

*maybe add a  
line return here to  
better distribute text  
among the two lines*

Reinforcement Learning Based Mechanical Design and Control of Flexible-Legged  
Jumping Robots

A Thesis  
Presented to the  
Graduate Faculty of the  
University of Louisiana at Lafayette  
In Partial Fulfillment of the  
Requirements for the Degree  
Masters of Science

Andrew Albright

Spring 2022

© Andrew Albright

2022

All Rights Reserved

Reinforcement Learning Based Mechanical Design and Control of Flexible-Legged  
Jumping Robots

Andrew Albright

APPROVED:

---

Joshua E. Vaughan, Chair  
Assistant Professor of Mechanical  
Engineering

---

Anthony S. Maida  
Associate Professor of Computer  
Science

---

Alan A. Barhorst  
Department Head of Mechanical  
Engineering

---

Mary Farmer-Kaiser  
Dean of the Graduate School

*To all the poor souls using Word, one day you will see the light that is L<sup>A</sup>T<sub>E</sub>X.*

*“Before we work on artificial intelligence why don’t we do something about natural  
stupidity?”*

— Steve Polyak

## Acknowledgments

I would like to firstly thank my advisor Dr. Joshua Vaughan for his leadership, both academic and personal, during my time here at the university. His support has been invaluable in helping me to understand the underlying ideas behind this material. Additionally, his constant drive to produce high quality material has kept me motivated in my attempts to do the same. I would not be here without him.

Additionally, I would like to thank my committee members, Dr. Alan Barhorst, Dr. Anthony Maida and Dr. Brian Post for their input in regards to this work. Along with my lab members during my time in the C.R.A.W.L.A.B., Gerald Eaglin, Adam Smith, Brennan Moeller, Darcy LaFont, Thomas Poche and Y (Eve) Dang. Their conversations and advice have greatly assisted me in my efforts to complete this work.

Lastly, I would like to thank the Louisiana Crawfish Board for their financial support, <sup>which</sup> in the form of a grant from the University of Louisiana at Lafayette. This grant has allowed me to work knowing my fiscal security was intact.

## Table of Contents

Dedication . . . . .	iv
Epigraph . . . . .	v
Acknowledgments . . . . .	vi
List of Tables . . . . .	ix
List of Figures . . . . .	x
<b>I Introduction and System Description . . . . .</b>	<b>1</b>
1.1 Improving Performance with Flexible Components . . . . .	2
1.2 Controlling Flexible Systems . . . . .	3
1.3 Concurrent Design . . . . .	3
1.4 Reinforcement Learning . . . . .	4
1.5 Twin Delayed Deep Deterministic Policy Gradient . . . . .	5
1.6 Contributions . . . . .	10
<b>II Learning Efficient Jumping Strategies for the Monopode System . . . . .</b>	<b>12</b>
2.1 Monopode Jumping System . . . . .	12
2.2 Training Environment . . . . .	14
2.3 Efficient Control Strategies . . . . .	16
2.4 Deploying TD3 . . . . .	17
2.5 Average Performance of Network Controller . . . . .	19
2.5.1 Input Commands . . . . .	19
2.5.2 Jumping Height Performance . . . . .	20
2.5.3 Height Reached vs. Power Used . . . . .	21
2.6 Optimal Performance of the Proposed Controller . . . . .	23
2.6.1 Input Commands . . . . .	23
2.6.2 Jumping Height Performance . . . . .	24
2.7 Conclusion . . . . .	25
<b>III Using Input Shaping to Validate RL Controller . . . . .</b>	<b>26</b>
3.1 Input Shaping Controller Input . . . . .	27
3.2 RL Controller Input . . . . .	28
3.3 Conclusion . . . . .	28
<b>IV Mechanical Design of a the Monopode Jumping System . . . . .</b>	<b>29</b>
4.1 Learning a Mechanical Design . . . . .	29
4.2 Environment Definition . . . . .	31
4.3 Rewards for Learning Designs . . . . .	32
4.4 Design Space Variations . . . . .	32
4.5 Deploying TD3 . . . . .	34

<b>4.6</b>	<b>Jumping Performance . . . . .</b>	35
4.6.1	Narrow Design Space . . . . .	35
4.6.2	Wide Design Space . . . . .	36
4.6.3	Average Design Performance . . . . .	38
<b>4.7</b>	<b>Conclusion . . . . .</b>	39
<b>V</b>	<b>Concurrent Design of the Monopode System . . . . .</b>	42
<b>5.1</b>	<b>Concurrent Design Architecture . . . . .</b>	42
<b>5.2</b>	<b>Mechanical Design Update . . . . .</b>	43
5.2.1	Discrete vs. Continuous . . . . .	43
5.2.2	Averaging Design Policies . . . . .	44
5.2.3	Differing Reward Types . . . . .	44
5.2.4	Design Update Rate . . . . .	44
5.2.5	Design Space Limitations . . . . .	45
<b>5.3</b>	<b>Environment Definition . . . . .</b>	45
5.3.1	Learning the Controller . . . . .	45
5.3.2	Learning the Design . . . . .	46
<b>5.4</b>	<b>Deploying the Algorithm . . . . .</b>	47
<b>5.5</b>	<b>Discrete vs Continuous Designs . . . . .</b>	48
5.5.1	Reward vs Learning Step . . . . .	48
5.5.2	Designs Learned . . . . .	49
5.5.3	Resulting Input and Jumping Performance . . . . .	51
<b>5.6</b>	<b>Effects of Differing Update Rate . . . . .</b>	53
5.6.1	Reward vs Learning Step . . . . .	54
5.6.2	Designs Learned . . . . .	55
5.6.3	Resulting Input and Jumping Performance . . . . .	56
<b>5.7</b>	<b>Conclusion . . . . .</b>	59
<b>VI</b>	<b>Conclusion and Discussion . . . . .</b>	60
6.0.1	Conclusion . . . . .	60
6.0.2	Future Research . . . . .	62
<b>VII</b>	<b>Appendix: Concurrent Design Algorithm . . . . .</b>	64
7.0.1	Averaging $n$ Learned Designs . . . . .	64
7.0.2	Discrete vs. Continuous . . . . .	64
7.0.3	Design Update Rate . . . . .	65
<b>Bibliography</b>	. . . . .	69
<b>Abstract</b>	. . . . .	73
<b>Biographical Sketch</b>	. . . . .	74

## List of Tables

<b>Table 1.</b>	Monopode Model Parameters . . . . .	13
<b>Table 2.</b>	TD3 Training Hyperparameters . . . . .	18
<b>Table 3.</b>	TD3 Training Hyperparameters . . . . .	34
<b>Table 4.</b>	Learned Design Parameters . . . . .	38
<b>Table 5.</b>	Outer Loop TD3 Training Hyperparameters . . . . .	47
<b>Table 6.</b>	TD3 Training Hyperparameters . . . . .	48

## List of Figures

<b>Figure 1.</b> Flexible Robotics System . . . . .	1
<b>Figure 2.</b> Rotary Style Series Elastic Actuator . . . . .	2
<b>Figure 3.</b> Tendon Like Flexibility from [1] . . . . .	3
<b>Figure 4.</b> Reinforcement Learning Process . . . . .	5
<b>Figure 5.</b> Twin Delayed Deep Deterministic Policy Gradient Block Diagram with Monopode as Environment . . . . .	7
<b>Figure 6.</b> Monopode Jumping System . . . . .	13
<b>Figure 7.</b> Jumping Types for the Monopode Jumping System . . . . .	15
<b>Figure 8.</b> Reward vs Time Step During Training . . . . .	18
<b>Figure 9.</b> Average and Standard Deviation Inputs to Monopode . . . . .	19
<b>Figure 10.</b> Average and Standard Deviation Heights of Monopode . . . . .	20
<b>Figure 11.</b> Height Reached vs Power Consumed of Monopode . . . . .	22
<b>Figure 12.</b> Optimal Inputs to Monopode . . . . .	23
<b>Figure 13.</b> Optimal Heights of Monopode . . . . .	24
<b>Figure 14.</b> Jumping Command [2] . . . . .	26
<b>Figure 15.</b> Resulting Actuator Motion [2] . . . . .	27
<b>Figure 16.</b> Decomposition of the Jump Command into a Step Convolved with an Impulse Sequence [2] . . . . .	27
<b>Figure 17.</b> Learning a Mechanical Design . . . . .	30
<b>Figure 18.</b> Reference Jumping Performance of the Monopode . . . . .	33
<b>Figure 19.</b> Reward vs. Episode for Learning Mechanical Design . . . . .	35
<b>Figure 20.</b> Height Reached During Training Given Narrow Design Space . . . . .	35
<b>Figure 21.</b> Designs Learned for the Narrow Design Space . . . . .	36

<b>Figure 22.</b> Height Reached During Training Given Wide Design Space . . . . .	37
<b>Figure 23.</b> Designs Learned for the Wide Design Space . . . . .	37
<b>Figure 24.</b> Height vs Time of Average Optimal Designs . . . . .	38
<b>Figure 25.</b> Reference Jumping Performance of the Monopode . . . . .	41
<b>Figure 26.</b> Concurrent Design Architecture . . . . .	43
<b>Figure 27.</b> Reward During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	49
<b>Figure 28.</b> Spring Constant During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	49
<b>Figure 29.</b> Damping Ratio During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	50
<b>Figure 30.</b> Average Controller Performance for Efficient Concurrent Designs . . .	51
<b>Figure 31.</b> Average Controller Performance for High Jumping Concurrent Designs	52
<b>Figure 32.</b> Reward During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	54
<b>Figure 33.</b> Spring Constant During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	55
<b>Figure 34.</b> Damping Ratio During Training for Discrete and Continuous Implementation Methods of Concurrent Design . . . . .	56
<b>Figure 35.</b> Average Controller Performance for Efficient Concurrent Designs . . .	56
<b>Figure 36.</b> Average Controller Performance for High Jumping Concurrent Designs	58

## I Introduction and System Description

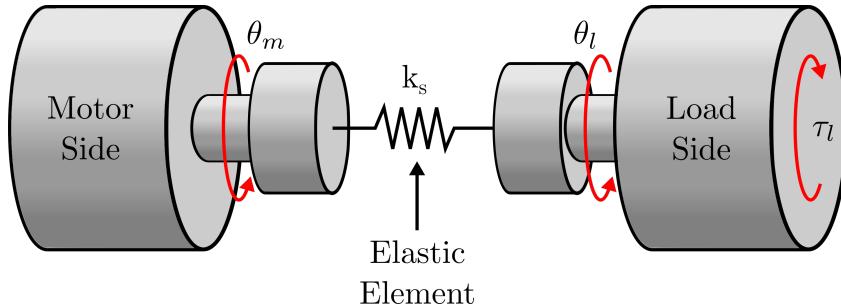
A legged locomotive robot can have many advantages over a wheeled or tracked one, particularly in regards to their ability to navigate uneven and unpredictable terrain [3, 4]. They can achieve this advantage because of the numerous movement types they can deploy. Abilities such as independently placing their feet within highly rigid terrain and jumping or bounding over obstacles have been shown to be effective ways of locomoting [5]. These advantages do ~~not~~ come at ~~no~~ cost, however. Legged systems are traditionally power inefficient compared to wheeled vehicles, making them a less attractive option for applications where power conservation is required. Research has ~~been conducted~~ showing the usefulness of adding flexible components, like the legs seen on the robot in Figure 1, for combating efficiency and other issues [4, 6, 7]. The addition of these components in legged robots has been shown to increase system performance measures such as running speed, jumping capability, and power efficiency [8]. However, the addition of flexible components creates a system that is highly nonlinear, and thus requires a more complex control system.



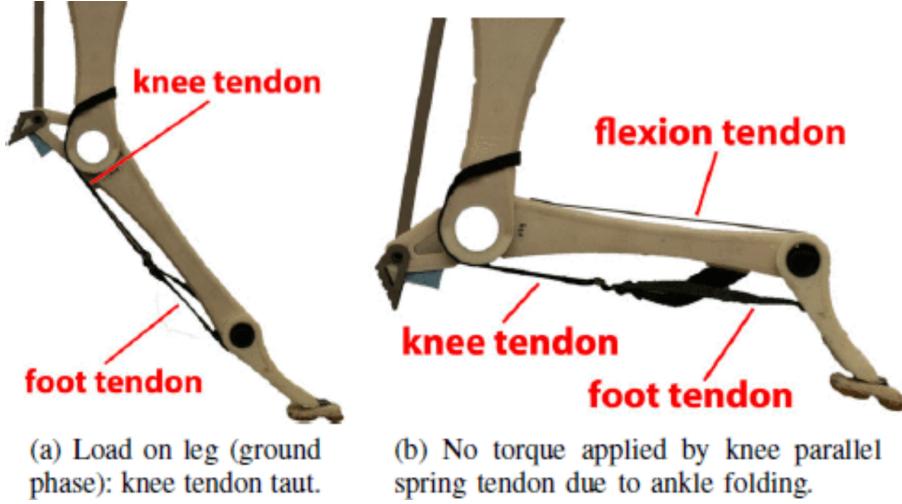
**Figure 1.** Flexible Robotics System

## 1.1 Improving Performance with Flexible Components

The use of flexible components within robotic systems has been shown to be an effective way of improving performance metrics such as movement velocity and power efficiency [4, 8]. Of the different techniques that have been deployed, the use of series elastic actuators (SEAs) has been shown to be an effective for increasing energy efficiency [9, 10]. Storing energy in the non-rigid parts of motor joints, such as the elastic element seen in Figure 2, have proven to be an effective way in increasing efficiency. The addition of flexible joints is not the only technique that has been used to improve performance, however; utilizing tendon like elastic members to connect actuators to links has also been shown to be an effective way of improving efficiency [1]. The use of tendons, being an example of replicating what is found in nature, is a common method of finding unique mechanical designs that perform well in the real world. An example of this type of design can be seen in Figure 3. Following a similar idea, research has also been conducted finding the usefulness of including flexibility in the spine of 2D running robots, leading to dramatic increases in velocity [11]. Research studying the effects of flexible links, like the ones shown in Figure 1, is limited though, particularly in the realm of legged-robots. Still, it has been shown as a viable method of increasing performance in these types of robots [12].



**Figure 2.** Rotary Style Series Elastic Actuator



**Figure 3.** Tendon Like Flexibility from [1]

## 1.2 Controlling Flexible Systems

Control methods developed for flexible systems have been shown to be effective for position control and vibration reduction [10, 13]. Due to the challenges seen in scaling the controllers to highly nonlinear systems, methods utilizing reinforcement learning are of interest. This method has been used in simple planar cases, where it was compared to a PD control strategy for vibration suppression and proved to be a higher performing method [14]. Additionally, it has also been shown to be effective at defining control strategies for flexible-legged locomotion. The use of actor-critic algorithms such as Deep Deterministic Policy Gradient [15] have been used to train running strategies for a flexible-legged quadruped [16]. Much of the research is based in simulation, however, and often the controllers are not deployed on physical systems, which leads to the question of whether or not these are useful techniques in practice.

## 1.3 Concurrent Design

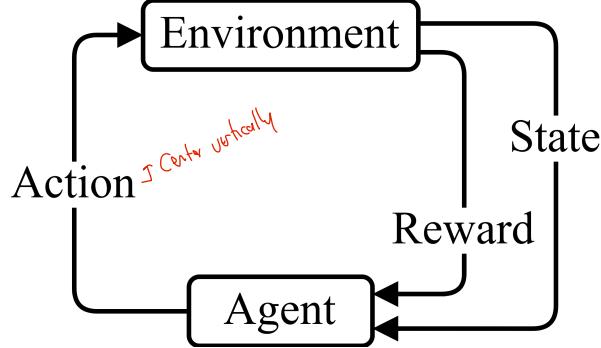
Defining an optimal controller for a system can be difficult due to challenges such as mechanical and electrical design limits. This is especially true when the system is flexible and the model is nonlinear. A solution to this challenge is to concurrently

design a system with the controllers so that the two are jointly optimized. Defining the design process such that the robot's design results in a simple dynamic model has been shown to improve the performance of mechatronics systems [17]. Additionally, in more recent work, the utilization of complex deep learning methods have shown to be an effective strategy for finding optimal concurrent designs [18]. Deep learning has been used to find concurrent designs for simulated legged robotic systems, leading to improved performance in regards to movement velocity [19]. Some research has even been completed where the designs were deployed on physical hardware, validating that this area of research is an effective one for learning how best to define system/controller architectures [20]. Little research exists<sup>①</sup>, however, utilizing these techniques on legged-robotic systems, particularly ones with flexible components.

#### 1.4 Reinforcement Learning

With the recent successes seen in utilizing reinforcement learning (RL) to define control strategies, design parameters, and concurrent designs for robot systems, it is of interest to apply this technique in a unique way to flexible-legged jumping systems. Firstly, it is important to understand the generalities regarding a reinforcement learning problem.

Reinforcement Learning is the process of training a policy to define a series of commands using an environment where those commands can be applied. This is an iterative process<sup>②</sup> which is shown in Figure 4. A policy, often referred to as an agent, from a controls theory perspective, is synonymous with a controller. The environment the controller is deployed in, again from a controls theory perspective, is synonymous with a robotic system. Training the controller requires iteratively deploying the controller's commands, or actions, to the environment and observing the results. The results are often in the form of the state of the environment and a reward resulting from the action that was applied. The reward function is defined by the designer so



**Figure 4.** Reinforcement Learning Process

that the controller is trained to accomplish a desired task. Other than the reward, the controller has no way to discern what commands are good when the environment is in some state.

Learning an optimal control strategy is accomplished by deploying a gradient-decent-based learning algorithm utilizing information such as the state of the environment and the reward. For general robotics applications, at each discrete time step  $t$ , the environment will be in a state  $s \in \mathcal{S}$ , and the controller will select an action  $a \in \mathcal{A}$  according to the current policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  and apply said action within the environment. The environment will transition to a new state  $s'$  and will generate a reward  $r$  based on the user's definition.<sup>1</sup> The return, being the value the algorithm is trying to optimize, is defined as a discounted sum of rewards,  $R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$ , where  $\gamma$  is a discount factor for assigning the level of importance for near-term or long-term rewards.

The challenge of an RL algorithm is to optimize a policy,  $\pi_\phi$ , with parameters,  $\phi$ , such that actions generated at each time step will maximize the return. Ultimately, an optimized policy will maximize the expected return,  $J(\phi) = \mathbb{E}_{s_i \sim p_\pi, a_i \sim \pi}[R_0]$ .

## 1.5 Twin Delayed Deep Deterministic Policy Gradient

There are many algorithms used to train a neural-network-based controller in an RL application, some of which have shown their ability to learn high-performing control

strategies for robotic systems [21–23]. Of the different algorithms used in current research, the one selected and tested in this work is Twin Delayed Deep Deterministic Policy Gradient (TD3) [24]. This is an actor-critic ~~type~~<sup>succesor?</sup> learning algorithm which is widely considered the ~~predecessor~~ to the popular and proven Deep Deterministic Policy Gradient (DDPG) algorithm [15].

Figure 5 displays the flow of information for this algorithm. In general, this algorithm learns both a Q-function and a policy, being the *critic* and the *actor*. For algorithms such as TD3, the ultimate goal is to find a policy,  $\pi_\theta$ , which maximizes the expected return:

$$\nabla_\phi J(\phi) = \mathbb{E}_{s \sim p_\pi} [\nabla_a Q^\pi(s, a)|_{a=\pi(s)} \nabla_\phi \pi_\phi(s)] \quad (1)$$

where  $Q^\pi(s, a) = \mathbb{E}_{s_i \sim p_\pi, a_i \sim \pi}[R_t | s, a]$  is the Q-function (sometimes called the value function) and the critic in the case of the TD3 architecture. This function is based on the Bellman Equation [25] and returns a numerical value from being in a state  $s$ , taking action  $a$ , and following policy  $\pi$  from there after:

$$Q^\pi(s, a) = \mathbb{E}_{s' \sim P} [r(s, a) + \gamma \max_{a'} Q^\pi(s', a')] \quad (2)$$

Using a differentiable function approximator,  $Q^\pi(s, a)$  can be represented and estimated by  $Q_\phi(s, a)$ , with parameters  $\phi$  [26]. Updating the Q-function is accomplished using the temporal difference error between the Q-function and a target Q-function [27, 28]. To maintain a fixed objective over multiple policy updates, the target Q-function approximator is instantiated separately as  $Q_{\phi_{targ}}(s, a)$ . The target does depend on the same parameters that are being trained,  $\phi$ , so there exists an issue when trying to use it as a target. To solve this issue the target network is updated at a delayed pace following the main Q-function approximator by either matching the parameters or by polyak averaging,  $\phi_{targ} \leftarrow \tau\phi + (1 - \tau)\phi_{targ}$ , where  $\tau$  is a tunable hyperparameter.

In summary, the critic side of the TD3 algorithm is responsible for minimizing

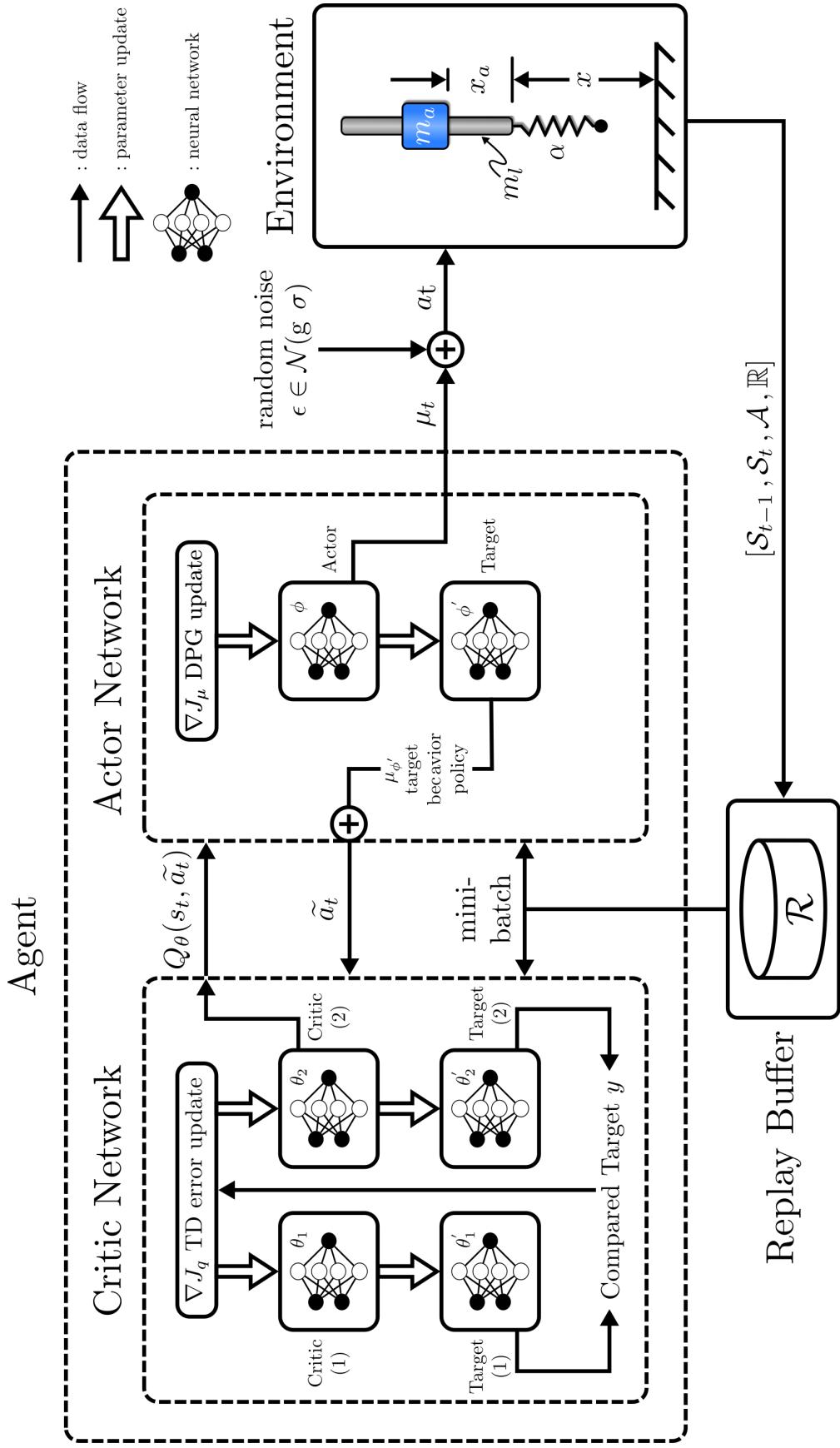


Figure 5. Twin Delayed Deep Deterministic Policy Gradient Block Diagram with Monopode as Environment

the difference between the value of the current state/action pair using the main Q-function, and the reward of the current state/action pair plus the discounted value of the next state/action pair using the target Q-function. The loss function takes the form:

$$L(\phi, \mathcal{D}) = \underset{(s,a,r,s',d) \sim D}{E} \left[ (Q_\phi(s, a) - (r(s, a) + \gamma(1-d) Q_{\phi_{targ}}(s', \pi_{\theta_{targ}}(s'))))^2 \right] \quad (3)$$

where  $\pi_{\theta_{targ}}(s')$  is a target policy that, in a similar manner to the target Q-function, follows the main policy,  $\pi_\theta$ , at a delayed pace either by directly copying the values or by Polyak averaging. Additionally,  $d$  represents a boolean value which depends on the terminal status of the next state,  $s'$ .

As for updating the policy for the actor critic-type algorithm, this aspect is rather simple. Because DDPG, and therefore TD3, are built to accommodate only continuous action spaces, the Q-function is assumed to be differentiable with respect to action. Therefore ~~to~~ to find optimal policy parameters,  $\theta$ , for the policy,  $\pi_\theta$ , the solution of the following equation must be found:

$$\max_{\theta} \underset{s \sim D}{E} [Q_\phi(s, \pi_\theta(s))] \quad (4)$$

The reason that TD3 is considered the successor to DDPG is that there are some additional tricks deployed in addition to the description thus far. The first is the addition of noise to the target policy. It can be seen in (3) that the target policy,  $\pi_{\theta_{targ}}$ , is required to generate an action to evaluate the target Q-function. Noise is added to the policy, taking the form:



$$a'(s') = \text{clip} (\pi_{\theta_{targ}}(s') + \text{clip}(\epsilon, -c, c), a_{low}, a_{high}), \quad \epsilon \sim \mathcal{N}(0, \sigma) \quad (5)$$

where  $\epsilon$  represents the noise sampled using some user specified method. This method of adding noise was shown to reduce the issue of the Q-function approximator developing large values for certain state/action pairs, therefore smoothing ~~the~~ the Q-function.

The second trick the TD3 algorithm employs is the addition of a second Q-function approximator and target Q-function approximator. A known potential issue

of the Q-function is that it can suffer from overestimation of the value of action/state pairs. This, of course, leads to the policy learning actions that the Q-function assumes are better than they actually are. To alleviate this issue, the authors suggest instantiating two Q-functions and two target Q-functions! Calculating the target Q-function is then completed by evaluating the two target Q-functions:

$$y(r, s', d) = r(s, a) + \gamma (1 - d) Q_{\phi_i, \text{targ}}(s', \pi_{\theta_{\text{targ}}}(s')), \quad \text{for } i = 1, 2 \quad (6)$$

and taking the lower of the two targets to update both the main Q-functions:

$$L(\phi_1, \mathcal{D}) = \underset{(s, a, r, s', d) \sim D}{E} \left[ (Q_{\phi_1}(s, a) - y(r, s', d))^2 \right] \quad (7)$$

$$L(\phi_2, \mathcal{D}) = \underset{(s, a, r, s', d) \sim D}{E} \left[ (Q_{\phi_2}(s, a) - y(r, s', d))^2 \right] \quad (8)$$

The last trick that the TD3 algorithm employs is the addition of a delay between the update of the Q-functions and the policy. It has been shown that in doing this the Q-function is able to converge to a better solution before updating the policy. Ultimately, the addition of a policy update delay was done to reduce coupling between the Q-function and the policy. The recommended delay is updating the policy every two Q-function updates.

There are many implementations of the TD3 algorithm that are available. Of these, the StableBaselines3 implementation is used to complete the work in this thesis [29]. StableBaselines3 is a widely used library of RL algorithm implementations and is composed of well written and understandable documentation for the supported implementations. The differentiable function approximators used to estimate the policies and Q-functions are built within StableBaselines3 using PyTorch [30], which is also a widely used framework for machine learning and more specifically reinforcement learning.

## 1.6 Contributions

The purpose of the work presented in the remainder of the document is to propose and evaluate the performance of a concurrent design architecture that utilizes RL techniques for flexible jumping systems. A concurrent design system in this work ~~(X)~~ is one that concurrently learns a mechanical design for a system and an associated control policy for said system. There is a void in the literature surrounding RL-based concurrent design, particularly regarding locomotive robotics applications. It is of interest in this work to evaluate if a technique can be developed to generate better performing systems than ones that implement only controller policy learning. The method will be evaluated in simulation on a simplified flexible-legged jumping monopode. The components which build the concurrent design architecture will be split across the next few chapters wherein additional findings will be presented.

*Cultured*

In the next chapter, a RL-based controller will be trained on the monopode jumping system to evaluate the effectiveness of training for efficient control. Power use is often considered when designing RL controllers for rigid systems, typically taking the form of a weighted negative reward when deploying an RL algorithm. It is of interest to evaluate if defining strategies for flexible systems, where efficiency is the primary objective, if the resulting control strategy takes advantage of system flexibility. To determine if the learned policies are approaching what current literature supports regarding optimal control, the performance will be evaluated against input shaping techniques in Chapter 3.

*try to use this planning less throughout*

Additionally, it is of interest to incorporate mechanical design into the controller learning process. Therefore, in Chapter 4, a RL problem is defined in a unique way such that the environment the RL policy is deployed in is a simulation of an environment where the actions sent by the policy are mechanical design updates. The method will be evaluated on the monopode jumping system to learn mechanical design parameters related to flexibility. Furthermore, using a single fixed control input, it is of

interest to determine if this technique can be used to define designs to accomplish multiple tasks. Therefore, using a single control input generated from the above input shaping techniques, it will be evaluated if an RL learning technique can be used to learn designs that cause the monopode to jump to multiple heights.

Next, in Chapter 5, the methods of learning a control policy and a design will be combined to create a concurrent design architecture. Two methods of implementing the mechanical design update will be shown, and the effects of the two methods will be discussed. Furthermore, a newly introduced hyperparameter when implementing the constructed concurrent design will be evaluated and the results will be discussed. The methods will be tested to evaluate efficient jumping ~~vs~~ non-efficient jumping concurrent designs for the monopode jumping system. Ultimately, the resulting concurrent design performances ~~X~~ will be presented and ~~D~~ compared to the performances ~~X~~ of control policies trained on static designs.

Lastly, in Chapter 6, the work presented in this document will be ~~concluded~~ and the results of the proposed concurrent design process will be highlighted. Accompanying the conclusive remarks, future research that has been enabled by the work presented will be discussed. ~~and recommendations for next steps given.~~

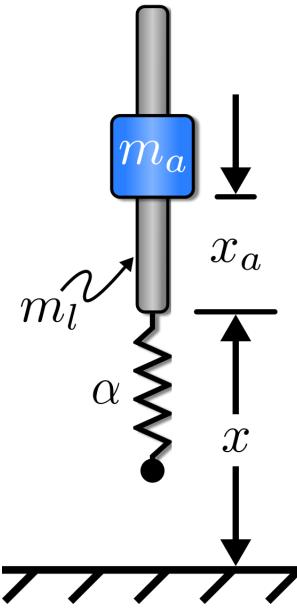
## II Learning Efficient Jumping Strategies for the Monopode System

Utilizing reinforcement learning to train a neural network based controller has been shown to be useful for controlling many robotic systems [22, 23]. It has been used to successfully control rigid-legged robots both in simulation and on physical hardware [21, 31]. Reinforcement learning has been shown to be capable of defining more effective and efficient-jumping techniques for a single-legged robot with SEAs [32]. It has also been shown to be an effective method for controlling multi-legged robots both in simulation and on physical hardware [33–35]. Furthermore, it has been shown to be useful for defining energy efficient strategies for multi-legged robots that have been deployed on physical hardware [36]. However, the body of work demonstrating the use of RL to train controllers for flexible systems is limited, particularly in regards to legged locomotive systems. In this chapter, RL is deployed to define an energy efficient-jumping strategy for a monopode jumping system. The purpose of this work is to validate the use of RL for defining the control aspect of a concurrent design architecture for flexible jumping systems.

### 2.1 Monopode Jumping System

To evaluate the methods discussed in this chapter, a monopode system like the one shown in Figure 6 was used to represent a flexible jumping system. This system has been studied and has been ~~proven~~ <sup>shown</sup> to be an effective base for modeling the jumping gaits for many different animals [37].

The monopode is controlled by accelerating the actuator mass,  $m_a$ , along the rod mass,  $m_l$ , causing a hopping-like motion. The system contacts the ground through a nonlinear spring, represented by the variable  $\alpha$  in the figure. Also included in the model is a damper parallel with the spring, having a damping coefficient of  $c$ , though it is not shown in the figure. Variables  $x$  and  $x_a$  represent the rod's global position and the actuator's local position with respect to the rod, respectively. The equation of



**Figure 6.** Monopode Jumping System

**Table 1.** Monopode Model Parameters

Model Parameter	Value
Mass of Leg, $m_l$	0.175 kg
Mass of Actuator, $m_a$	1.003 kg
Spring Constant, $\alpha_{nominal}$	5760 N/m
Natural Frequency, $\omega_n$	$\sqrt{\frac{\alpha}{m_l+m_a}}$
Damping Ratio, $\zeta_{nominal}$	1e-2 $\frac{N}{m/s}$
Gravity, $g$	9.81 m/s <sup>2</sup>
Actuator Stroke, $(x_a)_{max}$	0.008 m
Max. Actuator Velocity, $(\dot{x}_a)_{max}$	1.0 m/s
Max. Actuator Acceleration, $(\ddot{x}_a)_{max}$	10.0 m/s <sup>2</sup>

Move to next page if no reference  
stay there

motion for the system is:

$$\ddot{x} = \frac{\gamma}{m_t} (\alpha x + \beta x^3 + c \dot{x}) - \frac{m_a}{m_t} \ddot{x}_a - g \quad (9)$$

where  $x$  and  $\dot{x}$  are position and velocity of the rod, respectively, the acceleration of the actuator,  $\ddot{x}_a$ , is the control input, and  $m_t$  is the mass of the complete system.

Constants  $\alpha$  and  $c$  represent the linear spring and damping coefficient, respectively, and constant  $\beta$  is set to 1e8. Ground contact determines the value of  $\gamma$ , so that the spring

mass  $m_a$

*apply only*  
and damper do not supply force while the leg is airborne:

$$\gamma = \begin{cases} -1, & x \leq 0 \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

Additionally, the spring compression limit, or the systems position in the negative  $x$  direction, is limited to 0.008m. The system is also confined to move only vertically ~~in regards to Figure 6~~ so that controlling balance is not required. The values of the parameters shown in Figure 6 are displayed in Table 1.

## 2.2 Training Environment

Using the monopode model, an environment aligning with the standards set by OpenAI for a Gym environment was created [38]. The observation and action spaces were defined, respectively, as follows:

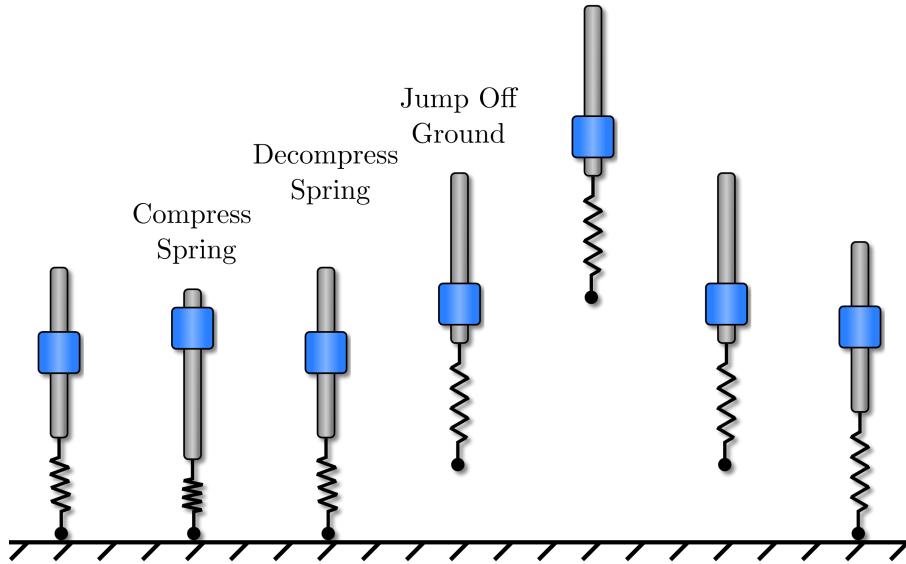
$$\mathcal{S} = [x_{a_t}, \dot{x}_{a_t}, x_t, \dot{x}_t] \quad (11)$$

$$\mathcal{A} = [\ddot{x}_{a_t}] \quad (12)$$

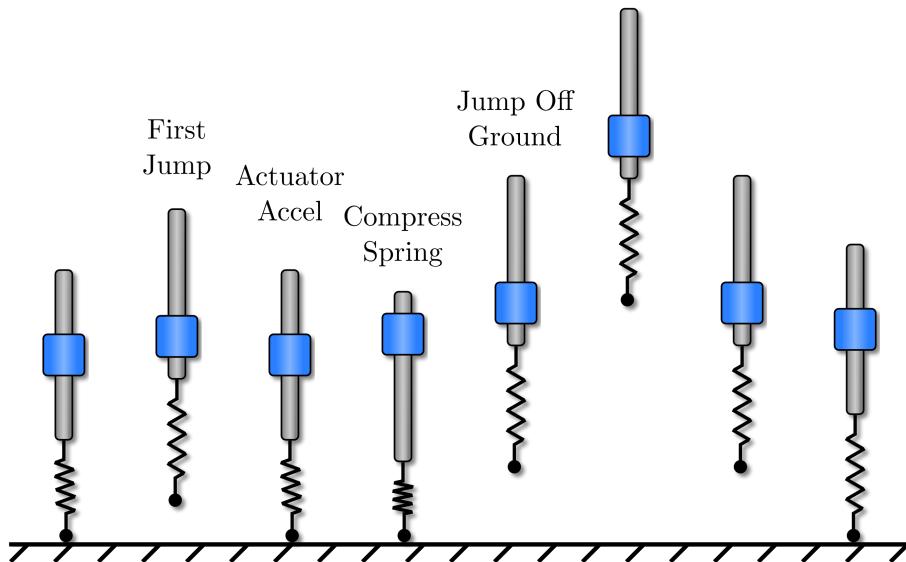
where  $x_t$  and  $\dot{x}_t$  were the monopode's position and velocity at time  $t$ , and  $x_{a_t}$ ,  $\dot{x}_{a_t}$  and  $\ddot{x}_{a_t}$  were the actuator's position, velocity and acceleration, respectively.

Two separate stopping conditions were defined for the environment to evaluate two different jump types and therefore two different ~~input~~ <sup>jumping</sup> commands. The first was defined as the monopode's position being greater than zero then returning to zero once. The second was defined like the first but with the monopode's position being greater than zero and then less than zero twice.

Two different jumps were created from these stopping conditions. The first was referred to as a single jump command, and the second a stutter jump command. The intent of utilizing two different jumping commands was to determine if a RL algorithm was more or less effective in learning differing strategies depending on the complexity of the desired command.



(a) Example Single Jump



(b) Example Stutter Jump

**Figure 7.** Jumping Types for the Monopode Jumping System

An example single jump can be seen in Figure 7a. The intended command from the learned controller would be one that would jump the monopode once. This type of command would ideally compress the spring/damper by accelerating the actuator in the positive direction. This would cause the rod mass to accelerate downward compressing the spring to store energy that could be used to cause the system to jump. The

actuator mass should then accelerate downward forcing the spring to decompress. At this point, the system should be accelerating upwards and the actuator downwards such that the monopode leaves the ground completing a single jump.

An example stutter jump can ~~be seen~~ <sup>is shown</sup> in Figure 7b. The intended command from the learned controller would be one that would jump the monopode twice. This type of command would firstly complete an optimal single jump. Following that motion, the actuator should ~~assume an acceleration direction~~ <sup>accelerate</sup> to recompress the spring, storing more energy with a farther compression. When the spring is compressed to its maximum value or the system's total acceleration reaches zero, the actuator mass should accelerate downwards forcing the spring to decompress. At this point, the system should be accelerating upwards and the actuator downwards, similar to the single jump, such that the monopode leaves the ground completing a stutter jump.

### 2.3 Efficient Control Strategies

Efficient control of a robotic system is often one of the most important aspects of a controller's design. Applications where a robotic system is deployed and relies on a limited power source, such as a mobile walking robot, will often require an efficient control strategy. Modern, traditional methods, such as model predictive control, have been shown to produce energy efficient locomotion strategies for wheeled and legged systems [39, 40]. It is of interest in this work to utilize RL, a modern neural network based control method, to find strategies which are designed with power efficiency as the primary objective.

Two different reward functions were designed to accomplish the task of determining how well RL learns efficient-jumping strategies. The purpose of defining two different reward functions was to compare the ~~input~~ <sup>performance</sup> commands and resulting jumping shapes of the two controller types to determine if the efficient controller was learning to conserve power.

The first reward function was one that ignored power usage and focused solely on the height of the jump:

$$R = x_t \quad (13)$$

where  $x_t$  was the height of the monopode system at any given time step. The second reward function was one that was defined to accomplish the same task, but also consider power consumption. It was defined as:

$$R = \frac{x_t}{\sum_{t=0}^t P_t} \quad (14)$$

where  $P_t$  was the power consumption of the monopode system at any given time step defined mechanically as the product of the actuator's acceleration, velocity and mass:

$$P_t = m_a \dot{x}_a \ddot{x}_a \quad (15)$$

where  $m_a$  was the mass of the actuator, and  $\dot{x}_a$  and  $\ddot{x}_a$  where the actuators velocity and acceleration, respectively.

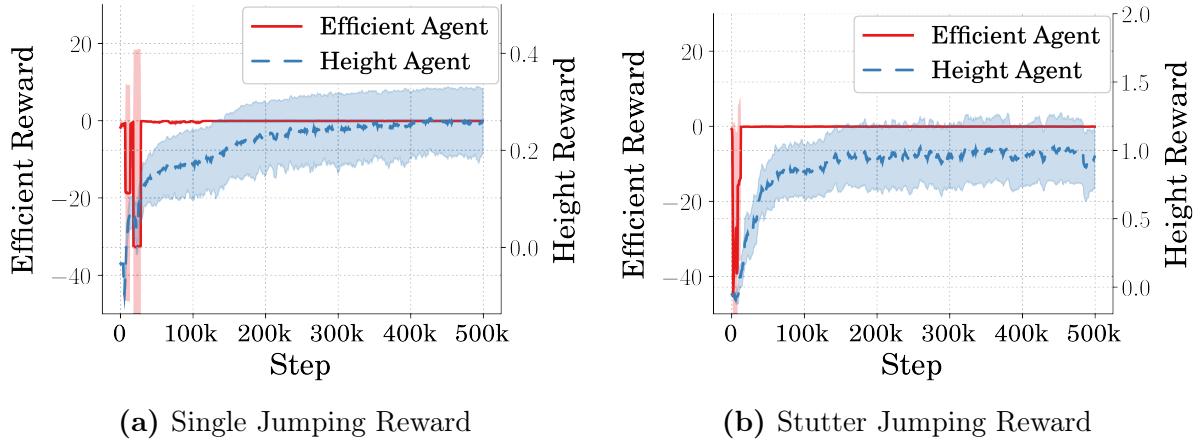
## 2.4 Deploying TD3

Because RL may generate a locally optimal controller instead of a globally optimal controller, training more than one controller is common practice for evaluating performance. In this work, fifty different controllers where trained, each with a different random network initialization. Each controller was trained for a total 500k time steps. The remaining hyperparameters set using the TD3 algorithm are defined in Table 2.

The average and standard deviation of the rewards during training for the efficient and high jumping strategies for both the single and stutter jumping commands are shown in Figure 8. They represent the controllers being trained to accomplish their respective goals. Looking at Figure 8a, which shows the rewards for learning a single jumping command, it is clear that there are definite differences between the efficient and high-jumping reward types. Firstly, the high-jumping strategy does converge after 500k steps of training. The reward for the efficient-jumping controller, having been

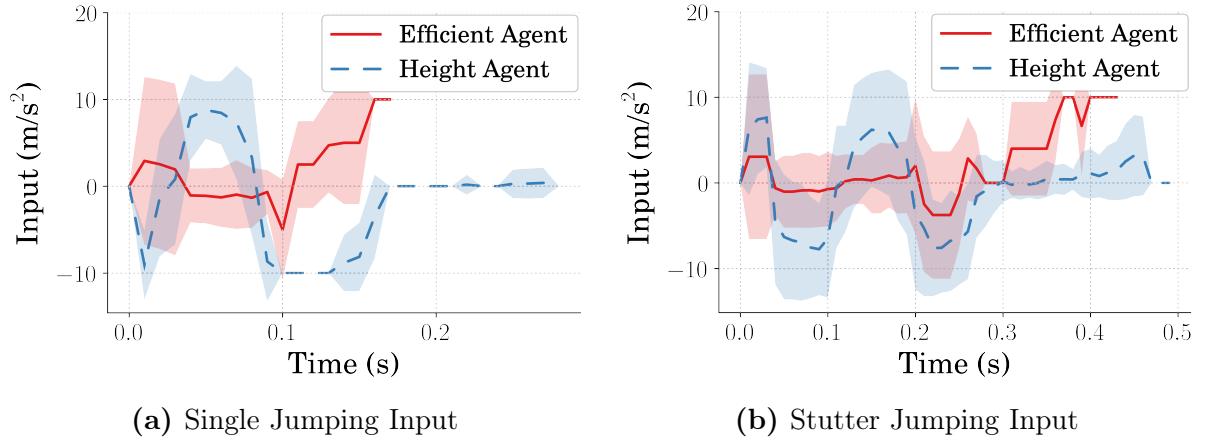
**Table 2.** TD3 Training Hyperparameters

Hyperparameter	Value
Learning Rate, $\alpha$	0.001
Learning Starts	1000 Steps
Batch Size	100 Transitions
Tau, $\tau$	0.005
Gamma, $\gamma$	0.99
Training Frequency	1:Episode
Gradient Steps	$\propto$ Training Frequency
Action Noise, $\epsilon$	None
Policy Delay	1 : 2 Q-Function Updates
Target Policy Noise, $\epsilon$	0.2
Target Policy Clip, $c$	0.5
Seed	50 Random Seeds



**Figure 8.** Reward vs Time Step During Training

defined drastically different than the reward from the high-jumping controller, not surprisingly, looks drastically different than the high-jumping reward. The reward for the efficient agent is heavily punished for using power without gaining height, which is a frequent occurrence in the beginning of training. This forces the policy to learn that using less power will result in higher rewards within very few learning steps. In Figure 8b, it is also apparent that the height controller ~~is~~ converging to a solution. Additionally, the efficient controller can be seen to have learned in the same rapid form as the single jumping command type.



**Figure 9.** Average and Standard Deviation Inputs to Monopode

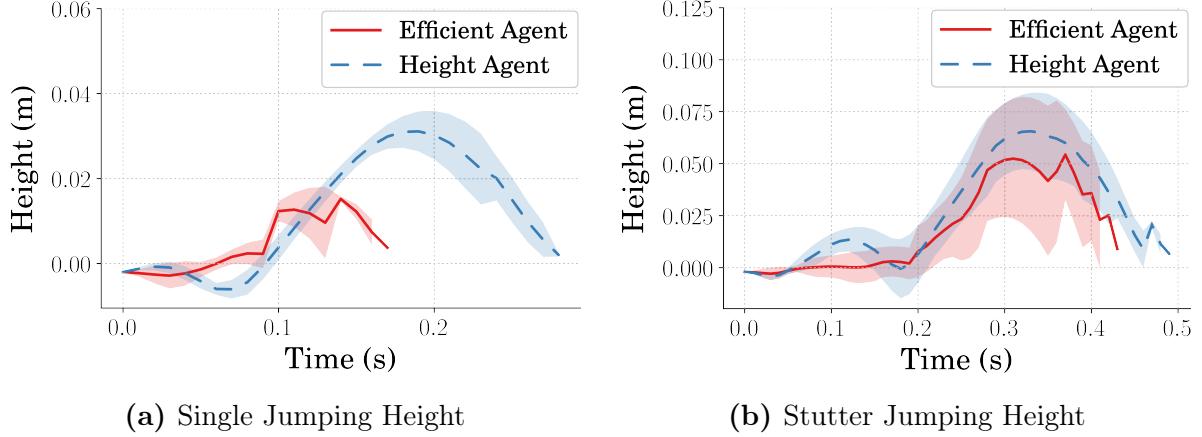
## 2.5 Average Performance of Network Controller

**2.5.1 Input Commands.** The average and standard deviation of the final controller's ~~input~~<sup>Jumping</sup> commands were evaluated for both the single and stutter jumping cases, and the results are shown in Figure 9. The individual subfigures are to compare the controllers that were trained to jump efficiently to those trained to jump high. At first glance, there are obvious differences regarding timing, magnitude, and direction. There are also slight differences in variance seen between the two controller types.

Starting with Figure 9a, which displays the commands for the single jumping case, it is most obvious that the direction for the initial acceleration of the actuator mass differs between the efficient controllers and height controllers. In the case where the controller is learning to jump high, an initial acceleration of the actuator mass in the negative direction is learned, which contrasts the case where the controller is learning an efficient command. Further, the magnitude of the commands is drastically different which may be an indicator for conserving power. For the stutter jumping case shown in Figure 9b, it is immediately apparent that the magnitudes of the commands differ greatly. They are, however, more similar in regards to their timings and directions than the single jumping command.

How could you determine if this is true?

In both the single and stutter jumping cases, it can be seen that there is upward



**Figure 10.** Average and Standard Deviation Heights of Monopode

acceleration command towards the end of the jump, which again, might be an indicator of a more efficient-jumping strategy. Furthermore, it can be observed that in the single jumping case, there exists more variance across different instances of the trained efficient controllers in comparison to the height controllers. This does not seem to be the case for the stutter jumping command type, though both cases do seem to generate controllers with high variance inputs across instances.

**2.5.2 Jumping Height Performance.** The average and standard deviation of the final controller’s jumping performance were evaluated for both the single and stutter jumping cases, and the results are shown in Figure 10. In both the single and stutter jumping cases, there are differences in jumping ability when comparing the efficient and height controller types. It is apparent that when increasing the complexity of the command from a single jump to a stutter jump, the efficient controllers are better able to match the performance of the height controllers.

Shown in Figure 10a, the height controllers for the single jumping case learned a command ~~input~~ that outperformed the efficient controllers in terms of jump height. The resulting motion from the input discussed in the previous section can be seen in that the efficient controller learned to simply compress the spring, then jump the monopode. The height controllers, in contrast, disregarding power consumption,

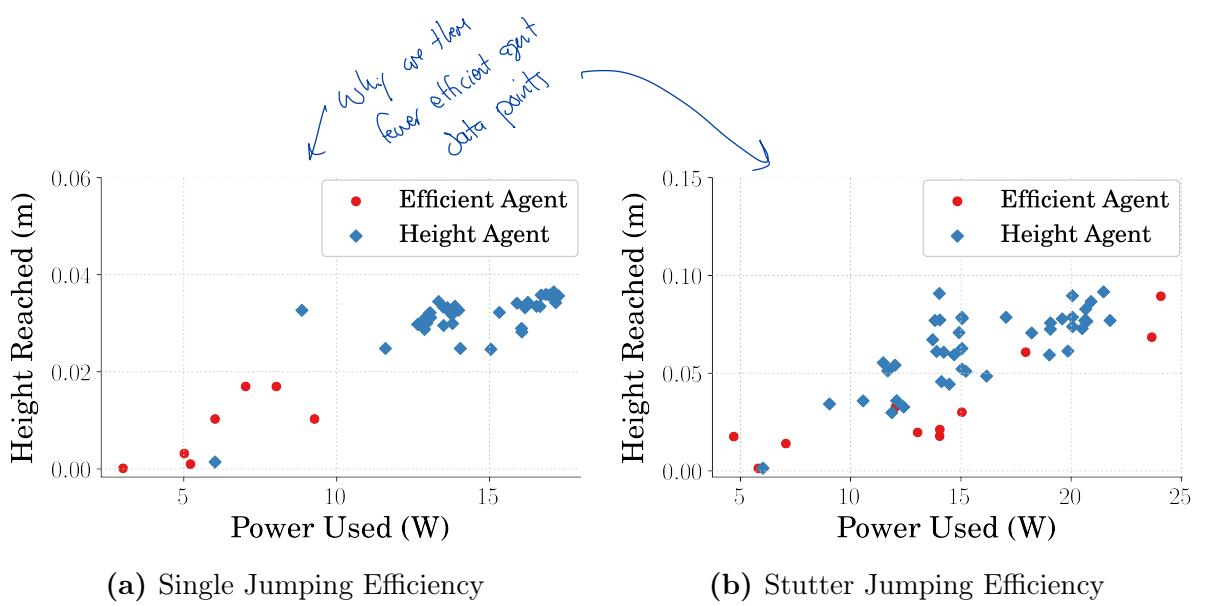
learned to decompress the spring from its nominal position, keeping it below the point of leaving the ground, then recompressing for a much higher jump. For the single jumping commands, the high-jumping strategy trained controllers that jumped the monopode 104.53% higher than the efficient-jumping strategy ~~on average~~.

Figure 10b compares the jumping performance of the efficient and height strategies for the stutter jumping command. They differ, but less drastically than the single jumping case. The similarity of the command shapes shown in Figure 9b ~~(X)~~ results in jumping responses that also share a similar form. The large differences seen in the stutter jumping performances ~~(X)~~ are similar to those seen in the inputs, in that they differ mostly regarding the magnitudes. For the stutter jumping command, the high-jumping strategy trained controllers that jumped the monopode 20.69% higher than the efficient-jumping strategy.

In both the single and stutter jumping cases, the upward acceleration from efficient controllers toward the end of the command can be seen in that the monopode regains height after the start of its final decent from maximum height. This can be explained in that the efficient control method, which punishes power consumption, discovered a way to maintain height throughout a jump where the utilization of additional power is less costly. It is less costly to influence the monopode's position while airborne because there is no resistance from the spring and damper. Additionally, regarding variance, the single jumping controllers seem to produce jumping shapes with similar levels of variance. Whereas in the stutter jumping case, though similar in high levels of input variance, the jumping height variance for the efficient controllers is noticeably higher than that of the height controller.



**2.5.3 Height Reached vs. Power Used.** Height reached versus power consumed data for both the single jumping and stutter jumping cases, is shown in Figure 11. In both the single and stutter jumping cases, the efficient controllers utilized

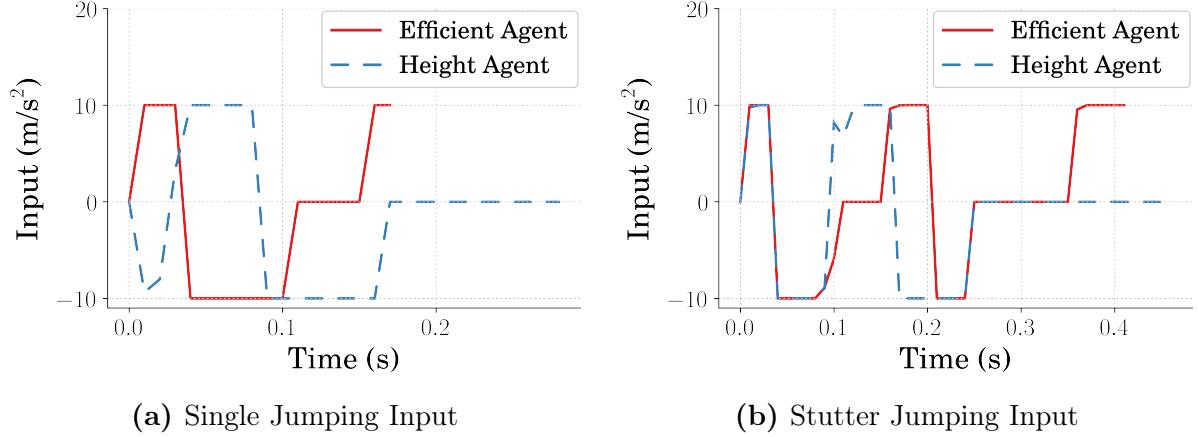


**Figure 11.** Height Reached vs Power Consumed of Monopode

less power and therefore suffered regarding jump height. This matches what was seen in the previous sections regarding input commands and jumping shapes. It ~~is ultimately observed, that~~ in both jumping cases, the power conservation gain is greater than the jumping height cost.

In the single jumping case, shown in Figure 11a, there is an apparent separation between the two controller types where the height controllers learn control strategies that use more power and jump higher. On average, for the single jumping command type, the efficient-jumping strategy learned jumping commands that were 126.27% more efficient at the cost of 104.53% in average jump height.

As for the stutter jumping case, shown in Figure 11b, the difference in performance is less obvious. This can be explained in that more complex jumping strategies give an RL algorithm more opportunity to learn a control policy that can better take advantage of system flexibility. The variance of the two controller types, being quite high, matches what is seen in the previous sections and results in more mixing of the data. On average, for the stutter jumping command type, the efficient strategy learned commands that were 101.45% more efficient with the average maximum jumping height only being punished by 20.69%.



**Figure 12.** Optimal Inputs to Monopode

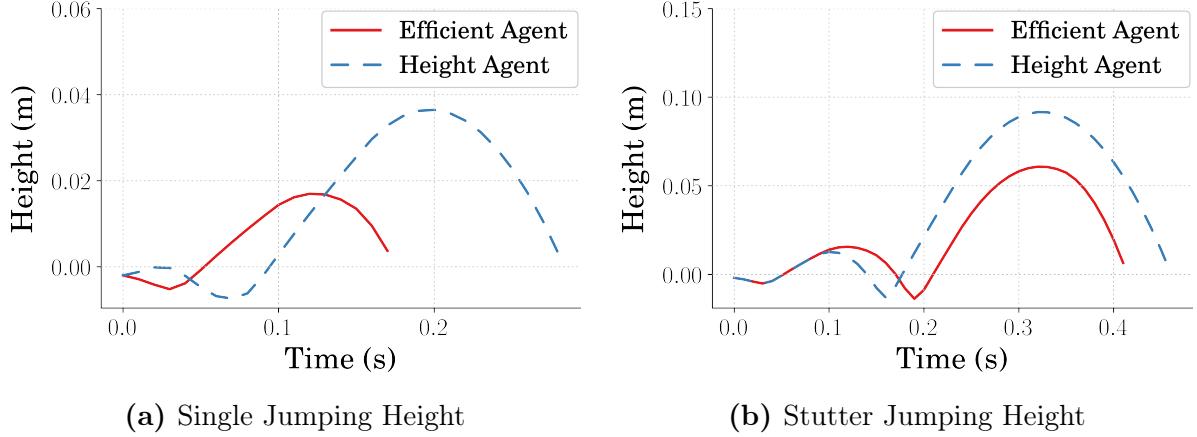
## 2.6 Optimal Performance of the Proposed Controller

**2.6.1 Jumping Input Commands.** Taking the best of the fifty different controllers trained for both the single and stutter jumping cases and comparing the efficient and height controller's performance can show what is possible with a properly defined RL problem. Figure 12 shows the differences in the input commands generated when selecting the highest performing controller in terms of reward received. It can be seen that there are less differences between the the efficient and height controllers in comparison to the average results from Section 2.5. A major similarity is that magnitudes are similar across all cases such that the controllers utilize the actuator's maximum acceleration.

Looking at Figure 12a, which compares the efficient and height controllers for the single jumping input, the major differences are the timing and direction of the commands. This is similar to the average performance evaluation, from Section 2.5.1, in that the efficient controller does not take advantage of the slight decompression of the spring before the monopode leaves the ground. Because of this, the efficient controller learns a different timing for a single jump.

As for the stutter jumping case, shown in Figure 12b, the differences between the efficient and height controller ~~are~~ is less drastic. The initial timing is largely the same

*Does this really  
what you'd  
expect, for a  
traditional control  
perspective?*



**Figure 13.** Optimal Heights of Monopode

as both controller types learn to utilize the decompression of the spring. The differences begin when decompressing the spring a second time and completing the first jump. The efficient controller learns a command similar in form to a bang-coast-bang command, where, in contrast, the height controller learns a command similar to that of a bang-bang shaped input.

**2.6.2 Jumping Height Performance.** In line with the previous section, it is of interest to evaluate the jumping performance curves of the best controllers trained. Figure 13 displays the jumping performance for both jump types as well as both controller types when utilizing the inputs shown in Section 2.6.1. These curves show that the efficient controllers, in the best case scenario, do not generate commands that jump the monopode as high as the high jumping controllers.

Figure 13a, which compares the efficient and height controllers for the single jump, shows that the efficient controller does not learn to utilize the allowable decompression in the spring. In Figure 13b, it can be seen that when utilizing a command more similar in form to a bang-coast-bang command, like the efficient controller learned, the timing of the jump sequence is shifted and the resulting final height is less than the height controller who's command is more similar to a bang-bang shaped command. The efficient controller having learned to coast between commands

shows it is a useful method for conserving power, but not a good strategy for optimizing jumping height.

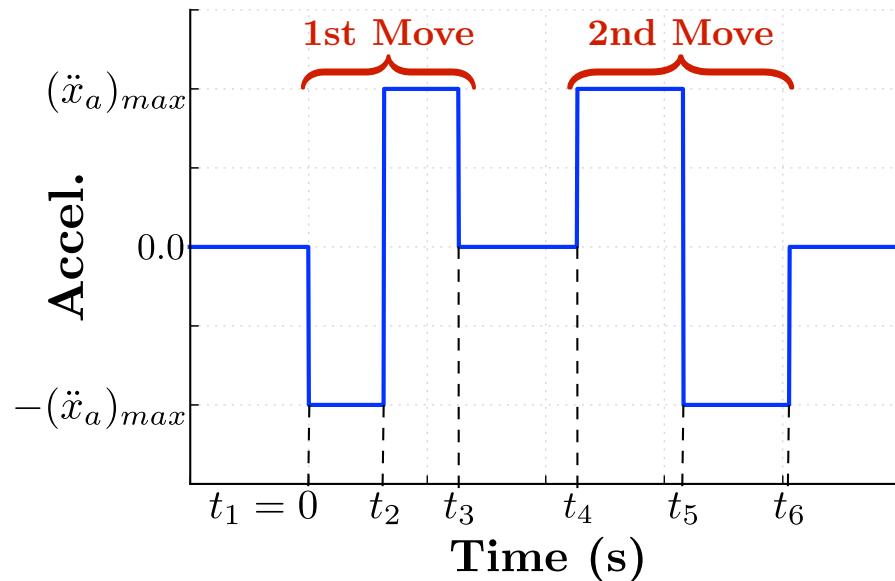
## 2.7 Conclusion

Two different controller types were trained to generate two different jumping commands for the simplified monopode jumping system. The first type of controller was one that would command the monopode system to jump high where the reward was based on nothing other than system height. The second type of controller was one which controlled the monopode to jump high but at the cost of power consumed, such that high jumps that consumed high amount of power were less desirable than high jumps that consumed less power. It was shown that the rewards passed to RL algorithms that are training controllers can be manipulated so that the learned input commands take advantage of the spring/damper that exists within the monopode jumping system. Furthermore, the timing of the commands, as well as the input magnitude and direction are all affected when defining a reward strategy that seeks to increase power efficiency. When considering the average performance of the different control strategies, for both the single and stutter jumping cases, the heights reached were less for the efficient strategies. However, they were significantly more efficient, particularly when scaling the complexity of the command from the single jump to the stutter jump. It should be concluded that RL might serve as a useful method for defining control strategies for flexible-legged jumping systems, particularly when energy efficiency is of interest. Additionally, when considering more complex control strategies, which might be difficult to define efficiently, RL might serve as a useful method for effective efficient strategies.

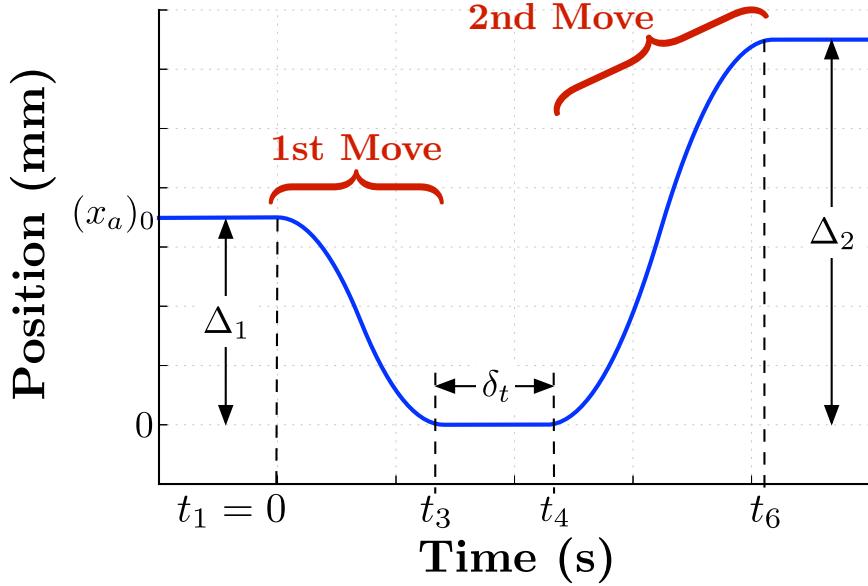
*Edits to this chapter coming soon.*

### III Using Input Shaping to Validate RL Controller

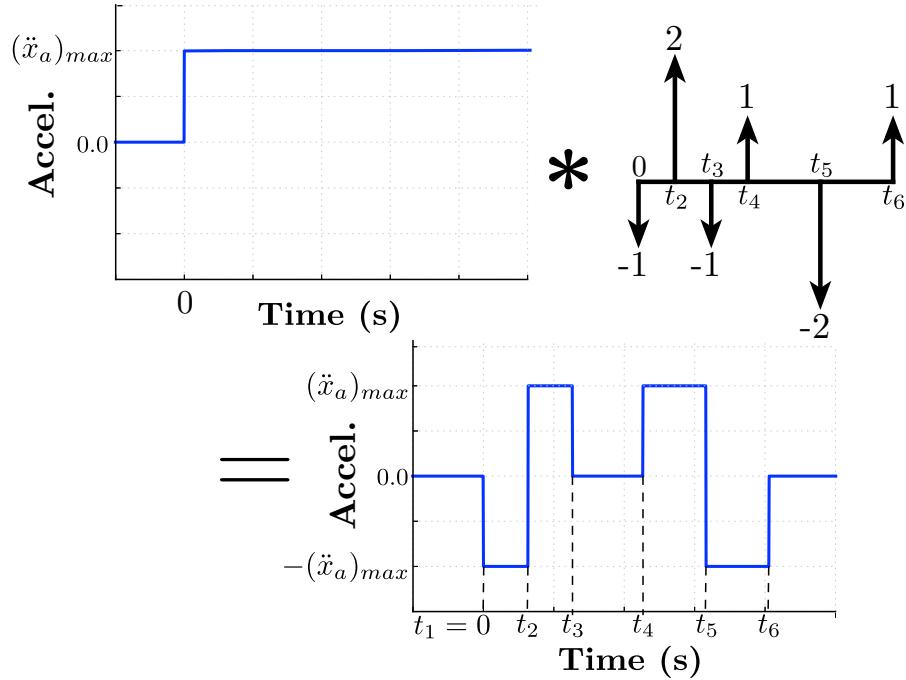
In utilizing RL to define a control strategy for a robotic system, the resulting commands sent to the system are often described as optimal, or at least approaching optimal. They are described as such due to the nature of RL problems in that the techniques used to learn a policy are optimization theory based so the policy being trained is one that is approaching an optimal solution in regards to the reward defined. Interpreting the commands claimed to be optimal, in the context of control theory, is an important part of utilizing RL for defining control strategies for robotic systems however. In the case of a flexible jumping robot, where the command is to jump the system as high as possible or as efficient as possible, are the commands generated truly approaching an optimal solution? Methods such as input shaping, having been shown to be an effective method for defining optimal control strategies for flexible jumping systems, can be used to evaluate the commands generated by the RL generated policies [2].



**Figure 14.** Jumping Command [2]



**Figure 15.** Resulting Actuator Motion [2]



**Figure 16.** Decomposition of the Jump Command into a Step Convolved with an Impulse Sequence [2]

### 3.1 Input Shaping Controller Input

Bang-bang based jumping commands like the one shown in Figure 14 are likely to result in a maximized jump height [2]. For these command types, regarding the

monopode jumping system, the actuator mass travels at maximum acceleration within its allowable range, pauses, then accelerates in the opposite direction. Commands designed to complete this motion are bang-bang in each direction, with a selectable delay between them. The resulting motion of the actuator along was rod is shown in Figure 15. Starting from an initial position,  $x_{a_0}$ , the actuator moves through a motion of stroke length  $\Delta_1$ , pauses there for  $\delta_t$ , then moves a distance  $\Delta_2$  during the second portion of the acceleration input.

This bang-bang-based profile can be represented as a step command convolved with a series of impulses, as was shown in Figure 16 [41]. Using this decomposition, input-shaping principles and tools can be used to design the impulse sequence [42, 43]. For the bang-bang-based jumping command, the amplitudes of the resulting impulse sequence are fixed,  $A_i = [-1, 2, -1, 1, -2, 1]$ . The impulse times,  $t_i$ , can be varied and optimal selection of them can lead to a maximized jump height of the monopode system [2]. Commands of this form will often result in a stutter jump like what was shown in Figure 7b of Chapter 2, where the small initial jump allows the system to compress the spring to store energy to be used in the final jump.

### 3.2 RL Controller Input

Discussion and figures from the inputs defined by the RL algorithms for the monopode system.

Waiting for some data from DV.

### 3.3 Conclusion

Discuss the results.

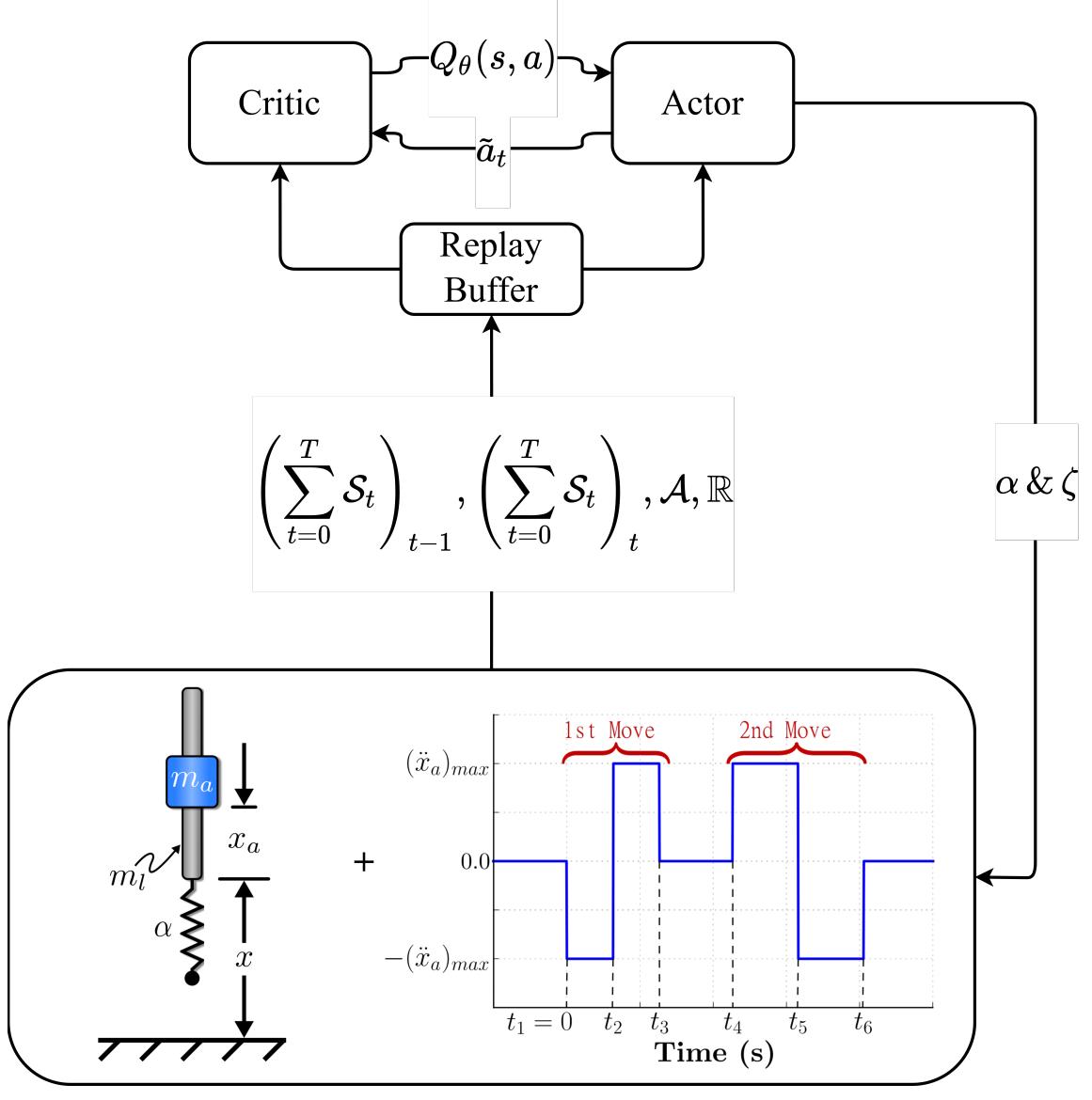
## IV Mechanical Design of a the Monopode Jumping System

Often it is the goal of a controls engineer to design a controller to accommodate and manipulate systems according to the system description provided. However, research has been conducted showing the value of studying the manipulation of mechanical design parameters in order to achieve a desired system behavior [17]. In this chapter, reinforcement learning is shown to be useful as a tool to learn mechanical designs given a predefined system controller for the monopode jumping system. RL has been shown to be an effective strategy for finding optimal concurrent designs for many different types of systems [19, 44, 45]. It has even shown *it's* ability to define designs that are successfully deployed on physical hardware [18]. It is comparable to work where evolutionary algorithms are deployed to optimize physical parameters of systems for improved energy efficiency [46, 47]. Here, RL is used to define an optimal design for the monopode jumping system described in Chapter 2.

### 4.1 Learning a Mechanical Design

Section 1.4 describes the common deployment methodology of an RL problem where defining a control policy for a robotic system is often the primary goal. In this chapter, rather than finding a control policy for a defined robotic system, RL is deployed to find a mechanical design for a defined control input. To do this, the general methods in setting up the problem have to change. Figure 17 shows the general flow of information for the algorithm.

The RL problem is similar to the common use case in that the algorithm utilizes the same information types to optimize a design, such as the state of the environment and the reward. The algorithm's interaction with the environment is drastically different, however. The action space, instead of being a command type input, is a range of design choices for a set of parameters within a simulation of a system. Applying these actions within the environment, rather than changing the state of a robotic



**Figure 17.** Learning a Mechanical Design

system from state  $s$  to  $s'$ , instead simulates a system within the environment from time  $t = 0$  to time  $T$  using a predefined control input. Because of this, the state saved as a transition is a matrix of states rather than a single vector state. The reward does not differ greatly in that it is based on the state of the system. A difference being that the information stored in a single state transition is much greater so the reward can be defined to utilize all the information.

For the work presented in this chapter, as is shown in Figure 17, the predefined control input used to simulate the monopode system within the environment at each step was an optimized controller generated using the input shaping techniques discussed in Section 3.1. Having been shown to be a useful technique for generating optimal control strategies, it was of interest to evaluate if an optimal mechanical design could be found to accompany the input.

## 4.2 Environment Definition

To allow the agent to find a mechanical design, a reinforcement learning environment conforming to the OpenAI Gym standard [38] was created. The monopode model described in Chapter 2 was used as the simulation, and the fixed controller input was based on the work described in Section 3.1. The mechanical parameters the agent was tasked with optimizing were the spring constant and the damping ratio of the monopode system. At each episode during training, the agent selected a set of design parameters from a distribution of available designs. The actions applied,  $\mathcal{A}$ , and transitions saved,  $\mathcal{S}$ , from the environment were defined as follows:

$$\mathcal{A} = \{\{a_\alpha \in \mathbb{R} : [-0.9\alpha_{nom}, 0.9\alpha_{nom}]\}, \{a_\zeta \in \mathbb{R} : [-0.9\zeta_{nom}, 0.9\zeta_{nom}]\}\} \quad (16)$$

$$\mathcal{S} = \left\{ \sum_{t=0}^{t_f} x_t, \sum_{t=0}^{t_f} \dot{x}_t, \sum_{t=0}^{t_f} x_{at}, \sum_{t=0}^{t_f} \dot{x}_{at} \right\} \quad (17)$$

where  $\alpha_{nom}$  and  $\zeta_{nom}$  are the nominal spring constant and damping ratio of the monopode, respectively;  $x_t$  and  $\dot{x}_t$  are the monopode's rod height and velocity steps, and  $x_{at}$  and  $\dot{x}_{at}$  are the monopode's actuator position and velocity steps, all captured during simulation. The action space was set to  $\pm 90\%$  of the nominal values as that percentage allowed for a large variance in performance across the range of designs.

### 4.3 Rewards for Learning Designs

The RL algorithm was utilized to find designs for two different reward cases. Time series data was captured during the simulation phase of training and was used to evaluate the designs performance through these rewards. The first reward case used was:

$$R_1 = \left( \sum_{t=0}^{t_f} x_t \right)_{max} \quad (18)$$

where  $x_t$  was the monopode's rod height at each step during simulation. The goal of the first reward was to find a design that would cause the monopode to jump as high as possible.

The reward for the second case was:

$$R_2 = \frac{1}{\frac{|R_1 - x_s|}{x_s} + 1} \quad (19)$$

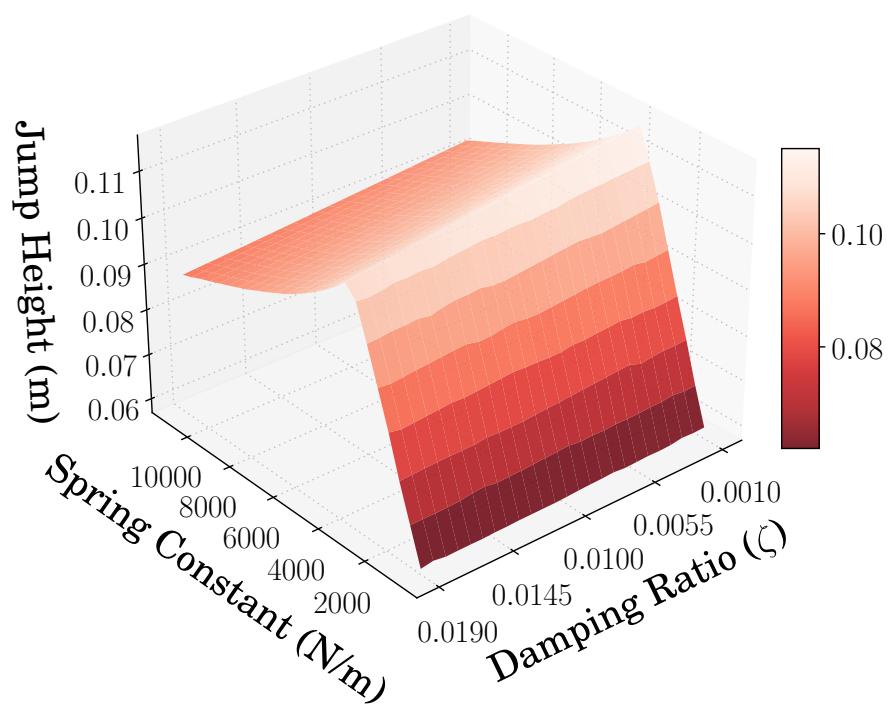
where  $x_s$  was the desired jump height, which was set to 0.01 m. The second case was utilized to test RL's ability to find a design that minimized the error between the maximum height reached and the desired maximum height to reach.

### 4.4 Design Space Variations

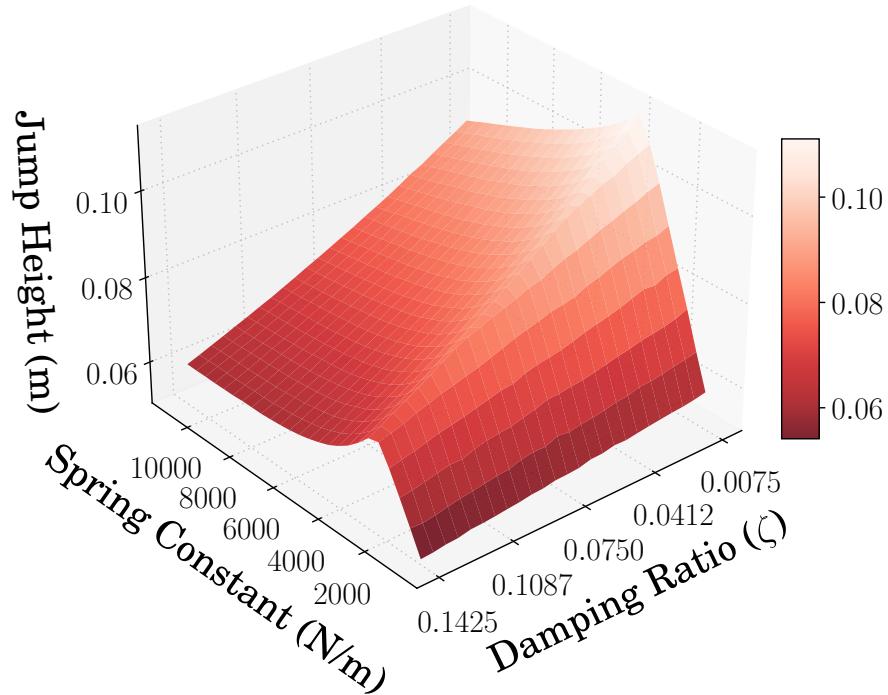
Figure 18 represents the heights the monopode could reach for two different design spaces. The design space provided for the first case, shown in Figure 18a, represents a space where the nominal damping ratio,  $\zeta_{nom}$ , is 0.01, such that  $\pm 0.9 \zeta_{nom}$  creates a more narrow design space. This design space also has more values closer to the optimal value therefore making it more difficult to optimize quickly. The design space provided for the second case, shown in Figure 18b, represents a space where the nominal damping ratio,  $\zeta_{nom}$ , is 0.075, such that  $\pm 0.9 \zeta$  creates a wider design space. Additionally, there is a more obvious maximum within the design space, making it easier to ascend towards an optimal design. In both cases there are many values that would satisfy the specified height jumping strategy.

Why not look at efficiency again?

Still not sure I follow this logic



(a) Jumping Performance of Narrow Design Space



(b) Jumping Performance of Wide Design Space

**Figure 18.** Reference Jumping Performance of the Monopode

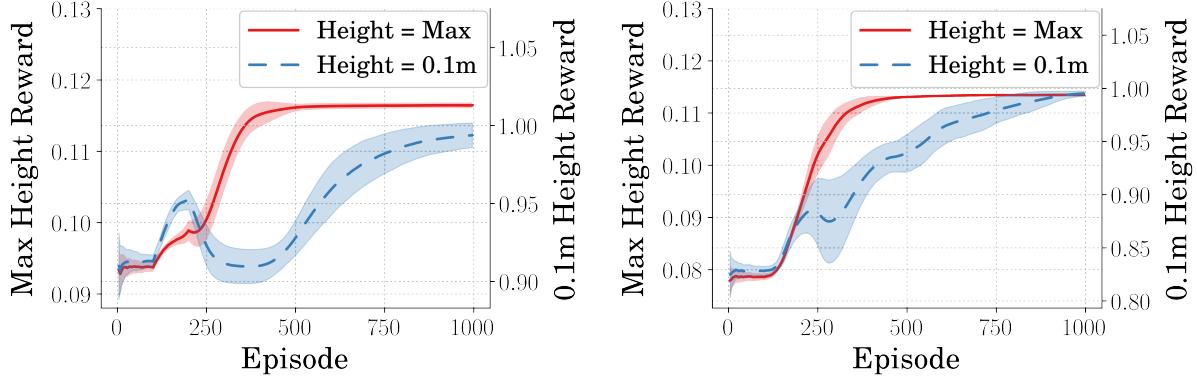
**Table 3.** TD3 Training Hyperparameters

Hyperparameter	Value
Learning Rate, $\alpha$	0.001
Learning Starts	100 Steps
Batch Size	100 Transitions
Tau, $\tau$	0.005
Gamma, $\gamma$	0.99
Training Frequency	1:Episode
Gradient Steps	$\propto$ Training Frequency
Action Noise, $\epsilon$	None
Policy Delay	1 : 2 Q-Function Updates
Target Policy Noise, $\epsilon$	0.2
Target Policy Clip, $c$	0.5
Seed	100 Random Seeds

#### 4.5 Deploying TD3

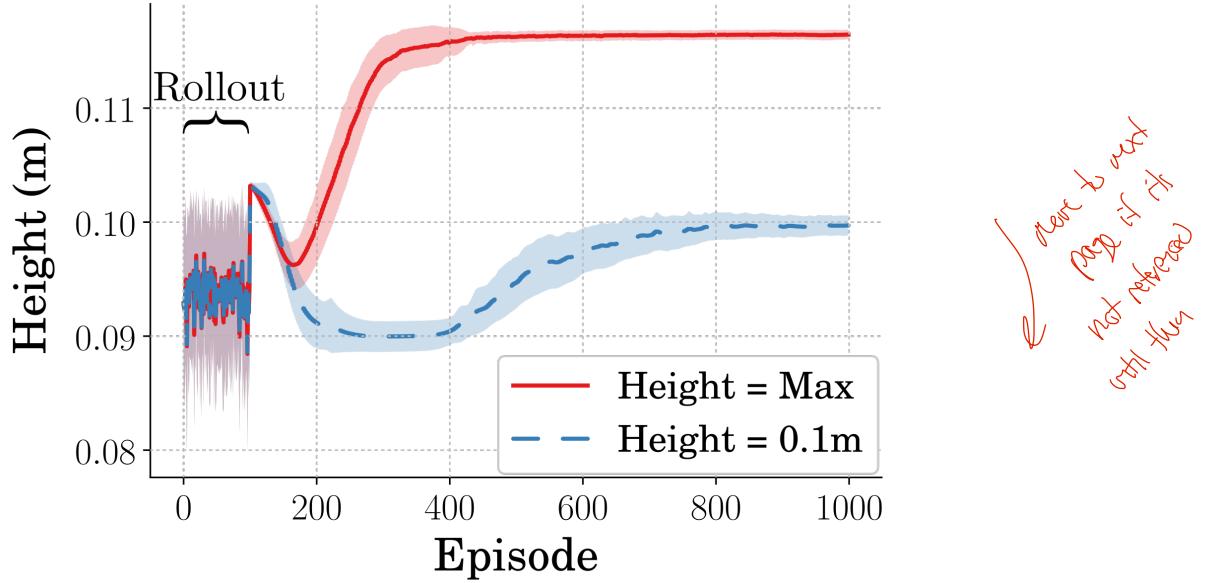
The training hyperparameters were selected based on TD3’s author recommendations and StableBaselines3 [29] experimental findings and are displayed in Table 3. All of the hyperparameters, with the exception of the rollout (Learning Starts) and the replay buffer, were set according to StableBaselines3 standards. The rollout setting was defined such that the agent could search the design space at random, filling the replay buffer with enough experience to prevent the agent from converging to a design space that was not optimal. The replay buffer was sized proportional to the number of training steps due to system memory constraints.

The average rewards for both the narrow and the wide design space agents are shown in Figure 19. They represent the agents learning a converging solution to the problem of finding optimal design parameters. Looking at Figure 19a, it is apparent that given a more narrow design space, both the high and the specified jumping agents were still able to learn a converging solution. ~~It can also be observed that there~~ <sup>also</sup> exists more variance for the specified-height agent type compared to the agents assigned with maximizing jump height. Looking at Figure 19b, it is apparent that the agents given a



(a) Reward vs. Episode: Narrow Design Space    (b) Reward vs. Episode: Wide Design Space

**Figure 19.** Reward vs. Episode for Learning Mechanical Design

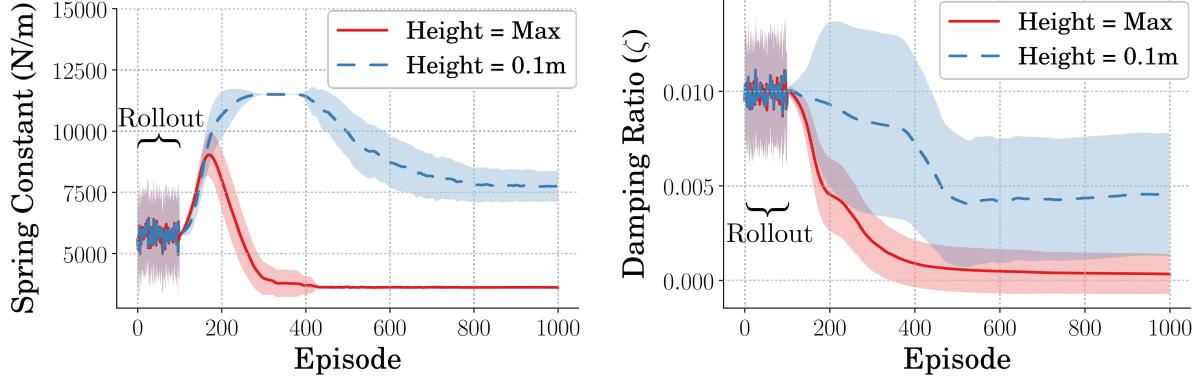


**Figure 20.** Height Reached During Training Given Narrow Design Space

wider design space where both able to learn a converging design solution. Though, given more designs to choose from, it appears the specified height agents are taking longer to converge. Additionally, once converged, because there are less values that allow the controller/design architecture to accomplish the tasks, there is less variance seen in designs learned.

#### 4.6 Jumping Performance

) don't leave  
"new" or  
"old" or  
a page... surprise  
(Lex did this.)



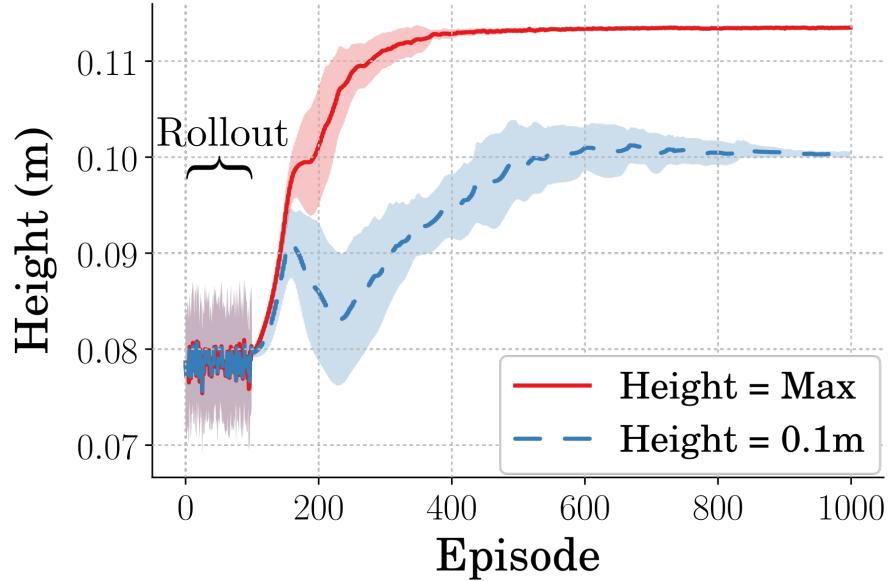
(a) Spring Constant Selected During Training    (b) Damping Ratio Selected During Training

**Figure 21.** Designs Learned for the Narrow Design Space

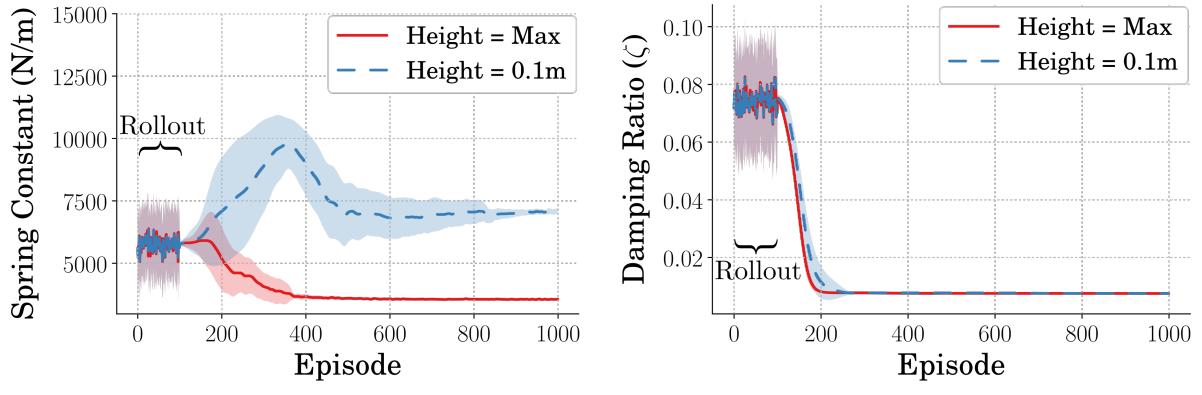
**4.6.1 Narrow Design Space.** Figure 20 shows the height achieved by the learned designs for the agents given the narrow range of possible damping ratio values. For the agents learning designs to maximize jump height, Figure 20 can be compared with Figure 18a showing that the agent learned a design nearing one which would achieve maximum performance. Additionally, looking at the agents learning designs to jump to the specified 0.1 m, the designs learned accomplish this with slightly more variance than that of the maximum height case.

The average and standard deviation of the spring constant and damping ratio design parameters the agents selected during training are shown in Figure 21. These plots represent the learning curves for the agents learning design parameters to maximize jump height and the agents learning design parameters to jump to 0.01 m. There is a high variance in both the spring constant and the damping ratio found for the agents that learned designs to jump to a specified height. The agents which were learning designs which maximized height found designs with very little variance in terms of spring constant and significantly less variances in terms of damping ratio. *less than what?*

**4.6.2 Wide Design Space.** Figure 22 shows the height achieved by the learned designs for the agents given a wide range of damping ratios. For the agents



**Figure 22.** Height Reached During Training Given Wide Design Space



(a) Spring Constant Selected During Training    (b) Damping Ratio Selected During Training

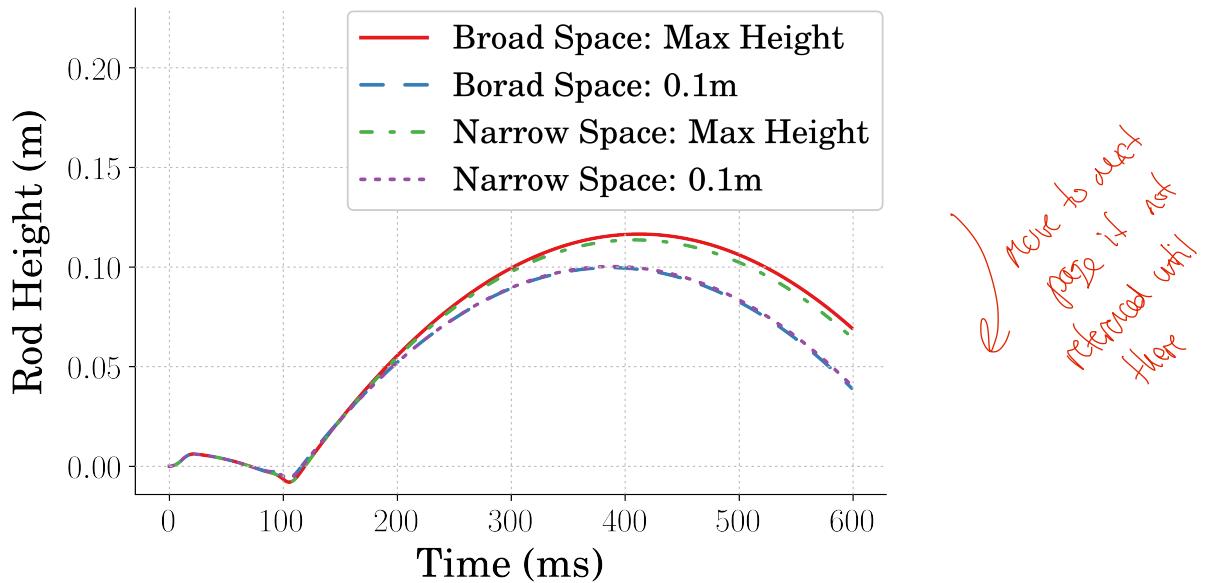
**Figure 23.** Designs Learned for the Wide Design Space

learning designs to maximize jump height, Figure 22 can be compared with Figure 18b showing that the agents learned a design nearing one which would achieve maximum performance. Additionally, looking at the agents learning designs to jump to the specified 0.1 m, the designs learned to accomplish this, again, only with slightly more variance than what is seen in the maximum height agents.

The average and standard deviation of the spring constant and damping ratio design parameters the agents selected during training are shown in Figure 23. For the

**Table 4.** Learned Design Parameters

Training Case	Design Parameter	Mean	STD
Narrow Design Space	Max Height	Spring Constant	3.62e03
		Damping Ratio	3.37e-04
	Specified Height	Spring Constant	7.74e03
		Damping Ratio	4.55e-03
Broad Design Space	Max Height	Spring Constant	3.55e03
		Damping Ratio	7.53e-03
	Specified Height	Spring Constant	7.07e03
		Damping Ratio	7.54e-03



**Figure 24.** Height vs Time of Average Optimal Designs

agents that learned designs to jump to a specified height, it can be seen that there is a high variance in spring constant throughout training. However, the majority of agents converge to a specific design, lowering the variance. The same can be seen in the damping ratio; however, the variance is mitigated significantly earlier in training. The agents which were learning designs that maximized height found them with very little variance in terms of spring constant and damping ratio.

**4.6.3 Average Design Performance.** The final mean and standard deviation of the design parameters for the two different cases are presented in Table 4. Figure 24 shows the jumping performance of the mean designs learned for both cases tested. The agents tasked with finding designs to jump to the specified 0.1 m, did so with minimal error. The difference seen in maximum height reached between the two cases represents the difference in the nominal damping ratio within the design spaces. The peak heights from Figure 24 for the maximum height designs can be compared to Figures 18a and 18b to show that the agents learned designs nearing those achieving maximum performance.

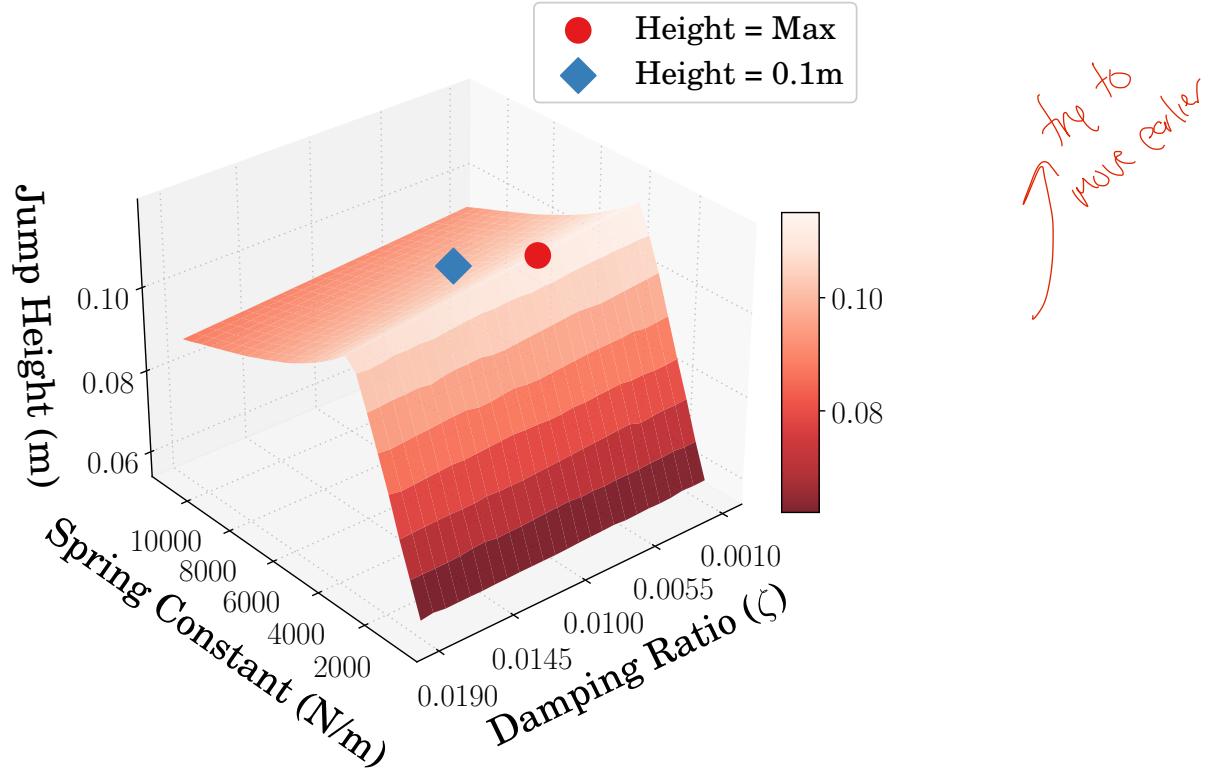
Additionally, Figure 25 shows the average designs, and their performance, against the design space performance data. It is apparent that in the case of the narrow design space, where the optimal design for maximizing height is more prominent, the design found is one that is approaching optimal performance. In the case of the wide design space, where the design that optimizes height is less prominent, it is apparent that the optimal design found by the agent is only nearing the optimal design. In both design space cases, the design found to force the monopode to jump to the specified height was found within the higher values of the spring constant range, even though values did exist in the lower range that would satisfy the jumping condition.

Why do you think that happened?

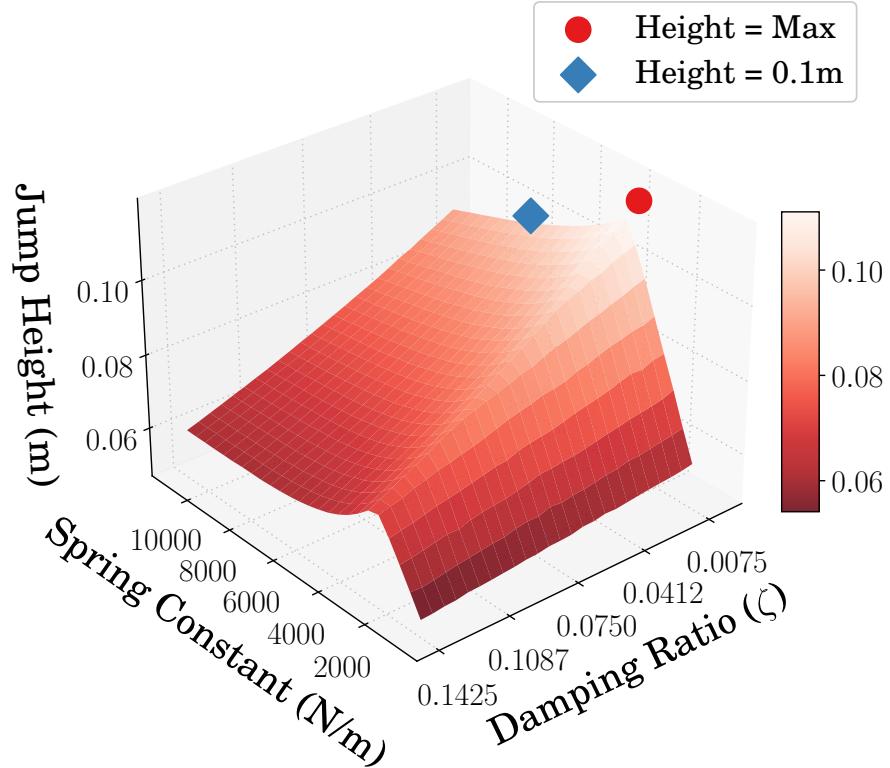
## 4.7 Conclusion

The monopode model was used in conjunction with a predetermined control input to determine if a reinforcement learning algorithm (TD3) could be used to find optimal performing design parameters regarding jumping performance. This work was done in part to determine if reinforcement learning could be used as the mechanical design learner for an intelligent concurrent design algorithm. It was shown that when providing an agent with a design space that was smaller in size with a more prominent optimal value, the agents performed well in finding design parameters which met the

performance constraints. The designs found were lower in variance as well, even in the case where the algorithm was tasked with finding a design for a specific performance within the range of possible performances. It was additionally shown that when provided with larger design space, that additionally had many values closer to the optimal value, the agents still excelled at finding design parameters that performed close to optimal. The parameters found were higher in variance however, as expected, particularly in the case where a design was to be found to generate a specific performance. This was due to the number of viable design options that would satisfy the performance requirements. It should be concluded ultimately that utilizing an RL algorithm, such as TD3, for the mechanical design aspect of a concurrent design method, is a viable solution.



(a) Average Design Performance Within Narrow Design Space



(b) Average Design Performance Within Wide Design Space

**Figure 25.** Reference Jumping Performance of the Monopode