\mathcal{I}_{1,2}, \mathbf{q}, \mathbf{r}_{1\dots 4}, c_{1\dots 4}, \phi_{1\dots 4}, \mathbf{a}_{prev})$, i.e the two depth images $\mathcal{I}_{1,2}$ from depth sensors shown in Figure 1 (to cover terrains near both front and rear legs), the CoM pose \mathbf{q} (defined in Appendix A.5),
 112 the robot's feet positions $\mathbf{r}_{1\dots 4}$, the current feet contact states $c_{1\dots 4}$ (one if in contact and zero otherwise), the phase ϕ , of each leg in its respective gait cycle (Appendix A.3), and the previous action.
 113 The action space is $(\mathbf{q}^d, \dot{\mathbf{q}}^d, \mathbf{r}_{1\dots 4}^d)$, i.e. the desired base pose, velocity, and i th swing foot's target
 114 landing position (r_{xi}, r_{yi}) . The z component of the landing location is directly inferred from the
 115 depth readings: we compute the point cloud from the depth image near the currently predicted (x, y)
 116 coordinate and take the average height of four nearest points. By inferring the z component of the
 117 118

119 landing location from depth images, we eliminate invalid foothold targets (e.g. inside an obstacle or
 120 high in the air) and thus create a better action space for the policy to explore within. We find this to
 121 be critical for the policy to traverse highly uneven terrains such as stairs.

122 **Reward function.** We design a reward function:

$$R(\mathbf{s}, \mathbf{a}) = \text{clip}(\dot{p}_x, -\dot{p}_x^{\max}, \dot{p}_x^{\max}) - w_1(|p_y| + |\dot{p}_y|) - w_2|\Psi|, \quad (1)$$

123 where \dot{p}_x is the CoM velocity in the forward direction, p_y the CoM displacement in the lateral
 124 direction, and Ψ the base yaw angle. The first term rewards the robot to move forward with a
 125 maximum speed controlled by \dot{p}_x^{\max} , the second term penalizes the robot from moving sideways, and
 126 the last term encourages the robot to walk straightly; w_1, w_2 modulates the importance of different
 127 terms. In our experiments, we chose $\dot{p}_x^{\max} = 0.375$ m/s, $w_1 = 1.25$, and $w_2 = 0.125$.

128 **Early termination.** A training episode is terminated if: 1) the robot loses balance (CoM height
 129 p_z below 0.15 m, pitch $|\Theta| > 1$ rad, or roll $|\Phi| > 0.3$ rad in our experiments), 2) the robot steps
 130 within 0.02 m of the boundary of the stepping stones or stairs, or 3) the robot reaches an invalid joint
 131 configuration, e.g. knee bending backwards.

132 3.3 Motion Controller

133 Given the desired landing positions of the swing legs and the target base poses from high level, the
 134 motion controller achieves these goals by computing appropriate motor position and torque com-
 135 mands. Our motion controller contains separate modules for swing and stance leg control. The
 136 swing leg controller computes the feet positions by interpolating a time based curve $\alpha(t)$ between
 137 the swing start and target landing positions. Instantaneous feet positions are converted to desired
 138 joint angles through inverse kinematics (IK) and tracked using proportional-derivative (PD) con-
 139 trol. The stance leg controller, on the other hand, achieves the desired base poses by computing
 140 sequences of contact forces between the feet and the ground. By approximating the robot dynamics
 141 using Centroidal Dynamics Model (CDM), we formulate this problem as a convex model predictive
 142 control (MPC) problem similar to [22]. The desired position and velocity of the robot base from
 143 the policy is then used as the tracking target in the MPC formulation. The optimized contact forces
 144 are mapped to stance leg joint torques using Jacobian Transpose. More details regarding the motion
 145 controller can be found in Appendix A.

146 3.4 Sim-to-real Transfer

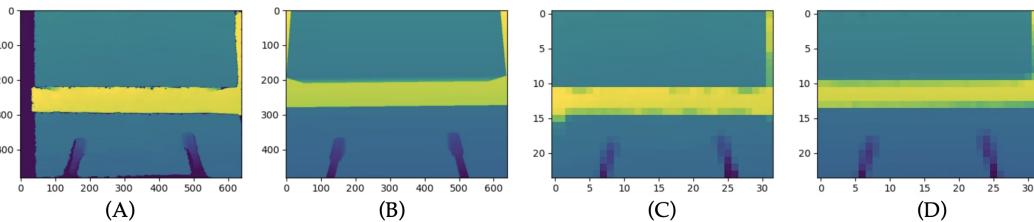


Figure 3: A post-processing technique to reduce the sim-to-real gap in perception. (A) An original depth map obtained from the D435 camera with an invalid band on the left of the image and noises near edges. (B) A simulated depth map. (C) The filtered and down-sampled D435 depth image. (D) The filtered and down-sampled simulation depth image.

147 Sim-to-real transfer is a persisting problem in robotics, and is especially hard for visual locomotion
 148 problems because of distribution shift in image space and discrepancies in robot dynamics.

149 Simulated images can be significantly different from those obtained from the real-world sensors,
 150 Figure 3, A and B. To bridge the gap, we adopt a post-processing procedure that maps both simulated
 151 and real depth images to a similar domain. As shown in Figure 3, C and D, we first add random
 152 Gaussian noise and randomly paint pixels black in the simulated depth image. For pixels that are
 153 along an edge in the depth image identified by a Canny edge detector, they have a higher probability
 154 of being dropped. This is to mimic the noise around the object edges in a stereo depth camera. We

155 then apply an in-painting operation for both simulated and real images to fill the missing pixels[23],
156 followed by a down-sampling operation. This transforms the depth images in simulation and real
157 world to a similar distribution.

158 Discrepancies in dynamics make precise foot-placements difficult. For example, if there is a state
159 estimation error, motors are weaker or joint frictions are higher, the real robot will likely step to
160 different positions than anticipated. Even if the same sequence of actions are executed, the average
161 difference of feet landing positions in the world frame can be $2 \sim 4$ cm between sim and real, which
162 is enough to trigger a failure in the step-stone task. To compensate this issue, we adopt the domain
163 randomization technique [24] to train a robust policy with randomized simulation parameters. More
164 details regarding the randomization parameters and ranges can be found in Appendix C. Adding
165 randomization not only allows the policy to be more robust, but also enables the vision policy to
166 observe more diverse states.

167 4 Experiment and Results

168 4.1 Experiment Setup

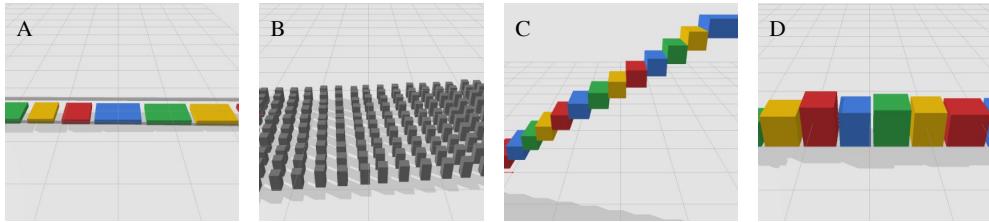


Figure 4: Examples of environments used for visual locomotion tasks: randomly placed stepping stones, quincuncial piles, staircases, and uneven terrains.

169 We evaluate the proposed framework on the Unitree Laikago [25] quadruped robot, which is 0.55 m
170 in length, 0.35 m in width, and about 0.45 m in height when walking. The robot is equipped with 12
171 actuated joints. To collect visual data, we install two depth cameras on the robot: one Intel D435 in
172 the front for a wider field of view and one Intel L515 on the belly with better depth quality in close
173 range (Figure 1). At inference time, we process all depth images as described in Section 3.4, which
174 results in two 32×24 depth images as inputs to the NN policy.

175 We train our hierarchical policies using the PyBullet physics simulator [26]. We use a distributed
176 implementation of augmented random search (ARS) [27], which optimizes the policy parameters
177 using an evolutionary strategy. We use a Multi-layer Perceptron (MLP)-based policy architecture
178 with 54,656 parameters. More training details can be found in Appendix D. For all experiments,
179 we run ARS with a grid search for four hyper-parameters, resulting in 16 trials. We then report the
180 performance of top-3 policies by testing them on 150 randomized environments for each task.

181 For experiments in the real world, we first train policies in simulation with perception noise only
182 and then fine-tune them with the dynamics randomization scheme as described in Section 3.4. We
183 empirically find this helpful for faster training convergence. The fine-tuned policy is deployed on
184 the robot without the need of additional hardware data. Videos of our trained policies in simulation
185 and real-world can be found in the supplementary video.

186 4.2 Simulation Results

187 We first evaluate the capability of our learning system on a variety of challenging visual locomotion
188 tasks, including walking over randomly placed stepping stones, quincuncial piles, uneven terrains,
189 and moving platforms. A subset of the simulation environments are shown in Figure 4. We train
190 separate policies for each environment and show the statistical results in Figure 5. To measure
191 the performance of each policy, we design a metric, *performance ratio*, as $\frac{p_x^T}{p_x^{max}}$, where p_x^T is the
192 distance travelled by the policy in the environment, and p_x^{max} is the maximum distance the policy
193 could reach, e.g. where the terrain ends.

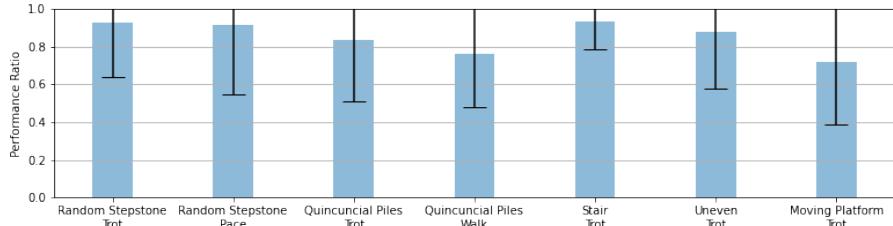


Figure 5: Performance ratio of our policy for different simulated environments.

194 **Randomly placed stepping stones.** The first task for Laikago robot is to walk over a series of
 195 randomly placed stepping stones, Figure 4A. The widths, lengths and gap sizes of stepping stones
 196 are sampled from $[0.55, 0.7]$, $[0.5, 0.8]$, and $[0.07, 0.2]$ meters respectively. The agent fails the task
 197 if the robot steps outside of the stone or into the gaps. Since the stones and their positions are
 198 randomly sampled, this task cannot be accomplished without vision. To successfully achieve this
 199 task, the robot needs to correctly identify the position and size of the stones and also plan for future
 200 footholds. We demonstrate that our method can obtain policies that solves this task using both
 201 trotting and pacing gaits.

202 **Quincuncial piles.** The random stepping stone task evaluates the robot’s ability to identify and
 203 handle terrain changes mainly in the forward direction. A natural extension is to include the lateral
 204 direction into consideration as well. So in the second task we create a grid of uniform quincuncial
 205 piles (i.e. 2 dimensional stepping stones), Figure 4B. Each stone has an area of $0.15 \times 0.15 \text{ m}^2$ with
 206 a standard deviation of 0.015 m in height, and is separated by $[0.13, 0.17] \text{ m}$ from each other in both
 207 x and y directions. At the beginning of each episode, we also randomly rotate entire stone grid by
 208 an angle sampled in $[-0.1, 0.1] \text{ rad}$. Despite significantly more difficult, using our framework we
 209 can obtain policies that can traverse the field using trotting or walking gaits.

210 **Uneven terrains.** So far the tasks focus on evaluating the robot’s ability to avoid undesired regions
 211 on relatively even terrains. To test the policy’s performance in environments with different heights,
 212 we designed an environment where the robot needs to climb up a sequence of stairs, Figure 4C.
 213 The depth of each stair is uniformly sampled in $[0.25, 0.33] \text{ m}$ and the height is in $[0.16, 0.19] \text{ m}$ to
 214 mimic the dimensions of real-world stairs. This task is quite challenging for Laikago robot because
 215 each stair is as tall as the robot’s knee joint. Figure 4D shows another environment we designed
 216 for evaluating the ability of our approach to handle uneven terrains. In this environment, the height
 217 offsets of neighboring stones are uniformly sampled in $[-0.13, 0.2] \text{ m}$, and a gap of $[0.05, 0.1] \text{ m}$ is
 218 added between the stones.

219 **Moving platforms.** One benefit of using camera images as inputs is that, policies can potential
 220 handle moving objects better when compared with methods that relies on terrain reconstruction, be-
 221 cause of latency and stationary assumptions in SLAM. To demonstrate the capability of our learning
 222 system, we take the random stepping stone environment (Figure 4A) and allow each piece to move.
 223 Each platform follows a periodic movement whose magnitude and frequency are randomly sampled
 224 in $[0.10, 0.15] \text{ m}$ and $[0.4, 1.0] \text{ Hz}$, respectively. Also, we randomly pick half of the platforms to
 225 move horizontally and the rest vertically. This task requires the robot to infer both the position and
 226 velocity of the platforms. To facilitate learning, we stack a history of three recent images as input
 227 to the policy for this task. As shown in the accompanying video, our policy learns to identify the
 228 moving objects in the scene and will slow-down and wait for the platforms to reach an ideal location
 229 before striding.

230 4.3 Comparison to baseline methods

231 We compare our proposed method to two baselines: end-to-end training and heuristics-based foot
 232 placement. For the sake of simplicity yet without loss of generality, all comparisons are done in the
 233 uneven terrain environment (Figure 6 Top) composed of stepping stones with large height differences
 234 and moderate gaps, which captures the difficulty of both the stepping stones and staircases. The
 235 results can be found in in Figure 7.

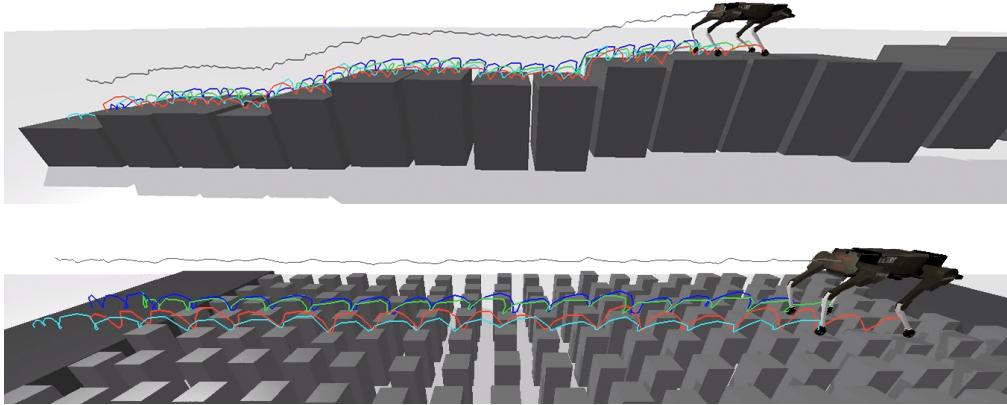


Figure 6: Laikago robot walking over uneven terrain and randomized pillar terrains. The black curve refers to the CoM trajectory, and colored curves represent the feet trajectories.

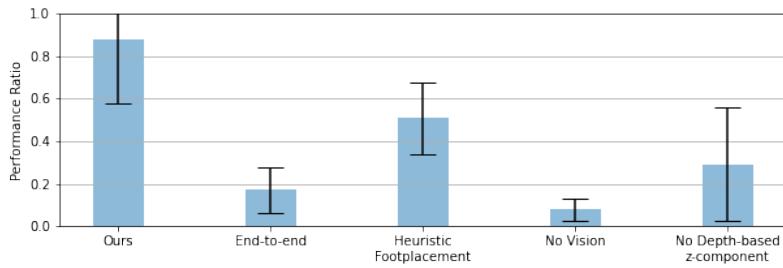


Figure 7: Comparison of performance between our method and the baseline methods.

236 4.3.1 End-to-end training

237 In the first baseline, we train an end-to-end neural network policy that takes the images and the
 238 robot states as input (same as our hierarchical policy) and outputs desired motor angles. We choose
 239 the architecture of Policy Modulating Trajectory Generator (PMTG) [28], which has demonstrated
 240 high-quality and transferrable locomotion policies [20]. Since PMTG generates the trajectory in the
 241 joint space, it was not able to achieve precise foot placement that is needed for many of our testing
 242 environments, as shown in Figure 7.

243 4.3.2 Heuristics-based foot placement

244 In the second baseline, we compare our method to prior approaches that adjust the robot foot placement
 245 by computing a score map for the surrounding terrains [2, 4, 18]. In our baseline implemen-
 246 tation, we define a point on the nearby terrain to be valid if 1) it is within reach of the swing leg,
 247 2) the z-component of the surface normal is larger than 0.9 (i.e. pointing upward), and 3) the stan-
 248 dard deviation of the heights in the nearby region (i.e. roughness) is less than 0.05m. For a fair
 249 comparison, we use the same MPC-based motion controller as in our approach and train a neural
 250 network policy to output the desired base pose and velocity using the same observation space and
 251 learning procedure. As shown in Figure 7, using heuristics alone to determine footplacement got
 252 lower score than our proposed approach. One important reason is that this heuristics does not take
 253 the base movement of the robot into consideration. As a result, it may propose landing positions that
 254 are incompatible with the robot’s CoM speed.

255 4.4 Ablation Study

256 A key hypothesis we make in this work is that, despite that the images have low resolution, they
 257 contain critical information that enables our policy to traverse a large variety of environments. To
 258 validate this hypothesis, we first perform an ablation run by removing the vision component from
 259 the observation space and retraining a vision-less policy. As seen in Figure 7, without vision input,
 260 the policy is not able to accomplish the task, indicating the importance of having vision as input.

261 Another key component in our method is to use the depth input for inferring the z-component of
 262 the foothold location. This creates a more meaningful action space for the policy to explore during
 263 learning. As shown in Figure 7, the performance drops significantly we do not use depth to infer the
 264 z-component of the foot placement.

265 **4.5 Validation on Real Robot**

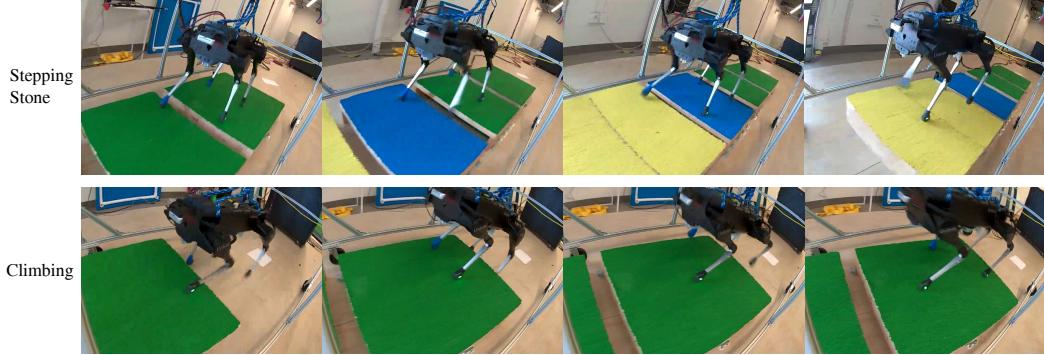


Figure 8: Laikago robot solving two challenging visual-locomotion tasks.

266 We deploy the trained policy to a Laikago robot for two tasks, walking over gaps and climbing onto
 267 stepping stones (see supplementary video). Our real-world setup (Figure 8 Top) consists of four
 268 stepping stones that are separated by three gaps with widths between $[0.12, 0.18]$ m.

269 In the first task, the performance is measured by the number of gaps that are crossed successfully,
 270 summed over legs. For instance, if all four legs cross one gap and the robot is completely on the next
 271 stepping stone, the score would be 4. Therefore, the upper bound of the score is 12: three gaps times
 272 four legs. An experiment episode is terminated if the robot falls, any leg steps into the gaps or the
 273 robot steps outside the stones. As a baseline, the score is 0.9 ± 1.4 for a blind-locomotion policy:
 274 it can rarely clear a single gap. Over eight real-world episodes, our policy is able to consistently
 275 achieve a score of 10.1 ± 2.2 . There are two episodes in which the robot completes the full course
 276 with a score of 12; For seven out of eight episodes, the robot reaches the last stepping stone.

277 In the second task, the robot needs to climb onto a stepping stone, which is 0.18 m above from
 278 the ground (Figure 8 Bottom). Similarly, we also define the performance score as the number of
 279 legs reached the top of the stone (4 is the max). The learned policy reaches a score of 3.7 ± 0.7
 280 over seven experiment episodes. The robot successfully steps onto the stone for six out of the seven
 281 times. In the only failure case, the robot loses balance because its rear foot hits the edge of the stone.

282 **5 Conclusion and Future Work**

283 In this work, we present a hierarchical learning system to tackle challenging visual-locomotion
 284 tasks. By using a high-level learned vision policy that consumes raw camera images, we eliminate
 285 the need to explicitly construct 3D terrain maps, and thus reduce the control latency and architecture
 286 complexities, and more importantly, can handle dynamically moving terrains. The low-level, model
 287 predictive control based motion controller greatly reduces the motion tracking error on the hardware
 288 and thus narrows the sim-to-real gap. We demonstrate that policies trained using our system can
 289 reliably walk over challenging terrains such as stepping stones that require fine visual-motor control.

290 One limitation of this work is that the gaits are manually specified and fixed for each task. In
 291 contrast, animals can dynamically modulate their walking pattern to adapt to changes in the terrain.
 292 In the future, we plan to extend our framework to enable the visual policy to adjust the gait pattern
 293 online. Additionally, animals can traverse complex terrains while paying attention to a small area,
 294 and can reason about their hind limbs that they cannot see. This makes adding memory [29] and
 295 attention [30] to the policy architecture another promising future direction.

296 **References**

- 297 [1] C. Mastalli, M. Focchi, I. Havoutis, A. Radulescu, S. Calinon, J. Buchli, D. G. Caldwell, and
298 C. Semini. Trajectory and foothold optimization using low-dimensional models for rough ter-
299 rain locomotion. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*,
300 pages 1096–1103. IEEE, 2017.
- 301 [2] O. A. V. Magana, V. Barasuol, M. Camurri, L. Franceschi, M. Focchi, M. Pontil, D. G. Cald-
302 well, and C. Semini. Fast and continuous foothold adaptation for dynamic locomotion through
303 CNNs. *IEEE Robotics and Automation Letters*, 4(2):2140–2147, 2019.
- 304 [3] O. Villarreal, V. Barasuol, P. M. Wensing, D. G. Caldwell, and C. Semini. MPC-based con-
305 troller with terrain insight for dynamic legged locomotion. In *2020 IEEE International Con-
306 ference on Robotics and Automation (ICRA)*, pages 2436–2442. IEEE, 2020.
- 307 [4] P. Fankhauser, M. Bjelonic, C. D. Bellicoso, T. Miki, and M. Hutter. Robust rough-terrain
308 locomotion with a quadrupedal robot. In *2018 IEEE International Conference on Robotics
309 and Automation (ICRA)*, pages 5761–5768. IEEE, 2018.
- 310 [5] F. Jenelten, T. Miki, A. E. Vijayan, M. Bjelonic, and M. Hutter. Perceptive locomotion in rough
311 terrain—online foothold optimization. *IEEE Robotics and Automation Letters*, 5(4):5370–5376,
312 2020.
- 313 [6] R. Grandia, A. J. Taylor, A. D. Ames, and M. Hutter. Multi-layered safety for legged robots
314 via control barrier functions and model predictive control. *arXiv preprint arXiv:2011.00032*,
315 2020.
- 316 [7] S. Gangapurwala, M. Geisert, R. Orsolino, M. Fallon, and I. Havoutis. RLOC: Terrain-
317 aware legged locomotion using reinforcement learning and optimal control. *arXiv preprint
318 arXiv:2012.03094*, 2020.
- 319 [8] P. Fankhauser, M. Bloesch, and M. Hutter. Probabilistic terrain mapping for mobile robots
320 with uncertain localization. *IEEE Robotics and Automation Letters (RA-L)*, 3(4):3019–3026,
321 2018. [doi:10.1109/LRA.2018.2849506](https://doi.org/10.1109/LRA.2018.2849506).
- 322 [9] J. Di Carlo, P. M. Wensing, B. Katz, G. Bledt, and S. Kim. Dynamic locomotion in the mit
323 cheetah 3 through convex model-predictive control. In *2018 IEEE/RSJ International Confer-
324 ence on Intelligent Robots and Systems (IROS)*, pages 1–9, 2018. [doi:10.1109/IROS.2018.8594448](https://doi.org/10.1109/IROS.2018.8594448).
- 326 [10] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J.
327 Leonard. Past, present, and future of simultaneous localization and mapping: Toward the
328 robust-perception age. *IEEE Transactions on robotics*, 32(6):1309–1332, 2016.
- 329 [11] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly,
330 M. Kalakrishnan, V. Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-
331 based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- 332 [12] L. He, N. Aouf, J. F. Whidborne, and B. Song. Integrated moment-based lgmd and deep
333 reinforcement learning for uav obstacle avoidance. In *2020 IEEE International Conference on
334 Robotics and Automation (ICRA)*, pages 7491–7497. IEEE, 2020.
- 335 [13] P. Fankhauser and M. Hutter. A Universal Grid Map Library: Implementation and Use
336 Case for Rough Terrain Navigation. In A. Koubaa, editor, *Robot Operating System (ROS)
337 – The Complete Reference (Volume 1)*, chapter 5. Springer, 2016. ISBN 978-3-319-26052-5.
338 [doi:10.1007/978-3-319-26054-9_5](https://doi.org/10.1007/978-3-319-26054-9_5). URL <http://www.springer.com/de/book/9783319260525>.
- 340 [14] D. Kim, D. Carballo, J. Di Carlo, B. Katz, G. Bledt, B. Lim, and S. Kim. Vision aided dy-
341 namic exploration of unstructured terrain with a small-scale quadruped robot. In *2020 IEEE
342 International Conference on Robotics and Automation (ICRA)*, pages 2464–2470. IEEE, 2020.

- 343 [15] C. Mastalli, I. Havoutis, A. W. Winkler, D. G. Caldwell, and C. Semini. On-line and on-board
 344 planning and perception for quadrupedal locomotion. In *2015 IEEE International Conference*
 345 *on Technologies for Practical Robot Applications (TePRA)*, pages 1–7. IEEE, 2015.
- 346 [16] H.-W. Park, P. M. Wensing, S. Kim, et al. Online planning for autonomous running jumps over
 347 obstacles in high-speed quadrupeds. 2015.
- 348 [17] X. Da, Z. Xie, D. Hoeller, B. Boots, A. Anandkumar, Y. Zhu, B. Babich, and A. Garg.
 349 Learning a contact-adaptive controller for robust, efficient legged locomotion. *arXiv preprint*
 350 *arXiv:2009.10019*, 2020.
- 351 [18] Z. Xie, X. Da, B. Babich, A. Garg, and M. van de Panne. Glide: Generalizable quadrupedal
 352 locomotion in diverse environments with a centroidal model. *arXiv preprint arXiv:2104.09771*,
 353 2021.
- 354 [19] M. Xie, A. Escontrela, and F. Dellaert. A factor-graph approach for optimization problems
 355 with dynamics constraints. *arXiv preprint arXiv:2011.06194*, 2020.
- 356 [20] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning quadrupedal locomotion
 357 over challenging terrain. *Science robotics*, 5(47), 2020.
- 358 [21] J. Siekmann, K. Green, J. Warila, A. Fern, and J. Hurst. Blind bipedal stair traversal via sim-
 359 to-real reinforcement learning. *arXiv preprint arXiv:2105.08328*, 2021.
- 360 [22] G. Bledt, M. J. Powell, B. Katz, J. Di Carlo, P. M. Wensing, and S. Kim. MIT Cheetah 3:
 361 Design and control of a robust, dynamic quadruped robot. In *2018 IEEE/RSJ International*
 362 *Conference on Intelligent Robots and Systems (IROS)*, pages 2245–2252. IEEE, 2018.
- 363 [23] M. Bertalmio, A. L. Bertozzi, and G. Sapiro. Navier-stokes, fluid dynamics, and image and
 364 video inpainting. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer*
 365 *Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I. IEEE, 2001.
- 366 [24] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke. Sim-
 367 to-real: Learning agile locomotion for quadruped robots. *arXiv preprint arXiv:1804.10332*,
 368 2018.
- 369 [25] Unitree Robotics. URL <http://www.unitree.cc/>.
- 370 [26] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation in robotics, games
 371 and machine learning, 2017.
- 372 [27] H. Mania, A. Guy, and B. Recht. Simple random search provides a competitive approach to
 373 reinforcement learning. *arXiv preprint arXiv:1803.07055*, 2018.
- 374 [28] A. Iscen, K. Caluwaerts, J. Tan, T. Zhang, E. Coumans, V. Sindhwani, and V. Vanhoucke.
 375 Policies modulating trajectory generators. In *Conference on Robot Learning*, pages 916–926.
 376 PMLR, 2018.
- 377 [29] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–
 378 1780, 1997.
- 379 [30] K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins,
 380 J. Davis, A. Mohiuddin, L. Kaiser, et al. Rethinking attention with performers. *arXiv preprint*
 381 *arXiv:2009.14794*, 2020.
- 382 [31] M. H. Raibert. *Legged robots that balance*. MIT press, 1986.