

Demonstrating a Walk in the Park: Learning to Walk in 20 Minutes With Model-Free Reinforcement Learning

Laura Smith^{*1}, Ilya Kostrikov^{*1}, Sergey Levine¹

^{*}Equal contribution ¹Berkeley AI Research, UC Berkeley

{smithlaura, kostrikov}@berkeley.edu, svlevine@eecs.berkeley.edu



Fig. 1: We demonstrate that deep reinforcement learning can be used to efficiently train a quadruped robot directly on various real world terrains, e.g., flat ground (blue); soft, irregular mulch (green); grass (red); and a hiking trail (yellow), acquiring effective gaits within 20 minutes of training.

Abstract—Deep reinforcement learning is a promising approach to learning policies in unstructured environments. Due to its sample inefficiency, though, deep RL applications have primarily focused on simulated environments. In this work, we demonstrate that the recent advancements in machine learning algorithms and libraries combined with careful MDP formulation lead to learning quadruped locomotion in only 20 minutes in the real world. We evaluate our approach on several indoor and outdoor terrains that are known to be challenging for classical, model-based controllers and observe that the robot consistently learns a walking gait on all of these terrains. Finally, we evaluate our design decisions in a simulated environment. We provide videos of all real-world training and code to reproduce our results on our website: <https://sites.google.com/berkeley.edu/walk-in-the-park>

I. INTRODUCTION

Agile, robust, and capable robotic skills require careful controller design and validation to work reliably in the real world. Reinforcement learning offers a promising alternative, acquiring effective control strategies directly through interaction with the real system, potentially right in the environment in which the robot will be situated. Although large-scale robotic RL experiments in the real world have been described in a number of prior works [1–5], many others have sought to sidestep the need for real-world training over concerns about sample efficiency, often going to great lengths to design sophisticated simulation systems [6–10] and transfer learning methods [11–20]. For example, learning to solve a Rubik’s cube with a robotic hand required 13 thousand years’ worth of experience [21], which amounted to several months of wall-clock time using distributed training in simulation. In robotic locomotion, Rudin et al. [22] also utilize distributed training to collect 80 hours’ worth of simulated experience to train an ANYmal robot to walk in 20 minutes. Using the A1 quadrupedal robot as we do, Kumar et al. [16] use 1.2 billion samples to train a robust controller, corresponding to

roughly 4.5 months’ worth of cumulative experience—without accounting for the often significant physical overhead required to train in the real world—that can be acquired in 1 day using an appropriate simulator. In this paper, we focus specifically on the task of robotic quadrupedal locomotion and ask: just how efficiently can we implement fully model-free deep RL algorithms?

Perhaps surprisingly, we find that with careful design decisions in the task setup and algorithm implementation, it is possible for a quadrupedal robot to learn to walk from scratch in under 20 minutes across a range of environments. This does not require novel algorithmic components or any other unexpected innovation, but rather careful implementation of one of several existing algorithmic frameworks (and indeed multiple algorithms can work well), combined with modern optimized deep learning packages and a number of careful design decisions for the MDP formulation of the locomotion task. This result runs counter to the principles articulated in several prior works, which suggest either that simulated training is critical for quadrupedal locomotion because the training times are too long [22, 32–35], that demonstration data is needed to overcome local optima challenges [15, 35], or that more sophisticated algorithmic frameworks with model-based RL [23, 29] or hand-designed movement primitives [29, 36–39] are necessary for real-world training.

Our main contribution is an empirical demonstration that current deep RL methods can effectively learn quadrupedal locomotion directly in the real world in under 20 minutes. While our results largely build on existing methods, we demonstrate for the first time that a careful combination of existing components can enable direct real-world learning of locomotion skills with drastically lower training periods than reported in prior work. Our evaluation includes real-world training in four different locations (one indoor, three outdoor,

	Experimental Design				Training Statistics			
	Hardware	Actions	Resets	Terrains	Simulation		Real World	
					Samples	Hours	Samples	Hours
Ours	A1	PD targets	Learned	In/Outdoor	0	0	$20 \cdot 10^3$	0.33
Wu et al. [23]	A1	PD targets	None	Indoor	0	0	$72 \cdot 10^3$	1
Smith et al. [24]	A1	PD targets	Learned	In/Outdoor	10^6	N/A	$22.5 \cdot 10^3$	1
Kumar et al. [16]	A1	PD targets	N/A	In/Outdoor	$1.2 \cdot 10^9$	24	0	0
Rudin et al. [22]	ANYmal	PD targets	N/A	Indoor	$14.7 \cdot 10^6$	0.33	0	0
Lee et al. [25]	ANYmal	PMTG [26] parameters	N/A	In/Outdoor	N/A	16	0	0
Ha et al. [27]	Minitaur	PD targets	Engineered	Indoor	0	0	$60 \cdot 10^3$	1.5
Haarnoja et al. [28]	Minitaur	PD targets	Manual	Indoor	0	0	$160 \cdot 10^3$	2
Yang et al. [29]	Minitaur	PMTG [26] parameters	Unknown	Indoor	0	0	$45 \cdot 10^3$	0.167

Table I: Overview of experimental details (hardware platform used, the kinds of actions, and access to resets) and data requirements of the works most relevant—in terms of setup and task—to ours (ordered chronologically). We list the approximate numbers reported for the tasks that most closely resemble ours (walking forward). For “Training Statistics” we list the data used for *training*, i.e., we do not include data used in evaluation (that may be used in few-shot adaptation [16]), and the wall-clock time associated with collecting that data either in simulation or in the real world. Some earlier works demonstrated learning in the real world on the Minitaur [30], a relatively simple quadrupedal robot with neither shoulder nor hip abduction, in controlled, lab settings. RL with more complex quadrupedal robots [31] often used large amounts of simulation data. More similar in physical capabilities to the ANYmal and in accessibility to the Minitaur, the A1 robot has also been used to study real-world deployment in recent works.

see Figure 1), and detailed simulated analysis to understand the relative importance of different design choices.

II. RELATED WORK

Roboticians have traditionally designed controllers for quadrupedal locomotion using a combination of footstep planning, trajectory optimization, and model-predictive control (MPC) [31, 40–42]. Learning provides an appealing alternative to classic model-based methods, as it avoids the need for intricate modeling and extensive domain knowledge about the locomotion task, instead allowing the learning algorithm to discover a gait that works well for a given robot. Furthermore, learning in the real world allows the robot to improve over time with “on-the-job” experience.

A major question in this line of work is whether RL methods could ever be efficient enough to learn dynamic locomotion in diverse real-world environments with more complex robots. Indeed, some works have sought to sidestep this challenge by using simulation to learn controllers that are directly transferred to the real world [7, 10, 13, 22, 25, 32, 43–46] or learning policies that can use small amounts of data when deployed in the real world to search among the space of policies it has for a suitable one [11–16, 47]. Still, these methods rely entirely on the strategies the robot was able to learn while training in simulation. So, the training conditions must be designed and implemented such that the training captures the behaviors that are transferable to the real-world conditions. While this approach to learning controllers is often sufficient, it is non-trivial to foresee all the possible situations the robot may encounter when deployed in the real world and engineer the appropriate training environments in simulation to prepare for them. Furthermore, these learned models are fundamentally limited to those training experiences—they cannot perfectly generalize when they are tested in situations that differ enough from their training experience, but these unexpected situations are likely to arise in the real world that is highly complex, unstructured, and unpredictable. Early work in learning *directly* in the real world explored utilizing higher level action spaces [29, 36–39, 48, 49] or fine-tuning [24] to do

real-world training. We show that careful implementation of already known model-free RL techniques can enable learning to walk from scratch using low-level control, such as PD targets, in several real-world environments. Next, we detail prior work in using model-free RL for real-world training of locomotion policies, followed by model-based work, and finally, systems aspects relevant to contextualizing these approaches. For a tabular overview comparing the assumptions, requirements, and results of the works most relevant to ours, see Table I.

a) *Model-free learning*: Much of the early work on learning from scratch in the real world applied model-free policy gradient methods to modulate pre-defined motions on relatively simple hardware. Kohl et al. [36] learn parameters of a pre-defined open-loop trajectory generator to control a Sony AIBO quadruped, and Tedrake et al. [50] learn a feedback controller to improve the control of a passive dynamic walker whose design is such that it walks when control is disabled, using policy gradient methods. Luck et al. [48] achieved sample-efficient learning by leveraging periodicity to similarly learn few parameters to tweak a finned robot’s pre-defined controller. Choi et al. [39] also use a stochastic policy gradient method to learn to control a Snapbot by learning a distribution over Gaussian random paths that define fixed length joint trajectories that are then realized by an open-loop controller. Yang et al. [49] learns a model-free RL policy in the real world which outputs gait parameters and foot placements that are realized by a low-level MPC controller. While very effective for learning walking motions in the real world by shaping exploration and ensuring safety, using high-level action spaces limits the types of skills that can be learned. Recent works [27, 28] have utilized off-policy methods [51] to perform sample-efficient learning on a Minitaur (2 active DOF per leg) by directly outputting target joint positions. Training requires roughly 2 hours (160k control steps) to learn to walk forward and 1.5 hours (60k control steps per direction) to learn to walk forward and backward, respectively. Our work also uses off-policy model-free RL to learn from scratch in the real world—uniquely, though, we train an A1 quadrupedal robot (3 active DOF per leg) not only in lab settings, but also on

outdoor irregular terrains, in just 20 minutes.

b) Model-based learning: Model-based RL methods have been applied to enable an AI robot to learn how to walk on flat ground [23, 29]. Yang et al. [29] use trajectory generators [26], which output smooth, periodic leg trajectories and use a model-based method to modulate these trajectories. Concurrent work by Wu et al. [23] uses a latent dynamics model to generate additional training data in order to reduce the burden on collecting real-world samples [52] in order to learn to walk with low-level PD targets as actions within an hour. In contrast, we use a learned policy to automatically provide resets, and our robot learns to walk within 20 minutes in five environments, three of which are natural outdoor terrains, using a simple model-free method.

c) Systems considerations: RL is known to require lots of data to learn complex behaviors. For example, a state-of-the-art blind quadrupedal locomotion policy trained completely in simulation required 12 hours [25]; however, real-world data is significantly more expensive to collect. Some have sought to eliminate the data collection bottleneck, using upwards of thousands of workers to collect experience simultaneously and consolidating the information for policy updates [53–55]. In the locomotion domain, Rudin et al. [22] use this parallelism to train a simulated ANYmal quadruped to walk on uneven terrain—in 20 minutes of wall-clock time—and then deploy the learned policy in the real world. However, massively parallel data collection is not often feasible in the real world. As such, the focus of another line of work has been on enabling sample-efficient methods with asynchronous pipelines. Due to computational costs, many prior works claim asynchronous training to be necessary for real-world training [23, 28, 51, 56, 57]. Most similar in spirit to our work, Harnoja et al. [28] have a three-part asynchronous training pipeline, with jobs dedicated to data collection, motion capture for state and reward estimation. Follow-up work [27] achieved superior sample efficiency (approximately half the samples required) with the same underlying algorithm by performing synchronous training at a per-step basis (as opposed to episodic, asynchronous training), thereby reducing overall wall-clock time by a similar factor. Recently, more efficient model-free methods have been enabled by placing a larger burden on the computation. In our work, we find that we are able to support synchronous, per-step training in our implementation (see Subsection III-B), achieving learning in 20 minutes on a real robot using a single GPU laptop.

III. FAST AND SIMPLE REAL-WORLD RL

In this section, we describe the algorithmic framework we use, which is based on standard Q-function actor-critic methods [58], building most directly on DroQ [59], which extends the SAC algorithm [51] with Dropout [60] and Layer Normalization [61]. We emphasize that our result is not enabled so much by any one algorithmic component (though the algorithm is also important!), but rather careful implementation and task setup, which we discuss in Section IV.

A. Preliminaries

Learning to walk can be formalized as a Markov Decision Process (MDP), defined by a tuple $(\mathcal{S}, \mathcal{A}, p_0, p, r, \gamma)$ where $\mathcal{S} \subset \mathbb{R}^n$ is the state space, $\mathcal{A} \subset \mathbb{R}^m$ is the action space, $p_0(\cdot)$ is the initial state distribution, $p(\cdot|s, a)$ governs dynamics, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defines rewards, and $\gamma \in [0, 1)$ is the discount factor. The goal of RL is to optimize the expected discounted cumulative return induced by the policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$:

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right].$$

We consider actor-critic methods consisting of interleaved policy evaluation and improvement steps. During policy improvement, we fit a critic to estimate the discounted returns of the current training policy starting at state s and action a :

$$J(\theta) = \mathbb{E}_{(s, a, s') \sim D} [(Q_\theta(s, a) - y(s, a, s'))^2]$$

$$y(s, a, s') = r(s, a) + \gamma Q_{\theta'}(s', a') \text{ where } a' \sim \pi_{\eta}(\cdot|s')$$

where $Q_{\theta'}$ is a target network with weights updated via exponential moving average. Then, the critic is used to improve the policy by maximizing the objective:

$$J(\eta) = \mathbb{E}_{\substack{s \sim D \\ a \sim \pi_\eta(\cdot|s)}} [Q_\theta(s, a)].$$

The most widely-used off-policy RL algorithms for continuous control, such as SAC and DDPG [51, 62], are commonly trained by making one critic update and one policy update after each environment step. Since the update time can exceed what is allowed by the robot’s control frequency, sometimes real-world training is performed by collecting an entire trajectory and then updating the critic and policy for a number of steps equal to the length of the trajectory [28]. Next we will discuss the simple changes to these common practices we make that has enabled our result.

B. Efficient Model-Free RL

Actor-critic methods have recently become significantly more sample-efficient by improving critic training, thereby allowing more learning on the same amount of training data. For example, REDQ [63] extends SAC [51] by utilizing a large ensemble of critics and randomly sampling subsets of them. DroQ [59] regularizes the critic with Dropout [60] and Layer Normalization [64]. The common thread is some sort of regularization or normalization (or both) that mitigates the tendency of off-policy, bootstrapped critic updates to overfit and lead to overestimation [59, 63]. These techniques allow for taking significantly more gradient steps on the critic (up to 20) after each environment step, which in turn leads to significantly more sample-efficient learning. As we will discuss in our analysis in Section VI, a variety of approaches lead to the critical jump in improvement over the baseline method. These results suggest that the key is not any one particular choice but the general principle of augmenting actor-critic RL with regularization or normalization.

While these methods speed up learning with respect to the number of samples collected, the increased computational cost actually supersedes data collection as the primary bottleneck for real-world training. That is, a naïve implementation cannot train as fast as the samples are collected. Prior work has addressed this either by performing asynchronous training [23, 28] or training in-between trials [24]. Both add a delay between the agent interacting with the environment and learning from the samples, which slows down training. Our choice of algorithm and implementation is aimed at enabling real-time synchronous training, which we expand on in Section VI. For implementation, we use JAX [65], a modern machine learning framework that performs just-in-time compilation to optimize execution significantly (see Section V for a discussion of the practical aspects).

IV. SYSTEM DESIGN

We design our system prioritizing fast training in unstructured real-world environments. We use a relatively low-level action space—rather than using pre-defined motion primitives as discussed is used in many prior works (see Section II)—and use only proprioceptive information—so as to be able to train anywhere, without an instrumented motion capture system. Of course, the MDP definition and implementation often has a large impact on learning, e.g., we found the design of the action space to be particularly important. In the remainder of this section, we detail these design decisions. For our robot platform, we use the A1 robot from Unitree and build our simulation using MuJoCo [66] and DM Control [67]. The policy π_i and Q-function $\{Q_{\theta_i}\}_{i=1}^{N_{\text{ensemble}}}$ are modeled using separate fully-connected neural networks that are constructed and trained using JAX [65].

A. State and Action Spaces

For learning locomotion controllers, the robot’s position is used in order to provide reward supervision [67], and privileged information is often used to train policies in simulation [16, 25]. As such, such policies cannot trivially be further trained in the real world. When training in the real world, motion capture systems have been used to localize the robot [27, 28]. However, in order to train in the wild, the robot must be able to receive feedback purely from its onboard sensors. In terms of actions, for generality, we parameterize the policy to directly output joint targets rather than rely on pre-defined gaits or trajectory generators.

The state \mathbf{s}_t contains the root orientation, root angular velocity, root linear velocity, joint angles, joint velocities, binary foot contacts, and the previous action. For the root orientation, we include the roll and pitch. For angular velocity, we include roll, pitch, and yaw information (as to penalize excessive turning). We use an onboard linear velocity estimator which combines integrated acceleration and leg velocity estimated with forward kinematics via Kalman filter that was shown to suffice for reward supervision in outdoor environments [24]. Actions \mathbf{a}_t are PD position targets for each of the 12 joints and applied at a frequency of 20Hz.

We define the action space for every leg as $[p - o, p + o]$ where p corresponds to initial robot pose and o is an action offset. Following Fu et al. [68], we use $o = [0.2, 0.4, 0.4]$. We confirm that constraining the action space is crucial for training the agent with RL in Section VI. Moreover, we found that the initial robot pose also impacts exploration significantly. In the simulator, we used $p = [0.05, 0.7, -1.4]$; however, during the early experiments in the real world, we found that $p = [0.05, 0.9, -1.8]$ promotes safer exploration on the robot since this configuration leads to fewer failures. We use a position controller where the torques commanded to the robot are computed as $\tau = K_p(q^* - q) - K_d \cdot \dot{q}$. In this case, q^* and q define the desired and current positions correspondingly, while \dot{q} defined motor velocities. K_p and K_d define motor gains and damping.

B. Reward Function

Prior works have reported using very complex reward functions consisting of upwards of tens of terms to give rise to a desired motion [15, 25, 68, 69]. We found that with the state and action space choices described above, even a simple reward function was sufficient to produce naturalistic gaits. We define the reward function following standard DM Control [67] locomotion tasks, where the agent is provided with a constant reward within a target velocity interval, outside of which the reward linearly decays to 0.

$$r(s, a) = r_v(s, a) - 0.1v_{yaw}^2$$

where v_{yaw} is an angular yaw velocity and

$$r_v(s, a) = \begin{cases} 1, & \text{for } v_x \in [v_t, 2v_t] \\ 0, & \text{for } v_x \in (-\infty, -v_t] \cup [4v_t, \infty) \\ 1 - \frac{|v_x - v_t|}{2v_t}, & \text{otherwise.} \end{cases}$$

where v_t is the target velocity, while v_x is a forward linear velocity in the robot frame. Our choice of the reward function is motivated by simplicity.

During early experiments with the real robot, we found that using the forward velocity in the robot’s local frame caused it to dive forward as falling quickly onto its head does correspond to high forward velocity from the robot’s perspective. However, our goal is for the robot to move forward in the global frame, parallel to the ground plane, so we instead project the forward velocity onto the ground plane. To keep the robot from leaving the training area with minimal interruption, we simply continue training while providing a reward of 0 when we detect that someone has lifted the robot to move it (i.e. when there are no foot contacts detected). Otherwise, the robot trains continuously, terminating only when the robot’s roll or pitch exceeds 30 degrees. To reset the robot in the real world, we use the open-source reset policy from Smith et al. [24].

V. LEARNING IN THE REAL WORLD

We aim to answer the following through our experiments:

- (1) How quickly and consistently can the robot learn to walk in the real world using model-free RL?

- (2) Can this approach enable a robot to learn to walk not only on flat ground in the lab, but also ‘in the wild’?

In order to understand the variance across random seeds and small changes in real-world conditions (e.g., the status of the hardware, how many times the robot was redirected, etc.), for one of our experiments, we train the robot four times in the same controlled environment. To study (2), we train the robot in four additional terrains, three of which are outdoors. Lastly, we discuss our findings from using the design decisions, tuned in simulation as described in section IV, to train in the real world. Videos of all real-world training along with code to reproduce our results can be found on the project website.

A. Setup

We conduct experiments with five terrains (see Figure 2), each of which possesses unique characteristics:

- (1) *Flat, Solid Ground*: We placed tiles of dense foam on the floor for protection that make for a high-friction surface.
- (2) *Memory Foam Mattress*: The other indoor terrain we test on is a 5cm-thick memory foam mattress. The robot’s feet sink fairly deep into the surface—see, e.g., the depression in the surface of the mattress made by the robot in the earliest frame in Figure 2 (second from left), making walking difficult and requiring a unique gait to attain ground clearance.
- (3) *Mulch*: A layer of bark (about a foot deep) makes walking especially difficult as the terrain is not only heavily irregular but also obstructive. As shown in Figure 2 (middle), the robot naturally sinks such that the lower half of its legs are submerged in the bark. Thus, it must learn to churn through it in order to move.
- (4) *Lawn*: The lawn presents a lower friction, cushioned, deformable walking surface.
- (5) *Hiking Trail*: Here the surface material is a compact, dry dirt. As shown in Figure 2 (right), the trail is at a slight incline and irregular (tree roots, pebbles, troughs, etc.), presenting additional challenges.

To be able to train outdoors, we use a laptop (Origin EON15-X) for training equipped with a single NVIDIA GeForce RTX 2070 GPU. For all experiments, we collect data for about 1 minute (1000 time-steps) by sampling from the action space uniformly at random. We then train synchronously, making a policy update (20 critic updates and an actor update) in between each action executed on the robot every 0.05 seconds. On a standard workstation equipped with a NVIDIA GeForce RTX 2070 GPU, our initial JAX implementation of DroQ was capable of performing 700 critic updates per second, roughly corresponding to 35 time-steps per second. However, this was insufficient for training the agent on a laptop, which reiterates the importance of implementation for real-world training. For our final implementation, we jit 20 critic updates corresponding to one environment step, resulting in 2400 critic updates per second which corresponds to running training at 120 Hz. In contrast, for a baseline comparison, our best PyTorch implementation was not fast enough for synchronous

training, only permitting an effective control frequency of 7.5 Hz.

B. Results

We report the average velocity over intervals of 1000 time-steps (corresponding to 1 minute of wall-clock time) during training in Figure 3 with respect to the number of real-world samples collected during which the robot was on the ground (not including the times the robot was lifted and reoriented). Per answering (1), we report the average and standard deviation (solid line and shaded region, respectively) across four runs on the flat, solid ground (pink). In all cases, the robot learns to walk in less than 20,000 samples (roughly 17 minutes’ worth of data which amounts to 20 total minutes of wall-clock time due to various minor sources of overhead: jitting (about 1 minute), resetting and reorienting the robot, etc.). At the start of training, the robot mostly shuffles, slowly edging backwards. In 5 minutes, the robot learns to inch forward but is unstable. Within 10 minutes, the robot learns to take fairly large steps, but it has not learned to maintain balance while taking these larger steps. Finally, after 15 minutes of training, the robot adopts more conservative behavior in order to both walk and remain balanced.

Next, we train on the variety of challenging terrains detailed in Subsection V-A. On the memory foam, the robot makes its way out of the initial position it had sunk into during initial data collection (similar to the situation shown in the earliest frame of Figure 2 (second from left)) within about 5 minutes. In contrast to the flat, solid ground, the robot falls much less frequently, but its feet often catch on the fabric cover and it needs to learn to walk in a way such that it avoids dragging the fabric with it in order to traverse it. On the thick layer of wood mulch, the initial data collection essentially digs the robot quite deep into the loose ground. In the first ten minutes, the robot learns to kick up its front feet in order to dig itself through, as seen most clearly in the latest frame of Figure 2 (middle) where the robot is kicking its front, right leg up to pull itself forward. Furthermore, the terrain is irregular, yet it is able to quickly adapt its behavior. In all cases the robot learns to traverse each given terrain with roughly 20,000 samples (20 minutes).

VI. ANALYSIS OF DESIGN DECISIONS

This section presents simulated comparisons of design decisions and SAC variants we considered in this work. Since our goal is to run training on a real robot, we aim for design decisions and algorithms that lead to improved stability and sample efficiency. To match the real-world setup, we simulate the official A1 model in MuJoCo, and used the same position controller and rewards as discussed in Section III-B. We provide the exact model we used on the project website.

a) *MDP formulation*: First, we observe that the value of damping for the position controller used for the robot significantly impacts learning (see Figure 4a). Small damping values ($Kd = 1$) lead to instabilities, which are not desirable for a policy executed on a real robot, while large values



Fig. 2: Examples of learned gaits acquired on a variety of real-world terrains. Left to right: flat, solid ground covered in dense foam mats; a 5cm memory foam mattress; loose ground comprised of eucalyptus bark; a grassy lawn; a gently sloped hiking trail.

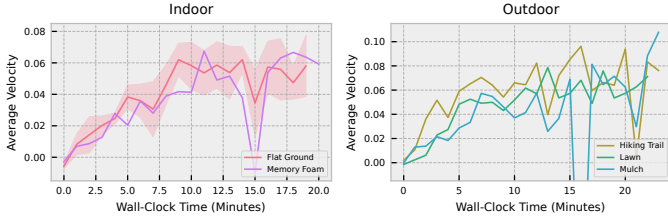


Fig. 3: Learning curves for all the real-world experiments showing the average velocity of the stochastic policy with respect to real-world, wall-clock time. Note that the robot runs continuously, training in between sending actions to the physical hardware. 1000 time-steps corresponds roughly to 1 minute of training time, all things considered.

prevent the agent from reaching the target velocity ($Kd = 20$). Therefore, for the remaining ablations, we used the value of damping set to 10. We also evaluate other design decisions in Figure 4b. In particular, we confirm the efficacy of constraining the action space: we observe that the simulated agent cannot make any progress in the unconstrained action space, while constraining the space leads to stable training and does not prevent the agent from reaching the target velocity. Finally, we note that applying a low-pass filter to the PD targets to promote smoothness, as is common practice [15, 24, 27], degrades learning efficiency. This is perhaps unsurprising, as the filter creates a dependency on action history that violates the Markovian assumption.

b) Algorithms: Our goal is to train a robot to walk in the real world as efficiently as possible, where efficiency includes computational complexity, sample complexity, and total wall-clock time. As discussed in Section III-B, we utilize off-policy, model-free actor-critic methods (basing all on SAC [51]), and investigate a variety of regularization and normalization techniques to accelerate them. We present an extensive comparison in Figure 4c, with the aim of understanding which ingredients are important for attaining the requisite sample efficiency for training in the real world. Standard SAC (purple) with an update-to-data (UTD) ratio of 1 takes one gradient step on the critic for every one time step of data collection. Efficiency can be increased by taking more gradient steps, and we show standard SAC with a larger UTD ratio of 20 (dark blue) as well. We see that naively increasing the number of critic updates made per time-step improves sample efficiency, but still requires roughly 30k samples, which would amount to roughly 30 minutes’ worth of data, to reach the target velocity. We implement REDQ [63] by modifying SAC with UTD ratio of 20 to use random subsets of a larger ensemble of networks in order to calculate target values, and we see this regularization indeed leads to improvement (roughly 5k fewer

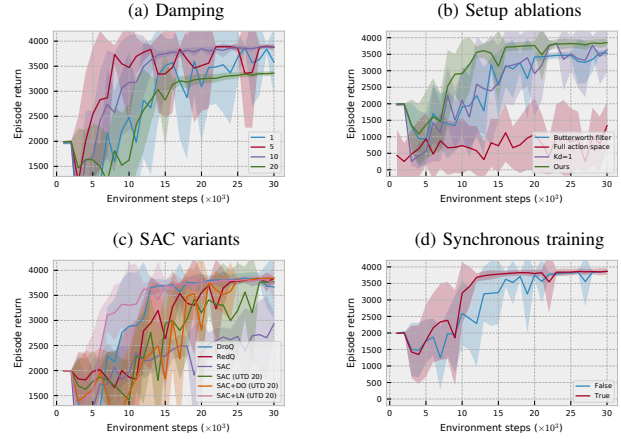


Fig. 4: Experimental evaluation of (a) performance for different value of the damping parameter for the position PD controller; (b) ablations of various task setup choices; (c) the effect of the frequency of policy updates (between time-steps versus episodes). Each curve and shaded region represents the average and standard deviation, respectively, across 10 random seeds.

time-steps, or 5 minutes wall-clock time, required). As noted in DroQ [59], REDQ is more computationally expensive due to the large ensemble, and regularizing the critic with a combination of LayerNorm and Dropout can lead to similar benefits at lower compute cost (blue). However, we also consider LayerNorm (pink) and Dropout (orange) in isolation (each also with a UTD ratio of 20). Perhaps surprisingly, Dropout alone already performs similarly to REDQ, whereas LayerNorm (even without Dropout) leads to even better performance. We conclude that a variety of regularization or normalization methods, if implemented and applied carefully, can all achieve a similar level of improvement in performance over their underlying algorithm *in our setup*. That is, the important thing is not any single specific regularization technique, but the use of any suitable regularization so as to enable SAC to effectively use higher UTD ratios. We also compare updating the agent between episodes and after every environment step and notice that getting immediate feedback leads to more stable training and faster convergence (see Figure 4d). In order to facilitate this kind of training synchronously, the updates must be inexpensive enough to be able to perform them between time-steps (of which there are 20 per second). As such, we favor using the less computationally expensive DroQ variants over others in the real world.

VII. CONCLUSION

We present our finding that we can train quadrupedal robots to walk via deep RL in real-world settings, such as grass, loose ground, forest trails, and mattresses, with about 20 minutes of training. We demonstrate that this can be enabled using a high-quality implementation of existing algorithmic ideas combining standard actor-critic algorithms with one of a number of regularization strategies. We compare the different design choices in simulation and show that a variety of designs can work well, and then demonstrate in the real world that this leads to highly efficient and successful learning on a range of different terrains. Our empirical results show that real-world training of locomotion policies via RL can be significantly more practical than previously believed, and does not necessarily require significant deviations from existing practice in RL, but rather careful combination of current best practices. We hope that our work will serve to further encourage investigation of real-world RL in robotics.

Acknowledgements

This work was supported by the Office of Naval Research and DARPA RACER. Laura Smith is supported by NSF Graduate Research Fellowship. We thank Kevin Zakka and Vikash Kumar for their help with the A1 MuJoCo model and Philipp Wu for designing and printing a protective shell for the physical robot.

REFERENCES

- [1] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, “Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation,” *Conference on Robot Learning (CoRL)*, vol. abs/1806.10293, 2018.
- [2] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen, “Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection,” *The International Journal of Robotics Research*, vol. 37, pp. 421 – 436, 2018.
- [3] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and C. Finn, “Robonet: Large-scale multi-robot learning,” *ArXiv*, vol. abs/1910.11215, 2019.
- [4] D. Kalashnikov, J. Varley, Y. Chebotar, B. Swanson, R. Jonschkowski, C. Finn, S. Levine, and K. Hausman, “Mt-opt: Continuous multi-task robotic reinforcement learning at scale,” *ArXiv*, vol. abs/2104.08212, 2021.
- [5] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and S. Levine, “Bridge data: Boosting generalization of robotic skills with cross-domain datasets,” *ArXiv*, vol. abs/2109.13396, 2022.
- [6] M. Cutler, T. J. Walsh, and J. P. How, “Reinforcement learning with multi-fidelity simulators,” *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3888–3895, 2014.
- [7] F. Sadeghi and S. Levine, “CAD²RL: Real single-image flight without a single real image,” *Robotics: Science and Systems (RSS)*, 2017.
- [8] A. Rajeswaran, S. Ghotra, S. Levine, and B. Ravindran, “Epopt: Learning robust neural network policies using model ensembles,” *International Conference on Learning Representations (ICLR)*, vol. abs/1610.01283, 2017.
- [9] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, “Domain randomization for transferring deep neural networks from simulation to the real world,” *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 23–30, 2017.
- [10] X. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-real transfer of robotic control with dynamics randomization,” *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1–8, 2018.
- [11] W. Yu, V. Kumar, G. Turk, and C. Liu, “Sim-to-real transfer for biped locomotion,” *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3503–3510, 2019.
- [12] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, “Learning agile and dynamic motor skills for legged robots,” *Science Robotics*, vol. 4, 2019.
- [13] Z. Xie, P. Clary, J. Dao, P. Morais, J. Hurst, and M. V. D. Panne, “Learning locomotion skills for cassie: Iterative design and sim-to-real,” in *Conference on Robot Learning (CoRL)*, 2019.
- [14] W. Yu, J. Tan, Y. Bai, E. Coumans, and S. Ha, “Learning fast adaptation with meta strategy optimization,” *IEEE Robotics and Automation Letters*, vol. 5, pp. 2950–2957, 2020.
- [15] X. Peng, E. Coumans, T. Zhang, T. Lee, J. Tan, and S. Levine, “Learning agile robotic locomotion skills by imitating animals,” *Robotics: Science and Systems (RSS)*, vol. abs/2004.00784, 2020.
- [16] A. Kumar, Z. Fu, D. Pathak, and J. Malik, “Rma: Rapid motor adaptation for legged robots,” *Robotics: Science and Systems (RSS)*, 2021.
- [17] J. X. Wang, Z. Kurth-Nelson, H. Soyer, J. Z. Leibo, D. Tirumala, R. Munos, C. Blundell, D. Kumaran, and M. M. Botvinick, “Learning to reinforcement learn,” *ArXiv*, vol. abs/1611.05763, 2017.
- [18] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International Conference on Machine Learning (ICML)*, 2017.
- [19] K. Rakelly, A. Zhou, D. Quillen, C. Finn, and S. Levine, “Efficient off-policy meta-reinforcement learning via probabilistic context variables,” *International Conference on Machine Learning (ICML)*, vol. abs/1903.08254, 2019.
- [20] X. Song, Y. Yang, K. Choromanski, K. Caluwaerts, W. Gao, C. Finn, and J. Tan, “Rapidly adaptable legged robots via evolutionary meta-learning,” *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3769–3776, 2020.

- [21] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. A. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. M. Zhang, “Solving rubik’s cube with a robot hand,” *ArXiv*, vol. abs/1910.07113, 2019.
- [22] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, “Learning to walk in minutes using massively parallel deep reinforcement learning,” *ArXiv*, vol. abs/2109.11978, 2021.
- [23] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel, “Daydreamer: World models for physical robot learning,” *ArXiv*, vol. abs/2206.14176, 2022.
- [24] L. Smith, J. C. Kew, X. B. Peng, S. Ha, J. Tan, and S. Levine, “Legged robots that keep on learning: Fine-tuning locomotion policies in the real world,” *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [25] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, “Learning quadrupedal locomotion over challenging terrain,” *Science Robotics*, vol. 5, 2020.
- [26] A. Iscen, K. Caluwaerts, J. Tan, T. Zhang, E. Coumans, V. Sindhwani, and V. Vanhoucke, “Policies modulating trajectory generators,” *Conference on Robot Learning (CoRL)*, vol. abs/1910.02812, 2018.
- [27] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan, “Learning to walk in the real world with minimal human effort,” *ArXiv*, vol. abs/2002.08550, 2020.
- [28] T. Haarnoja, S. Ha, A. Zhou, J. Tan, G. Tucker, and S. Levine, “Learning to walk via deep reinforcement learning,” *Robotics: Science and Systems (RSS)*, 2020.
- [29] Y. Yang, K. Caluwaerts, A. Iscen, T. Zhang, J. Tan, and V. Sindhwani, “Data efficient reinforcement learning for legged robots,” *Conference on Robot Learning (CoRL)*, vol. abs/1907.03613, 2019.
- [30] G. D. Kenneally, A. De, and D. E. Koditschek, “Design principles for a family of direct-drive legged robots,” *IEEE Robotics and Automation Letters*, vol. 1, pp. 900–907, 2016.
- [31] M. Hutter, C. Gehring, D. Jud, A. Lauber, D. Bellisico, V. Tsounis, J. Hwangbo, K. Bodie, P. Fankhauser, M. Bloesch, R. Diethelm, S. Bachmann, A. Melzer, and M. Höpflinger, “Anymal - a highly mobile and dynamic quadrupedal robot,” *IEEE International Conference on Intelligent Robots and Systems (IROS)*, pp. 38–44, 2016.
- [32] J. Tan, Z. Xie, B. Boots, and C. Liu, “Simulation-based design of dynamic controllers for humanoid balancing,” *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2729–2736, 2016.
- [33] L. Liu and J. Hodgins, “Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning,” *ACM Transactions on Graphics (TOG)*, vol. 37, pp. 1 – 14, 2018.
- [34] S. Lee, M. Park, K. Lee, and J. Lee, “Scalable muscle-actuated human simulation and control,” *ACM Transactions on Graphics (TOG)*, vol. 38, pp. 1 – 13, 2019.
- [35] X. Peng, P. Abbeel, S. Levine, and M. V. D. Panne, “Deepmimic: Example-guided deep reinforcement learning of physics-based character skills,” *ACM Trans. Graph.*, vol. 37, pp. 143:1–143:14, 2018.
- [36] N. Kohl and P. Stone, “Policy gradient reinforcement learning for fast quadrupedal locomotion,” *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA '04. 2004*, vol. 3, pp. 2619–2624 Vol.3, 2004.
- [37] R. Tedrake, T. Zhang, and H. Seung, “Stochastic policy gradient reinforcement learning on a simple 3d biped,” *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (IEEE Cat. No.04CH37566)*, vol. 3, pp. 2849–2854 vol.3, 2004.
- [38] G. Endo, J. Morimoto, T. Matsubara, J. Nakanishi, and G. Cheng, “Learning cpg sensory feedback with policy gradient for biped locomotion for a full-body humanoid,” in *AAAI*, 2005.
- [39] S. Choi and J. Kim, “Trajectory-based probabilistic policy gradient for learning locomotion behaviors,” *2019 International Conference on Robotics and Automation (ICRA)*, pp. 1–7, 2019.
- [40] H.-W. Park, P. M. Wensing, and S. Kim, “High-speed bounding with the mit cheetah 2: Control design and experiments,” *The International Journal of Robotics Research*, vol. 36, no. 2, pp. 167–192, 2017. [Online]. Available: <https://doi.org/10.1177/0278364917694244>
- [41] G. Bledt, M. J. Powell, B. Katz, J. Carlo, P. Wensing, and S. Kim, “Mit cheetah 3: Design and control of a robust, dynamic quadruped robot,” *IEEE International Conference on Intelligent Robots and Systems (IROS)*, pp. 2245–2252, 2018.
- [42] B. Katz, J. Carlo, and S. Kim, “Mini cheetah: A platform for pushing the limits of dynamic quadruped control,” *2019 International Conference on Robotics and Automation (ICRA)*, pp. 6295–6301, 2019.
- [43] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. D. Ratliff, and D. Fox, “Closing the sim-to-real loop: Adapting simulation randomization with real world experience,” *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 8973–8979, 2019.
- [44] Y. Du, O. Watkins, T. Darrell, P. Abbeel, and D. Pathak, “Auto-tuned sim-to-real transfer,” *ArXiv*, vol. abs/2104.07662, 2021.
- [45] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke, “Sim-to-real: Learning agile locomotion for quadruped robots,” *ArXiv*, vol. abs/1804.10332, 2018.
- [46] Y. Yang, T. Zhang, E. Coumans, J. Tan, and B. Boots, “Fast and efficient locomotion via learned gait transitions,” *ArXiv*, vol. abs/2104.04644, 2021.
- [47] Z. He, R. C. Julian, E. Heiden, H. Zhang, S. Schaal, J. J. Lim, G. Sukhatme, and K. Hausman, “Zero-shot skill composition and simulation-to-real transfer by learning task representations,” *ArXiv*, vol. abs/1810.02422, 2018.
- [48] K. S. Luck, J. Campbell, M. A. Jansen, D. M. Aukes, and H. B. Amor, “From the lab to the desert: Fast prototyping

- and learning of robot locomotion,” *Robotics: Science and Systems (RSS)*, vol. abs/1706.01977, 2017.
- [49] T.-Y. Yang, T. Zhang, L. Luu, S. Ha, J. Tan, and W. Yu, “Safe reinforcement learning for legged locomotion,” *ArXiv*, vol. abs/2203.02638, 2022.
- [50] R. Tedrake and H. Seung, “Learning to walk in 20 minutes,” 2005.
- [51] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *ICML*, 2018.
- [52] D. Hafner, T. P. Lillicrap, J. Ba, and M. Norouzi, “Dream to control: Learning behaviors by latent imagination,” *International Conference on Learning Representations (ICLR)*, 2020.
- [53] A. Nair, P. Srinivasan, S. Blackwell, C. Alcicek, R. Fearon, A. D. Maria, V. Panneershelvam, M. Suleyman, C. Beattie, S. Petersen, S. Legg, V. Mnih, K. Kavukcuoglu, and D. Silver, “Massively parallel methods for deep reinforcement learning,” *ArXiv*, vol. abs/1507.04296, 2015.
- [54] M. Babaeizadeh, I. Frosio, S. Tyree, J. Clemons, and J. Kautz, “Reinforcement learning through asynchronous advantage actor-critic on a gpu,” in *ICLR*, 2017.
- [55] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu, “Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures,” *ArXiv*, vol. abs/1802.01561, 2018.
- [56] S. S. Gu, E. Holly, T. P. Lillicrap, and S. Levine, “Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates,” *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3389–3396, 2017.
- [57] Y. Zhang, I. Clavera, B.-Y. Tsai, and P. Abbeel, “Asynchronous methods for model-based reinforcement learning,” *Conference on Robot Learning (CoRL)*, vol. abs/1910.12453, 2019.
- [58] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [59] T. Hiraoka, T. Imagawa, T. Hashimoto, T. Onishi, and Y. Tsuruoka, “Dropout q-functions for doubly efficient reinforcement learning,” *International Conference on Learning Representations (ICLR)*, 2022.
- [60] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *J. Mach. Learn. Res.*, vol. 15, pp. 1929–1958, 2014.
- [61] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [62] Y.-H. Xu, C.-C. Yang, M. Hua, and W. Zhou, “Deep deterministic policy gradient (ddpg)-based resource allocation scheme for noma vehicular communications,” *IEEE Access*, vol. 8, pp. 18 797–18 807, 2020.
- [63] X. Chen, C. Wang, Z. Zhou, and K. Ross, “Randomized ensembled double q-learning: Learning fast without a model,” *International Conference on Learning Representations (ICLR)*, 2021.
- [64] J. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *ArXiv*, vol. abs/1607.06450, 2016.
- [65] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, “JAX: composable transformations of Python+NumPy programs,” 2018. [Online]. Available: <http://github.com/google/jax>
- [66] E. Todorov, T. Erez, and Y. Tassa, “Mujoco: A physics engine for model-based control,” in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5026–5033.
- [67] S. Tunyasuvunakool, A. Muldal, Y. Doron, S. Liu, S. Bohez, J. Merel, T. Erez, T. Lillicrap, N. Heess, and Y. Tassa, “dm_control: Software and tasks for continuous control,” *Software Impacts*, vol. 6, p. 100022, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2665963820300099>
- [68] Z. Fu, A. Kumar, J. Malik, and D. Pathak, “Minimizing energy consumption leads to the emergence of gaits in legged robots,” *Conference on Robot Learning (CoRL)*, 2021.
- [69] J. Lee, J. Hwangbo, and M. Hutter, “Robust recovery controller for a quadrupedal robot using deep reinforcement learning,” *ArXiv*, vol. abs/1901.07517, 2019.