

A Hierarchical Reinforcement Learning Approach to Control Legged Mobile Manipulators

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Recent years have seen a Cambrian explosion of robotic systems yield-
2 ing ever more capable and affordable systems, with quadrupedal robotic platforms
3 emerging as a commercially-viable base to perform a wide variety of tasks across
4 uneven terrain. Augmenting these with a robotic arm allows the possibility of even
5 more complex interactions. At the same time, there has been a growing body of
6 research into using deep reinforcement learning (DRL) for embodied agent navi-
7 gation and object manipulation, which promises a more sample-efficient, flexible,
8 and robust approach to learning such policies than existing classical methods. Re-
9 cent works have shown a functional approach for learning a joint base and arm
10 policy with DRL but have not yet demonstrated how the result can be used in
11 downstream tasks. In this work, we investigate the problem of learning an ob-
12 ject manipulation and navigation policy for a quadrupedal robot with a mounted
13 robotic arm - specifically, we address the problem of fetching stationary and mov-
14 ing objects autonomously (“playing fetch” with the robot dog). Our method con-
15 sists of (a) a low-level policy that moves the base and arm and (b) a high-level
16 policy that generates the commands for the low-level policy. The low-level policy
17 is jointly learned for both the arm and the base which generates joint torques for
18 directional commands. The high-level policy is task-specific, translates the ball
19 position to directional commands for the low-level policy, and deals with accel-
20 eration/deceleration and stability. We demonstrate that our high-level policy can
21 outperform a tuned Proportional-Derivative (PD) controller.

22 **Keywords:** Reinforcement Learning, Quadruped Robots, Object Manipulation

23 1 Introduction

24 Mobile ground-based robots have been used to solve many tasks such as delivery, inspection, and
25 exploration [1]. More specifically, quadrupedal legged robots have the agility to navigate uneven
26 terrain and traverse obstacles commonly found in indoor settings such as stairs. At the same time,
27 robotic arms are very useful for manipulation tasks that require grasping and picking up objects,
28 but are most often stationary. Merging both of these technologies could yield very effective robotic
29 assistants, capable of accomplishing useful tasks in uncontrolled real-world environment. However,
30 this combination of quadruped base with robotic arm poses a difficult control problem.

31 Many classic control and Reinforcement Learning (RL) methods have been proposed for each of
32 these domains, quadrupedal locomotion [2, 3] and robotic arm manipulation [4, 5]. Far fewer meth-
33 ods have been proposed for the control of quadrupedal robots with robotic arms. Such methods
34 must take into account the interactions between the two robotic systems, especially when dealing
35 with moving objects or rapidly changing commands. As an example, [6] propose to add a correction
36 term to their state estimation pipeline to compensate for the fact that the arm is moving and affecting
37 the locomotion controller behavior. In this work, we target the problem of “playing fetch” with a

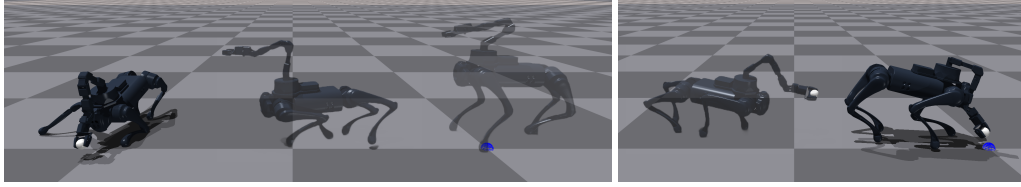


Figure 1: **Example trajectory.** Image sequence from our method of a robot dog running towards a thrown ball, picking it up before the ball comes to rest, and returning to the origin to drop off the ball.

38 quadruped robot because we believe it exemplifies a hard control task involving both locomotion
 39 and manipulation of dynamic objects.

40 We propose a hierarchical system with a task-independent low-level policy that learns to convert
 41 directional commands for the base and position offsets for the gripper to joint torques, and a high-
 42 level policy that translates object goals to commands for the low-level policy, balances speed and
 43 stability and learns task-specific knowledge.

44 Figure 1 shows an example of our method in action. Figure 2 shows the schematics of the proposed
 45 architecture, with the high-level and low-level control policy.

46 **Our contributions are as follows.** (a) A Hierarchical Reinforcement Learning (HRL) framework
 47 for jointly training a movable base and a robot arm to solve arbitrary control tasks. (b) Experi-
 48 ments demonstrating that the method can be applied to solve the problem of playing fetch with a
 49 quadrupedal robot. (c) Comparisons illustrating how our high-level policy outperforms a PD con-
 50 troller.

51 All experiments were conducted on the Unitree Go1 quadruped with the Interbotix ReactorX 150
 52 arm in simulation but with the manufacturer’s URDF robot models and realistic (manufacturer)
 53 dynamics. If this works is accepted, we will present a demo on a real robot at the workshop. Videos
 54 of the robot playing fetch can be found at our project website:

55 <https://robotdogfetch.github.io>

56 2 Related Work

57 **Control for Legged Mobile Manipulators.** Broadly speaking, control methods for mobile manipu-
 58 lators can be placed in two categories: Model Predictive Control (MPC) and learning-based control.
 59 MPC can produce a simple control policy for complex systems and provides generic consideration
 60 of constraints and complex control goals [6, 7]. However, MPC requires a faithful dynamical model
 61 based on forces and torques and can suffer from long re-planning times in dynamic scenes.

62 In this work, our goal is to provide a policy that can solve legged manipulation tasks with increased
 63 robustness and involving dynamic objects. However, manipulation so far has required more precise
 64 control, combining MPC and RL methods [8]. By contrast, our method does not require an exact
 65 specification of the dynamics of the robot to work. Only very recently, Fu et al. [9] showed that a
 66 joint policy can be learned for the base and arm. This work, however, leaves adding task-specific
 67 knowledge as an exercise to the reader. In our work, we address this explicitly through the high-level
 68 policy that gives the system the autonomy to solve a task without remote control.

69 **Hierarchical Reinforcement Learning.** Hierarchical Reinforcement Learning frameworks [10, 11,
 70 12] typically extend the usual notion of action to include options, a closed-loop policy for taking
 71 a sequence of actions over a period of time. The concept of an option is quite broad and includes
 72 turning around, picking up an object, and driving a commute, as well as lower-level actions such as
 73 joint torques. Many such methods train both low and high level policies jointly [13, 14, 15]. The
 74 high-level policy operates at a slower frequency than low-level policy options, which may have a
 75 fixed or learned duration. For example, [16] used a gradient-based approach to model the termination

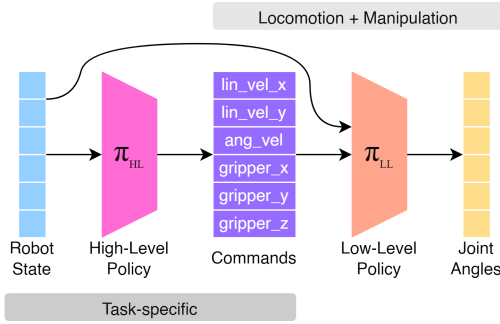


Figure 2: **Abstract policy network.** The high-level policy module takes as input a state vector that describes the robot and world state and produces a command vector. This command vector is then fed into the low-level policy to produce the robot’s actuation. The high-level policy and low-level operate at different frequencies, with the low-level policy having to produce joint angles for multiple timesteps to achieve a goal specified by the command vector.

76 of semi-Markov options by a logistic distribution on a cumulative measure of the features observed
 77 during the execution of that option. In this work, low-level policy options correspond to motor
 78 actuations which can benefit from very high-frequency control to manage high aleatoric uncertainty
 79 in real-world environments, and as a result we fix the duration of low-level policy options.

80 3 Method

81 Our method uses a 2-layer hierarchical policy (see Fig.2), where the low-level policy learns loco-
 82 motion and relative end effector control with respect to the base (based on external commands), and
 83 the high-level policy learns to generate commands for the low-level policy for a given task.

84 **Background.** We use Nvidia’s Isaac Gym simulator to train thousands of quadrupeds in parallel on
 85 a GPU to follow commands [17]. In their work, the agent and in our work, the low-level policy is
 86 trained via Proximal Policy Optimization (PPO) [18] to obey linear and angular velocity commands.
 87 Each command is a 3-dimensional vector, representing the desired linear velocities along the x and y
 88 axes (both of which are parallel to the ground) and the yaw rate (rotation speed around the “up” axis).
 89 At each timestep, the policy outputs a 12-dimensional vector corresponding to joint angles for the
 90 joints of the quadruped (4x hip, upper leg, lower leg). Similar to Rudin et al. [19], we use a weighted
 91 combination of 10 reward terms to teach the robot dog to walk in a physically plausible way, e.g.,
 92 by rewarding the feet to not drag along the floor, or to exert low torques. The most important reward
 93 terms are $r_{tracking_lin_vel}$, which incentivizes the agent to match the robot body’s linear velocity to
 94 the first two command terms and $r_{tracking_ang_vel}$, which does the same for the yaw rate.

95 **Task & Robot.** Our experiments are performed on a simulated Unitree Go1 robot with a manipulator
 96 attached to the top of the robot body (see Fig.1). The mounting location was chosen to optimize
 97 balance while maximizing forward reach. We use the Interbotix ReactorX 150 robot arm due to its
 98 light weight. To train the arm as part of the low-level policy, we added to the observation space
 99 the positions and velocities of all joints, to the action space the 5 degrees of freedom, and to the
 100 command vector 3 elements corresponding to the desired position of the end effector with respect to
 101 the base. For our fetch task, we use a rubber ball and a simple grasping model that closes the gripper
 102 when the ball is in the center of the gripper.

103 **High-level Policy.** Our high-level policy also uses a 3-layer MLP architecture for both actor and
 104 critic networks with [256, 256, 256] hidden units and ELU non-linearities in between. The observa-
 105 tions are task-specific and consist of linear/angular velocity of the base, projected gravity, distance
 106 between end effector and goal, the current relative position of the end effector, the direction of the
 107 goal, and the current heading of the base. The action space consists of the 6 commands comprising
 108 the goal conditioning of the low-level policy.

109 **Reward Terms.** The high-level policy follows the same PPO training procedure as the low-level
 110 policy but with different reward terms. We have reused the existing reward terms and coefficients for
 111 action rate (discouraging large changes between timesteps), and large action values (discouraging

	Success (%)	Ball Pickup (%)	Ball Pickup Time (s)	Goal Time (s)
PD Controller & π_{LL}	0 %	3 %	2.7±0.4	n/a
No Random Command	8 %	18 %	1.5±1.1	0.8±0.2
Ours	97 %	99 %	1.6±1.6	1.0±0.7

Table 1: **Results.** We show the results of our method compared to an ablation (“no command frequency randomization”) and a PD controller. We list the success rate, which includes ball pickup and return within 10s, ball pickup rate, average time to pick up the ball, normalized by distance of pickup, and average time to return the ball to the goal (also normalized by ball distance).

112 extreme joint positions), and we added the terms $r_{goal_distance}$ (Euclidean distance between gripper and goal), $r_{goal_distance_eps}$ (bounded exponential reward for minimizing the distance between gripper and goal), and r_{facing_goal} (encouraging facing towards the goal) - implementation details can be found in the appendix (App.A.1). We also note that the last reward term is both cosmetic and practical for this task as facing the ball reduces the risk of accidentally kicking it further away.

117 4 Experiments

118 **Experimental Setup.** We train our low-level policy for 1,500 timesteps with 4096 parallel environments with a fixed command sampling rate of $\frac{1}{10}$ Hz. We fine-tune the same low-level policy for 119 another 1,500 steps with a variable command sampling rate uniformly sampled from $[\frac{1}{4}, \frac{1}{16}]$ Hz. We 120 found that this domain randomization leads to better generalization. We compare the fine-tuned policy to one that was only trained on the default command resampling time (marked as “No Random 121 Command”). The high-level policy was trained for 3,000 timesteps, corresponding to approximately 122 2 hours of wall clock time on a single Nvidia Geforce RTX 2080 Ti GPU, which makes this approach 123 suitable to be retrained to a wide range of tasks. To illustrate the necessity of the high-level policy 124 in communicating the task goal to the low-level policy, we also compare to the setting where we 125 replace the high-level policy with the PD controller. The PD controller was manually tuned to strike 126 a balance between stability and speed and to face the goal in order to prevent accidental kicks. The 127 PD controller minimizes the distance between the end effector and the target ball. The target ball 128 is thrown from a position close to the robot in a random direction and the robot has to retrieve the 129 ball and return it to the origin point at $(0, 0, 0)$ (marked in blue in Fig.1) within 10s. We repeat this 130 experiment with 100 random throws for each approach. 131 132

133 **Results.** The results are listed in Tab.1. We can see that our method solves the task with high 134 efficiency. The few failure cases were due to sharp turns when the robot was moving at a high 135 speed. The results also highlight the necessity of randomizing the control frequency during training, 136 since the ball is a moving target and commands change very quickly, especially in the beginning of 137 an episode. When investigating the runs that use a PD controller, we noticed an interesting detail: 138 The PD controller is often able to approach the ball and lower the gripper to a position right above 139 but it is not able to bend the body of the robot towards the goal. And indeed, there is no way in which 140 the high-level policy can communicate this directly to the low-level policy. The goal command for 141 the gripper is relative to the body of the robot, i.e., moves with the base. In other words, the high- 142 level policy found a way to “trick” the low-level policy into lowering its body when approaching the 143 goal.

144 **Conclusion.** We presented a hierarchical control approach for combining quadruped locomotion 145 and robot arm manipulation and we demonstrated its efficacy on a difficult rolling ball pickup task. 146 There are still several avenues for improvement of this method that we would like to explore. So 147 far, we have only trained the method on flat ground but going forward, we want to investigate how 148 the method handles bumpy terrain and stairs. We are also excited to train the high-level policy on 149 different tasks with the same low-level policy to analyze how well our method generalizes to new 150 problems and environments.

References

- [1] F. Rubio, F. Valero, and C. Llopis-Albert. A review of mobile robots: Concepts, methods, theoretical framework, and applications. *International Journal of Advanced Robotic Systems*, 16(2):1729881419839596, 2019.
- [2] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning quadrupedal locomotion over challenging terrain. *Science robotics*, 5(47):eabc5986, 2020.
- [3] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter. Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*, 7(62):eabk2822, 2022.
- [4] S. Gu, E. Holly, T. Lillicrap, and S. Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 3389–3396. IEEE, 2017.
- [5] J. Mahler, M. Matl, V. Satish, M. Danielczuk, B. DeRose, S. McKinley, and K. Goldberg. Learning ambidextrous robot grasping policies. *Science Robotics*, 4(26):eaa4984, 2019.
- [6] S. Zimmermann, R. Poranne, and S. Coros. Go fetch! - dynamic grasps using boston dynamics spot with external robotic arm. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4488–4494, 2021. doi:10.1109/ICRA48506.2021.9561835.
- [7] J.-P. Sleiman, F. Farshidian, M. V. Minniti, and M. Hutter. A unified mpc framework for whole-body dynamic locomotion and manipulation. *IEEE Robotics and Automation Letters*, 6(3):4688–4695, 2021.
- [8] Y. Ma, F. Farshidian, T. Miki, J. Lee, and M. Hutter. Combining learning-based locomotion policy with model-based manipulation for legged mobile manipulators. *CoRR*, abs/2201.03871, 2022. URL <https://arxiv.org/abs/2201.03871>.
- [9] Z. Fu, X. Cheng, and D. Pathak. Deep whole-body control: Learning a unified policy for manipulation and locomotion. *arXiv preprint arXiv:2210.10044*, 2022.
- [10] R. Parr and S. Russell. Reinforcement learning with hierarchies of machines. In *Proceedings of the 1997 Conference on Advances in Neural Information Processing Systems 10, NIPS '97*, page 1043–1049, Cambridge, MA, USA, 1998. MIT Press. ISBN 0262100762.
- [11] T. G. Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *CoRR*, cs.LG/9905014, 1999. URL <https://arxiv.org/abs/cs/9905014>.
- [12] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999. ISSN 0004-3702. doi:[https://doi.org/10.1016/S0004-3702\(99\)00052-1](https://doi.org/10.1016/S0004-3702(99)00052-1). URL <https://www.sciencedirect.com/science/article/pii/S0004370299000521>.
- [13] P. Bacon, J. Harb, and D. Precup. The option-critic architecture. *CoRR*, abs/1609.05140, 2016. URL <http://arxiv.org/abs/1609.05140>.
- [14] S. Yu, S. Rammohan, K. Zheng, and G. Konidaris. Hierarchical reinforcement learning of locomotion policies in response to approaching objects: A preliminary study, 2022. URL <https://arxiv.org/abs/2203.10616>.
- [15] C. Li, F. Xia, R. Martin-Martin, and S. Savarese. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators, 2019. URL <https://arxiv.org/abs/1910.11432>.
- [16] G. Comanici and D. Precup. Optimal policy switching algorithms for reinforcement learning. volume 2, pages 709–714, 01 2010. doi:10.1145/1838206.1838300.

- 195 [17] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox. Gpu-accelerated
196 robotic simulation for distributed reinforcement learning. In *Conference on Robot Learning*,
197 pages 270–282. PMLR, 2018.
- 198 [18] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
199 algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- 200 [19] N. Rudin, D. Hoeller, P. Reist, and M. Hutter. Learning to walk in minutes using massively
201 parallel deep reinforcement learning. *CoRR*, abs/2109.11978, 2021. URL [https://arxiv.](https://arxiv.org/abs/2109.11978)
202 [org/abs/2109.11978](https://arxiv.org/abs/2109.11978).

203 **A Appendix**

204 **A.1 Implementation Details**

205 The reward terms for the high-level policy are as follows:

206 $r_{goal_distance}$ (normalized Euclidean distance between gripper and goal) is calculated as

$$r_{goal_distance} = \frac{d_{MAX} - d_{curr}}{d_{MAX}},$$

207 where

208 d_{MAX} .. maximal distance a goal can spawn from the robot’s base.

209 d_{curr} .. current distance between end effector and goal.

210 $r_{goal_distance_eps}$ (bounded exponential reward for minimizing the distance between gripper and
211 goal) is calculated as

$$r_{goal_distance_eps} = e^{\frac{-d_{curr}}{\sigma}},$$

212 where

213 σ .. hyperparameter, set to 0.1 in our experiments.

214 r_{facing_goal} (encouraging facing towards the goal) is calculated as

$$r_{facing_goal} = \|\vec{fwd} - \vec{goal}\|,$$

215 where

216 \vec{fwd} .. unit vector representing the forward direction of the robot body, parallel to the ground plane.

217 \vec{goal} .. unit vector representing the direction towards the goal, also ground-parallel.

218

219 **A.2 Hyperparameters**

Parameter	Value
σ	0.1
coefficient for $r_{goal_distance}$	1.0
coefficient for $r_{goal_distance_eps}$	1.5
coefficient for r_{facing_goal}	0.5