# LEARNING CONTROL POLICIES FOR FALL PREVENTION AND SAFETY IN BIPEDAL LOCOMOTION

A Dissertation Presented to The Academic Faculty

By

Visak Kumar

In Partial Fulfillment of the Requirements for the Degree Doctor of Philosophy in Robotics

School of Mechanical Engineering Georgia Institute of Technology

December 2021

© Visak Kumar 2021

# LEARNING CONTROL POLICIES FOR FALL PREVENTION AND SAFETY IN BIPEDAL LOCOMOTION

Thesis committee:

Dr. Karen Liu Computer Science Department *Stanford University* 

Dr. Sehoon Ha School of Interactive Computing *Georgia Institute of Technology* 

Dr. Gregory Turk School of Interactive Computing *Georgia Institute of Technology*  Dr. Gregory Sawicki School of Mechanical Engineering *Georgia Institute of Technology* 

Dr. Ye Zhao School of Mechanical Engineering *Georgia Institute of Technology* 

Date approved: September 7, 2021

#### ACKNOWLEDGMENTS

First, I would like to thank my parents and sister, Vijay Kumar, Gayithri and Apoorva, whose unconditional support and love have helped me immensely along this journey in graduate school.

Next, I want to thank the members of my committee: Greg Turk, Gregory Sawicki and Ye Zhao. I have learned a lot through our interactions and greatly appreciate your valuable comments and suggestions.

I would be remiss to not express my gratitude towards my labmates and friends I have made along the way: Wenhao Yu, Yifeng Jiang, Alex Clegg, Vikash Kumar, Kendall Lowerey, Nitish Sontakke, Maksim Sorokin and many others. I had to endure completing the final components of this doctoral dissertation amidst a global pandemic. During this challenging time, without the additional support and inspiration from my dear friends Sophia, Aakash and Laura obtaining this degree would have been extremely challenging. I would also like to thank Stan Birchfield and Johnathan Tremblay for valuable guidance during past internships.

Finally, I want to express my deepest gratitude to my advisors, Karen Liu and Sehoon ha. I am fortunate to have you both as my advisors and to have learned from you. Thank you for your patience and support as I navigated the ups and downs of research journey.

# TABLE OF CONTENTS

Acknov	vledgme	ents
List of '	<b>Tables</b>	
List of ]	Figures	ix
Summa	ry	
Chapte	r 1: Int	roduction
1.1	Thesis	overview
Chapte	r 2: Rel	ated work 5
2.1	Contro	ol of bipedal robots
	2.1.1	Balance recovery strategies
	2.1.2	Safe falling 6
	2.1.3	Overview of methods which address sample efficiency in reinforce- ment learning
2.2	Contro	ol of assistive devices
	2.2.1	Simulation of human motion
	2.2.2	Design of control algorithms
	2.2.3	Transfer of RL policies

	2.2.4	Adaptation for Assistive Devices	14
Chapte	r 3: Saf	e falling and fall prevention of Bipedal Robots	15
3.1	Learni	ng control policies for safe falling	15
	3.1.1	Motivation	15
	3.1.2	Mixture of actor-critic experts	17
	3.1.3	Results	22
	3.1.4	Conclusions	27
3.2	Fall pr	evention for bipedal robots	28
	3.2.1	Motivation	28
	3.2.2	Adaptive sampling to simplify learning	29
	3.2.3	Adaptive Sampling of Perturbations	30
	3.2.4	Results	34
	3.2.5	Postural Balance Controller	35
	3.2.6	Stepping Controller	37
	3.2.7	Correlation between RoA Size and Average Reward	39
	3.2.8	Conclusion	39
3.3	Expan	ding motor skills using relay networks	40
	3.3.1	Motivation	40
	3.3.2	Method	41
	3.3.3	Learning Relay Networks	42
	3.3.4	Computing Threshold for Value Function	45
	3.3.5	Applying Relay Networks	47

	3.3.6	Extending to Multiple Strategies	47
	3.3.7	Results	48
	3.3.8	Tasks	49
	3.3.9	Baselines Comparisons	50
	3.3.10	Analyses	51
	3.3.11	Conclusion	52
Chapter	r 4: Fall	l prevention using Assistive Devices	53
4.1	Motiva	tion	53
4.2	Metho	d	54
	4.2.1	Human Walking Policy	56
	4.2.2	Fall Predictor	57
	4.2.3	Recovery Policy	58
	4.2.4	Results	59
	4.2.5	Comparison of Policy and Human Recovery Behaviors	60
	4.2.6	Effectiveness of Recovery Policy	61
	4.2.7	Evaluation of Different Design Choices	64
	4.2.8	Conclusion	65
Chapter	r 5: Err	or-aware policy learning	67
5.1	Motiva	tion	67
5.2	Metho	d	68
	5.2.1	Problem Formulation	69
	5.2.2	Training an Error-aware Policy	70

	5.2.3 Training an Error Function
5.3	Results
	5.3.1 Baseline Algorithms
	5.3.2 Tasks
	5.3.3 Zero-shot Transfer with EAPs
	5.3.4 Ablation study
5.4	Hardware experiments
	5.4.1 Data collection, simulation environment and system identification . 84
5.5	In-place walking
5.6	Walking forward
5.7	Evaluation on real robot
5.8	Conclusion
Chapte	r 6: Conclusion and future work
6.1	Safe locomotion for bipedal robots
6.2	Fall prevention using assistive devices    94
Referen	nces
Vita .	

# LIST OF TABLES

3.1	Different falling strategies.	24
3.2	Problem and Learning Parameters	33
3.3	Comparison of Various LQR Parameters	35
5.1	Tasks and Network Architectures	74
5.2	Ranges of variation for observable parameters during training and testing in the assistive walking task.	74
5.3	Ranges of variation for unobservable parameters during training and testing in the assistive walking task.	75
5.4	Tasks and Network Architectures on the real robot	86

# LIST OF FIGURES

3.1	Abstract model of the humanoid used in policy training	15
3.2	Illustration of the method to compute control signal to execute a safe fall .	15
3.3	A schematic illustration of our deep neural network that consists of $n$ actor- critics. The numbers indicate the number of neurons in each layer	17
3.4	The average reward for 10 test cases	22
3.5	The histogram of rewards for the 1000 test cases. Our policy outperforms DP in 65% of the tests.	23
3.6	<b>Top:</b> A fall from a two-feet stance due to a 3 N push (Two-feet). <b>Middle:</b> A fall from an unbalanced stance due to a 5 N push (Unbalanced). <b>Bottom:</b> A fall from a one-foot stance due to a 6 N push (One-foot)	24
3.7	Comparison of the impulse profiles among the unplanned motion, the mo- tion planned by DP, and the motion planned by our policy.	25
3.8	Comparison of measured acceleration between motion computed by our policy and unplanned motion. Three trials for each condition are plotted. <b>Left:</b> A fall from a two-feet stance due to a 3 N push. <b>Right:</b> A fall from an one-foot stance due to a 5 N push.	26
3.9	An Linear inverted pendulum abstract model used in computing control signals for ankle and hip strategy	28
3.10	An abstract model to compute stepping distance when a large external push is applied	29
3.11	Polygon representation of a RoA (red line) and the range of perturbations the controller can handle (cyan area).	31
3.12	The function $f_{\alpha}$ that maps the success rate to the update rate $\alpha$	32

3.13	The learning curves for postural balance controllers. Note that the learning curve with adaptive sampling (blue) is measured from a different set of perturbations.	35
3.14	RoAs for postural balance controllers.	36
3.15	The learned motions with adaptive sampling. <b>Top:</b> Postural balancing with the 26.0 N backward perturbation. <b>Bottom:</b> Stepping with the 81.0 N forward perturbation.	36
3.16	The learning curves for stepping controllers. Note that the learning curve with adaptive sampling (blue) is measured from a different set of perturbations.	38
3.17	RoAs for stepping controllers.	38
3.18	Illustration of a graph of policies in a relay network. The policies in the graph as sequentially executed in order to complete a task.	41
3.19	Testing curve comparisons.	51
3.20	(a) The experiment with $\alpha$ . (b) Comparison of ONE with different numbers of neurons. (c) Confusion matrix of value function binary classifier (d) Confusion matrix after additional regression.	52
4.1	Left : We model a 29-Degree of Freedom(DoF) humanoid and the 2-DoF exoskeleton in PyDart. Right : Assistive device design used in our experiments.	55
4.2	Comparison between hip and knee joint angles during walking generated by the policy and human data [17].	60
4.3	(a) Comparison of torque loops of a typical trajectory generated by our policy and human data reported by [17] at the hip of stance leg during a gait cycle. The green dots indicate the start and the black dots indicate 50% of the gait cycle. The arrows show the progression of the gait from 0% to 100%. (b) Comparison of the forward foot step locations predicted by the policy and by the model reported by Wang <i>et al.</i> [69] as a function of the COM velocity.	61
4.4	Four different timing of the left leg swing phase during which we test the performance of the assistive device. First is at 10% of the phase and then subsequently 30%, 60% and 90% of the left swing leg.	61
		-

4.5	Stability region with and without the use of a recovery policy. A larger area shows increased robustness to an external push in both magnitude and direction.	62
4.6	Comparison of recovery performance when perturbation is applied at four different phases. <b>Top:</b> Comparison of stability region. <b>Bottom:</b> Comparison of COM velocity across five gait cycles. Perturbation is applied during the gait cycle 'p'. The increasing velocity after perturbation indicates that our policy is least effective at recovering when the perturbation occurs later in the swing phase.	63
4.7	<b>Top:</b> Successful gait with an assistive device. <b>Bottom:</b> Unsuccessful gait without an assistive device. Torques are set to zero	63
4.8	Average torques at the hip joints from 50 trials with various perturbations. The shaded regions represent the 3-sigma bounds. <b>Red:</b> Joint torques exerted by the human and the recovery policy. <b>Blue:</b> Human joint torques without a recovery policy. <b>Green:</b> Torques produced by a recovery policy.	64
4.9	Stability region for six policies trained with three sensor configurations and two actuator configurations.	65
5.1	Overview of An Error-aware Policy (EAP). An EAP takes the "expected" future state error as an additional input. The expected error is predicted based on the current state s, observable parameters $\mu$ , and an uncorrected action a that assumes zero error.	69
5.2	Left : A full state error representation input into the policy vs <b>Right</b> : Projected error representation as an input to the policy	73
5.3	Five different test subjects for the assistive walking experiment with vary- ing height, mass, leg length and foot length from the biomechanical gait dataset [144].	75
5.4	Learning curves for four tasks. The number of samples for EAP include the ones generated for training an error function.	79
5.5	Comparison of EAP and baselines DR and UP. The error bars represent the varia- tion in the average return of the policy in the target environment when trained with 4 different seeds.	80
5.6	Average stability region in five test subjects. The results indicate the better zero-shot transfer of EAP over DR and UP.	80

5.7	Ablation study with choosing different observable parameters as $\mu$ . The result indicates that our approach (EAP) shows more reliable zero-shot transfers for all different scenarios.	81
5.8	Ablation study with different reference dynamics. The results indicate that our algorithm is robust against the choice of different references.	81
5.9	Ablation study with different parameter setting for EAP training	82
5.10	Unitree's A1 quardupedal robot.	83
5.11	Comparison of contact profile generated by ground and soft foam mat	90
5.12	Evaluation of real world results.	91

#### **SUMMARY**

The ability to recover from an unexpected external perturbation is a fundamental motor skill in bipedal locomotion. An effective response includes the ability to not just recover balance and maintain stability but also to fall in a safe manner when balance recovery is physically infeasible. For robots associated with bipedal locomotion, such as humanoid robots and assistive robotic devices that aid humans in walking, designing controllers which can provide this stability and safety can prevent damage to robots or prevent injury related medical costs. This is a challenging task because it involves generating highly dynamic motion for a high-dimensional, non-linear and under-actuated system with contacts. Despite prior advancements in using model-based and optimization methods, challenges such as requirement of extensive domain knowledge, relatively large computational time and limited robustness to changes in dynamics still make this an open problem.

In this thesis, to address these issues we develop learning-based algorithms capable of synthesizing push recovery control policies for two different kinds of robots : Humanoid robots and assistive robotic devices that assist in bipedal locomotion. Our work can be branched into two closely related directions : 1) Learning safe falling and fall prevention strategies for humanoid robots and 2) Learning fall prevention strategies for humans using a robotic assistive devices. To achieve this, we introduce a set of Deep Reinforcement Learning (DRL) algorithms to learn control policies that improve safety while using these robots. To enable efficient learning, we present techniques to incorporate abstract dynamical models, curriculum learning and a novel method of building a graph of policies into the learning framework. We also propose an approach to create virtual human walking agents which exhibit similar gait characteristics to real-world human subjects, using which, we learn an assistive device controller to help virtual human return to steady state walking after an external push is applied. Finally, we extend our work on assistive devices and address the challenge of transferring a push-recovery policy to different individuals. As walking

and recovery characteristics differ significantly between individuals, exoskeleton policies have to be fine-tuned for each person which is a tedious, time consuming and potentially unsafe process. We propose to solve this by posing it as a transfer learning problem, where a policy trained for one individual can adapt to another without fine tuning.

# CHAPTER 1 INTRODUCTION

We often find ourselves startled by unexpected incidents like a push, slip or a trip while walking, yet we can successfully recover from these without having to put much thought into our actions. This ability to generate a recovery motion from external perturbations is one of the most important motor skills humans learn to prevent fall related injuries, these could be motion that either prevents a fall or enables us to fall in a safe manner. Studies have shown that we develop these abilities from a very young age, as infants transition from crawling into walking, an important step towards walking efficiently [1] is to learn how to fall. The quick reflexes and balance we develop over time help us avoid harmful injuries that can be sustained during a fall.

Understanding the neuro-muscular control strategies that humans employ to achieve these exemplary motor skills has been a long standing quest in the fields of computer graphics, bio-mechanics and robotics. For bipedal robots, providing these abilities will be a step closer towards ensuring safety during operation. A safety layer in the control framework could open doors for these robots to work in unstructured outdoor environments or perhaps even enable learning more challenging dynamic locomotion skills like gymnastics or yoga. This would fulfill some of the initial motivations for designing bipedal robots such as being deployed in a disaster management scenario [2] or perform labor intensive work such as delivery [3]. From a healthcare stand point, a deeper understanding of human motion can have lasting benefits for society. For example, methods used for helping disabled or injured people with gait rehabilitation can be vastly improved. In addition, better assistive devices can be designed for locomotion leading to improved the quality of life for many frail individuals.

Planning an effective response for a bipedal system is a challenging task because it

requires generating highly dynamic motion for a high-dimensional system with unstable, nonlinear and discontinuous dynamics due to contacts. To alleviate some of these issues researchers have drawn inspiration from human strategies to design control algorithms. Prior work in this area are typically built on model-based and optimization techniques [4, 5, 6, 7] which require significant domain knowledge, large computational time and have limited robustness to changes in dynamics leaving plenty of opportunities for improvement.

In this dissertation, we focus on developing learning based methods to teach robots safe locomotion skills. More specifically, we aim to teach effective push recovery response for two different kinds of robots 1) Humanoid robots and 2) Assistive robotic devices that aid in human walking. Our work with these two kinds of robots is along two closely related directions: 1) Learning safe falling and fall prevention strategies for bipedal robots and 2) Learning fall prevention strategies for humans using robotic assistive devices.

As bipedal robots are increasingly being tested in unstructured real-world environment rather than in a laboratory setting, it is paramount that they have similar human-like ability of safe locomotion to prevent damages to the robot or to its immediate surroundings. In the first few chapters of this thesis we focus on developing a set of reinforcement learning based algorithms that enables bipedal robots learn these skills efficiently. While healthy adults have remarkable motor skills, people with disability or older adults have a reduced ability to prevent falls. According the Center for Disease Control (CDC) [8], over 3 million older adults are treated for fall related injuries resulting in millions of dollars in medical bills [9]. One in five falls result in head injury [10, 11] which makes falls the leading cause of traumatic brain injuries [12] among geriatric patients. In latter chapters of this thesis, we shift our focus towards investigating the effectiveness of reinforcement learning methods in learning control policies for assistive devices such as hip exoskeleton to prevent falls.

#### **1.1** Thesis overview

This dissertation is organized as follows - In chapter 2 we provide a comprehensive overview of the most relevant prior work in designing control policies for humanoid robots and assistive devices. We highlight some prior reinforcement learning and transfer learning methods we draw inspiration from. Next, we present our work on designing safe falling and fall prevention control policies for bipedal robots in Chapter 3. We then describe our approach to design assistive device control policies for fall prevention in humans (Chapter 4). In Chapter 5, we present a novel method to enable transfer of policies from one human agent to another in a zero-shot manner. Finally, we conclude by summarizing our contributions and providing some exciting future directions this work can be extended to.

For bipedal robots, In our first work [13] we addressed the challenge of optimizing a safe falling motion for a humanoid robot. Falling in a safe manner involves minimizing the impulse when a body part makes contact with a surface. Ideally, if a robot could turn itself into a ball during a fall, the resulting rolling motion would produce minimal impulse, similar to how gymnasts or martial artists fall when they lose balance. Since this has physical constraints for a robot we solve for two quantities, planning the next contacting body location as well as the corresponding joint motion to minimize the impulse of contact. We designed an actor-critic DRL algorithm that learns n actors and critics simultaneously, where each critic corresponds to selecting the next contact location while the corresponding actor provides the necessary continuous joint motion in order to fall safely. We showed that using our approach we are able to solve for falling motion that generated lesser impulse on contact compared to prior approaches [5] while significantly speeding up the required computational time. In our next work [14], we developed a curriculum learning method to learn control policies for humanoid push recovery and balancing. Often solving a challenging task for high-dimensional robot requires a curriculum to enable efficient learning. During training, we maintain as estimate of the current push-magnitude the policy is capable of handling and provide pushes around those magnitudes. We showed that a policy learned using our approach outperforms other DRL baselines by handling larger push magnitudes. In [15], we introduce a novel framework to simplify learning complex tasks, such as 2D and 3D walking without falling, by building a directed graph of local control policies which automatically decomposes a complex task into sub-tasks. Using this approach, we show that control policies perform better while consuming lesser samples compared to naive DRL baselines.

For assistive devices, we introduce an algorithm in [16] to learn a hip exoskeleton policy that helps a virtual human return to steady state walking after an external push is applied. To achieve this, we first model steady state human walking in simulation using DRL and open source bio-mechanical motion capture data. We demonstrate that the virtual human exhibits similar dynamical behaviour to real humans by comparing joint kinematics, kinetics [17] and foot-step lengths to real world data. Next, we learn a recovery policy for exoskeleton to help the virtual human overcome external pushes applied to the hip. A thorough analysis is presented on the efficacy of the exoskeleton policy to help recovery and we also showed that our method can provide insights into the sensory and actuation capability required by exoskeletons to perform optimally.

In our final work, we design an algorithm to address the issue of transferring the learned exoskeleton policy to different individuals. As walking gait characteristics differ significantly between individuals, exoskeleton policies have to be designed for each person. This process can be tedious and time consuming, we propose to solve this by posing it as a transfer learning problem. More specifically, we extend our prior work [16] and take the following steps 1) Train multiple virtual human walking agents , each of them varying in physical characteristics like height, mass, etc. as well as ability to recover. 2) Learn an error function that predicts the difference in dynamics between them and 3) Train a policy which is explicitly aware of the difference enabling adaptation and efficient transfer.

# CHAPTER 2 RELATED WORK

In this chapter, we give an overview of the most relevant prior work. Our contributions can be broadly categorised into two topics - First, algorithms to teach safe falling and balance recovery for humanoid robots and second, fall prevention using assistive devices for humans. In the first section, we highlight prior work done at the intersection of human inspired control strategies for bipedal robots and reinforcement learning. Often, reinforcement learning suffers from the requirement of a large sample budget, we also present a summary of earlier work that addressed this shortcoming. Then, we provide an overview of common approaches to design control algorithms for assistive devices such as exoskeletons for humans. Finally, we outline prior work which focused on developing algorithms that enable transfer of control policies for robots from training to testing domains.

#### 2.1 Control of bipedal robots

Researchers in the field of robotics, computer graphics and bio-mechanics have worked on designing control policies for bipedal robots for a long time. Due to the inherent instability of bipedal locomotion, falling and loss of balance are a common occurrence for these robots, naturally, researchers have spent extensive efforts on developing methods to maintain stability or fall safely to prevent damage. Due to the challenging nature of designing these controllers, a common approach that has worked well is to draw inspiration from human strategies by first understanding human motion and then developing control algorithms that replicate the desired behaviour. There is a vast body of research in this area ranging from deriving linear controllers from simple inverted pendulum dynamical models to more complex model predictive controllers. We highlight some relevant ones in this section.

#### 2.1.1 Balance recovery strategies

As a response to external perturbations, humans employ multiple strategies to maintain balance. For example when pushed, bio=mechanical studies on postural balance control have shown that we regain balance by controlling the center of pressure on the foot using torques generated at the ankle and hip [18, 19, 20]. Inspired by such strategies, [4, 21, 22, 23] have derived control algorithms that are built strategies seen in humans. For larger perturbations, controllers that combine ankle, hip and foot-placement strategies have been proven successful in the humanoid robot community [24, 25, 26, 27, 28]. Many of these controllers rely on simplified dynamical models, such as Linear Inverted Pendulum (LIP) or 3D LIP, to compute control signals. While simplified models are beneficial for computational efficiency, they come at the cost of accuracy and inability to adapt to novel situations. Further, they also fail to capture arm motions that are sometimes crucial in maintaining balance. On the other hand, after the pioneering work of Mnih et al. [29] deep reinforcement learning methods such as [30] and [31] have shown promising results in learning effective control policies for high-dimensional systems. However, for challenging tasks, such as balance recovery, reinforcement learning methods can require a large sample budget to learn the skill. Our work [14] aims to combine the best of both worlds by building on top of these simplified models while leveraging the benefits of reinforcement learning to learn a control policy for a humanoid robot to recover from an external push. To reduce the required sample budget, we present a curriculum approach that simplifies learning process by providing external pushes of adequate difficulty during training.

#### 2.1.2 Safe falling

Recovery from an unstable state is important, but sometimes it is inevitable to prevent a fall, in such situations a controller which plans a safe falling motion is necessary to prevent damage to the robot and potentially to the surrounding environment as well. A plethora of motion planning and optimal control algorithms have been proposed to reduce the damage

of humanoid falls. These algorithms have also been largely inspired by human strategies, for example rolling as one falls is an effective strategy. One direct approach to generating falling motion is to use knee and torso flexion to generate a squatting motion and backward to minimize impulse upon contact with the ground [32, 33, 34, 35, 36], this strategy works best when the robot is falling backwards. Another heuristic based approach is to design a few falling motion sequences for a set of expected scenarios. When a fall is detected, the sequence designed for the most similar falling scenario is executed [37, 38]. This approach, albeit simple, can be a practical solution in an environment in which the types of falls can be well anticipated. To handle more dynamic environments, a number of researchers cast the falling problem to an optimization which minimizes the damage of the fall. To accelerate the computation, various simplified models have been proposed to approximate falling motions, such as an inverted pendulum [39, 40], a planar robot in the sagittal plane [41], a tripod [42], and a sequence of inverted pendulums [5]. In spite of the effort to reduce the computation, most of the optimization-based techniques are still too slow for real-time applications, with the exception of the work done by [6], who proposed to compute the optimal stepping location to change the falling direction. In contrast, our work takes the approach of policy learning using deep reinforcement learning techniques. Once trained, the policy is capable of handling various situations with real-time computation.

Our work is also built upon recent advancement in deep reinforcement learning (DRL). Although the network architecture used in this work is not necessarily "deep", we borrow many key ideas from the DRL literature to enable training a large network with 278976 variables. The ideas of "experience replay" and "target network" from Mnih *et al.* [43] are crucial to the efficiency and stability of our learning process, despite that the original work (DQN) is designed for learning Atari video games from pixels with the assumption that the action space is discrete. Lilicrap *et al.* [44] later combined the ideas of DQN and the deterministic policy gradient (DPG) [45] to learn actor-critic networks for continuous action space and demonstrated that end-to-end (vision perception to actuation) policies for

dynamic tasks can be efficiently trained.

The approach of actor-critic learning has been around for many decades [46]. The main idea is to simultaneously learn the state-action value function (the Q-function) and the policy function, such that the intractable optimization of the Q-function over continuous action space can be avoided. van Hasselt and Wiering introduced CACLA [47, 48] that represents an actor-critic approach using neural networks. Our work adopts the update scheme for the value function and the policy networks proposed in CACLA. Comparing to the recent work using actor-critic networks [44, 49], the main difference of CACLA (and our work as well) lies in the update scheme for the actor networks. That is, CACLA updates the actor by matching the action samples while other methods follow the gradients of the accumulated reward function. Our work is mostly related to the MACE algorithm introduced by Peng *et al.* [50]. We adopt their network architecture and the learning algorithm but for a different purpose: instead of using multiple actor-critic pairs to switch between experts, we exploit this architecture to solve an MDP with a mixture of discrete and continuous action variables.

#### 2.1.3 Overview of methods which address sample efficiency in reinforcement learning

Recently, many successful model-free reinforcement learning (RL) algorithms have been developed which show remarkable promise in learning challenging motor skills [43, 51, 31, 52, 30, 45, 47, 49]. However, one of the biggest shortcomings of these methods is requirement of a large sample budget. Additionally, policies trained using RL are typically effective only from a small set of initial states defined during training. To address these issues, one approach researchers have taken is to cast a large-scale task as a hierarchical reinforcement learning algorithm. This approach decomposes problems using *temporal abstraction* which views sub-policies as macro actions, or *state abstraction* which focuses on certain aspects of state variables relevant to the task. Prior work [53, 54, 55, 56, 57] that utilizes temporal abstraction applies the idea of parameterized goals and pseudo-rewards

to train macro actions, as well as training a meta-controller to produce macro actions. One notable work that utilizes state abstraction is MAXQ value function decomposition which decomposes a large-scale MDP into a hierarchy of smaller MDP's [58, 59, 60]. MAXQ enables individual MDP's to only focus on a subset of state space, leading to better performing policies. Our relay networks can be viewed as a simple case of MAXQ in which the recursive subtasks, once invoked, will directly take the agent to the goal state of the original MDP. That is, in the case of relay networks, the Completion Function that computes the cumulative reward after the subtask is finished always returns zero. As such, our method avoids the need to represent or approximate the Completion Function, leading to an easier RL problem for continuous state and action spaces.

Yet another approach is to use a set of policies chained together to solve a challenging task. Tedrake [61] proposed the LQR-Tree algorithm that combines locally valid linear quadratic regulator (LQR) controllers into a nonlinear feedback policy to cover a wider region of stability. Borno et al. [62] further improved the sampling efficiency of RRT trees [61] by expanding the trees with progressively larger subsets of initial states. However, the success metric for the controllers is based on Euclidean distance in the state space, which can be inaccurate for high dimensional control problems with discontinuous dynamics. Konidaris et al. [63] proposed to train a chain of controllers with different initial state distributions. The rollouts terminate when they are sufficiently close to the initial states of the parent controller. They discovered an initiation set using a sampling method in a low dimensional state space. In contrast, our algorithm utilizes the value function of the parent controller to modify the reward function and define the terminal condition for policy training. There also exists a large body of work on scheduling existing controllers, such as controllers designed for taking simple steps [64] or tracking short trajectories [65]. Inspired by the prior art, our work demonstrates that the idea of sequencing a set of local controllers can be Manipulating the initial state distribution during training has been considered a promising approach to accelerate the learning process. [66] studied theoretical foundation of using "exploratory" restart distribution for improving the policy training. Recently, [67] demonstrated that the learning can be accelerated by taking the initial states from successful trajectories. [68] proposed a reverse curriculum learning algorithm for training a policy to reach a goal using a sequence of initial state distributions increasingly further away from the goal. We also train a policy with reversely-ordered initial states, but we exploited chaining mechanism of multiple controllers rather than training a single policy.realized by learning policies represented by the neural networks.

## 2.2 Control of assistive devices

#### 2.2.1 Simulation of human motion

To study human motion during walking and recovery, biomechanics researchers often adopt an experimental approach. First, data is collected in the real-world, then control policies are synthesized using the real-world data. Winters et al [17] was among the first in the field to study human gait, and the data published in this work remains relevant to this day. Wang et al [69] and Hof et al [70] performed perturbation experiments and identified important relationships between COM velocity, step-lengths, center of pressure, stepping vs ankle strategy, etc.. We aim to leverage the finding of this research to validate some of the results. In Joshi et al [71], a balance recovery controller was derived using the results reported in [69], however, they use a 3D Linear Inverted pendulum model to approximate the human dynamics. A 3D LIPD does not capture the dynamics fully, for example, angular momentum about the center of mass. Most relevant to our work, Antoine et al [72], used a direct-collocation trajectory optimization to synthesize a walking controller for a 3D musculo-skeletal model in OpenSim (OpenSim gait2392 model). The gait generated by the controller closely matched experimental data. The proposed method, relies on understanding the basic principles that lead to walking, such as minimizing metabolic cost, muscle activations, etc. However, the proposed solution enforces left-right symmetry, which works for walking, but is not ideal for disturbance recovery. Hence its unclear how well this

approach will perform when there is an external disturbance to the human.

#### 2.2.2 Design of control algorithms

Many researchers have developed control algorithms for robotic assistive walking devices. However, majority of the work can be classified into tracking controllers [73, 74, 75] and model-based controllers [76, 77, 78]. There has been limited work at the intersection of DRL and control of assistive devices, especially for push recovery. Hamaya et al [79] presented model-based reinforcement learning algorithm to train a policy for a handheld device that takes muscular effort measured by electromyography signals (EMGs) as inputs. This method requires collecting user interaction data on the real-world device to build a model of the user's EMG patterns. However, collecting a large amount of data for lower-limb assistive devices is less practical due to the safety concerns. Another recent work, Bingjing et al [80], developed an adaptive-admittance model-based control algorithm for a hip-knee assistive device, the reinforcement learning aspect of this work focused on learning parameters of the admittance model. Our method [81] is agnostic to the assistive walking devices and can be used to augment any device that allows for feedback control.

#### 2.2.3 Transfer of RL policies

A popular approach to transfer control policies is Domain randomization (DR). DR methods [82, 83, 84, 85, 86] propose to train policies that are robust to variations in the parameters that affect the system dynamics. Although some of these methods have been validated in the real world [82, 85], DR often requires manual engineering of the range in which the parameters are varied to make sure that the true system model lies within the range of variation. For a complex robotic system, it is often challenging to estimate the correct range of all the parameters because a large range of variation could lead to lower task performance, whereas a smaller range leads to less robust policies. To address the demanding sample budget issue with domain randomization, [87] presented a data-efficient domain randomization algorithm based on bayesian optimization. The algorithm presented in Mehta et al [88] actively adapts the randomization range of variation to alleviate the need for exhaustive manual engineering. Ramos et al [89] proposed an approach to infer the distribution of the dynamical parameters and showed that policies trained with randomization within this distribution can transfer better.

Careful identification of parameters using data from the real world, popularly known as system identification, has also shown promising results in real-world robots. Tan et al [90] and Hwangbo et al [91] carefully identified the actuator dynamics to bring the source environment closer to the target, Xie et al [92] also demonstrated that careful system identification techniques can transfer biped locomotion policies from simulation to real-world. Jegorova et al [93] presented a technique that improves on existing system identification techniques by borrowing ideas from generative adversarial networks (GAN) and showed improved ability to identify the parameters of a system. Similarly, Jiang et al [94] presented a SimGAN algorithm that identifies a hybrid physics simulator to match the simulated trajectories to the ones from the target domain to enable policy adaptation. Yu et al [95] developed a method that combines online system identification and universal policy to enable identifying dynamical parameters in an online fashion. Citing the difficulty in obtaining meaningful data for system identification, [96] developed an algorithm that probes the target environment to provide more information about the dynamics of the environment. A few model based approaches have also been successful in transferring policies to a target domain [97, 98, 99].

Another popular approach of transferring policies includes utilizing data from the target domain to improve the policy. Chebotar et al [100] presented a method that interleaves policy learning and system identification, however this requires deploying the policy in the target domain every few iterations. This method would be impractical for a system that interacts closely with a human because of safety concerns. Yu et al [101] and Peng et al [102] presented latent space adaptation techniques where the policy is adapted in the target domain by searching for a latent space input to the policy that enables successful transfer. Exarchos et al [103] also presented an algorithm that achieved policy transfer using only kinematic domain randomization combined with policy adaptation in the target domain, similar to [101].

Yu et al [104] proposed Meta Strategy Optimization, a meta-learning algorithm for training policies with latent variables that can quickly adapt to new scenarios with a handful of trials in the target environment. Among the methods that use data from the target domain also include meta-learning approaches like Bhelkale et al [105], in which a model-based meta-reinforcement learning algorithm was presented to account for changing dynamics of an aerial vehicle carrying different payloads. In this approach, the parameters causing the variations in the dynamics are inferred by deploying the policy in the target domain, which in turn helps improve the policy's performance. In Ignasi et al. [106], the idea of model-agnostic meta-learning [107] was extended to modelling dynamics of a robot. The authors presented an approach to quickly adapt the model of the robot in a new test environment while using a sampling-based controller MPPI to compute the actions. [108] developed a zero-shot transfer for policy by combining reinforcement learning and a robust tracking controller with a disturbance observer in the target environment. The validated the approach on a vehicle driving task. Similarly, [109, 110] presented an approach to combine bayesian learning and adaptive control by learning model error and uncertainty.

For tasks such as assistive device control for human locomotion, it is potentially unsafe and prohibitive to collect sufficient task-relevant data in the real world which prevents us from using methods such as system identification or transfer learning approaches that need data in the target environment. In addition to this, human dynamics exhibit large variations due to many unobserved parameters, this makes it challenging to define the right parameters for the system model in simulation and also in finding the right range of parameter variation for an approach like DR.

#### 2.2.4 Adaptation for Assistive Devices

Assistive devices such as exoskeletons provide unique challenges for domain adaptation due to the large variations between individuals who pilot the device. Zhang et al [111] reported a human-in-the-loop optimization approach for ankle exoskeletons to account for this variability, however, this approach takes a few hours per individual to find the optimal control law. Jackson et al [112] presented a unique heuristic-based approach to design a control law that adapts to the person's muscle activity. While these methods work well for steady-state walking, the large number of data required to optimize for in the case of [111] and the complex muscle responses involved during push recovery make it an infeasible application. Several recent works have incorporated a learning-based approach to tackling the problem of adaptation, Peng et al [113] adopted a reinforcement learning approach to learn assistive walking strategies for Hemiplegic patients, which was tested on real human patients and showed robustness and adaptability. However, it requires online data to update the actor-critic network. This process involves deploying a policy on a patient to collect data, for a task like push recovery it might be challenging to collect relevant data required for updating the policy without compromising the patient's safety. Both [114] and [115] combined dynamic motion primitives (DMPs) and learning approaches to adapt control strategies for different individuals. Majority of the work with assistive devices have primarily focused on walk assistance and not on push-recovery.

## **CHAPTER 3**

# SAFE FALLING AND FALL PREVENTION OF BIPEDAL ROBOTS

## 3.1 Learning control policies for safe falling

## 3.1.1 Motivation



Figure 3.1: Abstract model of the humanoid used in policy training.

As research efforts to push bipedal robots out of a laboratory setting and into the real world are increasing, the ability for robots to tackle unexpected situations in unstructured environments becomes crucial. For bipedal robots, due to the inherently unstable nature of dynamics, falling and losing balance are situations where being equipped with control policies that can prevent damages to the robot are paramount for success in the real world. Unlike walking , which is periodic in nature, falling is more challenging problem as it involves reasoning about both discrete and continuous quantities.



Figure 3.2: Illustration of the method to compute control signal to execute a safe fall

First, the policy must decide which body part should make contact with the ground, the location on the ground as well as the timing of the contact. The policy must also output the

corresponding joint trajectory which minimizes impact upon contact. For example, humans employ a very effective strategy of rolling in order to reduce impact. Prior methods reported in generating walking motions which leverage the periodic nature of walking using finitestate machines (FSM) [116, 117, 118, 119] can be challenging to apply to the task of falling. Even approaches that mimic human motion using strong priors such as motion capture data [120, 121, 122, 123] are not entirely suitable due to lack of motion capture data of a wide variety of falls.

To this end, we developed a policy optimization approach [13] for minimizing the damage to a humanoid robot during a fall. In our approach, we first decide n candidate contacting body part such as hands, feet, or knees designed to be a contact point with the ground. Our algorithm trains n control policies (actors) and the corresponding value functions (critics) in a single interconnected neural network. Each policy and its corresponding value function are associated with a candidate contacting body part. During policy execution, the network is queried every time the robot establishes a new contact with the ground, allowing the robot to re-plan throughout the fall. Based on the current state of the robot, the actor corresponding to the critic with the highest value will be executed while the associated body part will be the next contact with the ground. With this mixture of actor-critic architecture, the discrete contact sequence planning is solved through the selection of the best critics while the continuous control problem is solved by the optimization of actors. To simplify learning, the policy uses an abstract model representation of a humanoid robot, shown in Figure 3.1, consisting of an inverted pendulum and a telescopic rod to plan actions which are then mapped back to a full joint space using inverse kinematics. Figure 3.2 illustrates the sequence of operations involved in planning a safe fall. We use CACLA policy learning algorithm [47, 48] to optimize the control policy. We show that our proposed approach is able to generate motions that lead to less impact while falling down compared to prior approaches [5] while also significantly speeding up the computational time (50 to 400 times faster).



Figure 3.3: A schematic illustration of our deep neural network that consists of n actorcritics. The numbers indicate the number of neurons in each layer.

# 3.1.2 Mixture of actor-critic experts

Figure 3.2 illustrates the workflow of the overall policy. We first define a function p:  $\mathcal{X} \mapsto \mathcal{S}$  which maps a state of robot ( $\mathbf{x} \in \mathcal{X}$ ) to a state of abstract model ( $\mathbf{s} \in \mathcal{S}$ ). The mapping can be easily done because the full set of joint position and velocity contains all the necessary information to compute  $\mathbf{s}$ . As the robot receives an external force initially, the state of the robot is projected to  $\mathcal{S}$  and fed into the abstract-level policy. The action (a) computed by the abstract-level policy is passed into the joint-level policy to compute the corresponding joint torques ( $\tau$ ). If no new contact is detected after executing  $\tau$ , the new state of the robot will be fed back into the joint-level policy and a new  $\tau$  will be computed (the lower feedback loop in Figure Figure 3.2). If the robot encounters a new contact, we re-evaluate the contact plan by querying the abstract-level policy again (the upper feedback loop in Figure 3.2).

## Abstract-level Policy

To overcome the challenge of optimizing over both discrete and continuous action variables, we introduce a new policy representation based on a neural network architecture inspired by MACE [50]. The policy takes as input the state of the abstract model and outputs the action, as well as the next contacting body part. The state space is defined as  $\mathbf{s} = \{c_1, r_1, \theta_1, \dot{r}_1, \dot{\theta}_1\}$ , where the total elapsed time t is removed from the state space previously defined by Ha and Liu [5]. As a side benefit of a feedback policy, we no longer need to keep track of the total elapsed time. For the action space, we remove the discrete action variable,  $c_2$  so that the new action is a continuous vector in  $\mathbb{R}^3$ :  $\mathbf{a} = \{\theta_2, \Delta, \dot{r}_1^d\}$ . The network combines n pairs of actor-critic subnets, each of which is associated with a contacting body part. Each actor,  $\Pi_i(\mathbf{s}) : \mathbb{R}^5 \mapsto \mathbb{R}^3$ , represents a control policy given that the *i*-th contacting body part is chosen to be the next contact. Each critic,  $V_i(\mathbf{s}) : \mathbb{R}^5 \mapsto \mathbb{R}$ , represents the value function that evaluates the return (long-term reward) of using the *i*-th contacting body part as the next contact and taking the action computed by  $\Pi_i(\mathbf{s})$  as the next action. We fuse all n actor-critic pairs in one single network with a shared input layer (Figure 3.3).

At each impact moment when a new contact is established, the network evaluates all the  $V_i(\mathbf{s})$ ,  $1 \le i \le n$ , and selects the policy corresponding to the highest critic. This arrangement allows us to train n experts, each specializes to control the robot when a particular contacting body part is selected to be the next contact. As a result, we cast the discrete contact planning into the problem of expert selection, while simultaneously optimizing the policy in continuous space.

#### Reward

Since our goal is to minimize the maximal impulse, we define the reward function as:

$$r(\mathbf{s}, \mathbf{a}) = \frac{1}{1+j},\tag{3.1}$$

where j is the impulse induced by the contact. Suppose the COM of the pendulum and the tip of the stopper are  $(x_1, y_1)$  and  $(x_2, y_2)$  respectively at the impact moment, the impulse can be computed as:

$$j = -\frac{\dot{y}_2^-}{\frac{1}{M} + \frac{1}{I}(x_2 - x_1)^2},$$
(3.2)

where M is the mass and I is the inertia of the abstract model (see details in [5]).

With this definition of the reward function, the objective of the optimization is to maximize the minimal reward during the fall.

#### Learning Algorithm

Algorithm 1 illustrates the process of learning the abstract-level policy. We represent the policy using a neural network consisting n pairs of actor-critic subnets with a shared input layer (Figure 3.3). Each critic has two hidden layers with 32 neurons each. The first hidden layer is shared among all the critics. Each actor has 3 hidden layers with 128 neurons each. All the hidden layers use tanh as the activation functions. We define weights and biases of the network as  $\theta$ , which is the unknown vector we attempt to learn.

The algorithm starts off with generating an initial set of high-reward experiences, each of which is represented as a tuple:  $\tau = (\mathbf{s}, \mathbf{a}, \mathbf{s}', r, c)$ , where the parameters are the starting state, action, next state, reward, and the next contacting body part. To ensure that these tuples have high reward, we use the dynamic-programming-based algorithm described in [5] (referred as DP thereafter) to simulate a large amount of rollouts from various initial states and collect tuples. Filling the training buffer with these high-reward experiences accelerates the learning process significantly. In addition, the high-reward tuples generated by DP can guide the abstract-level policy to learn "achievable actions" when executed on a full body robot. Without the guidance of these tuples, the network might learn actions that increase the return but unachievable under robot's kinematic constraints.

In addition to the initial experiences, the learning algorithm continues to explore the action space and collect new experiences during the course of learning. At each iteration, we simulate K(=10) rollouts starting at a random initial state sampled from a Gaussian distribution  $\mathcal{N}_0$  and terminating when the abstract model comes to a halt. A new tuple is generated whenever the abstract model establishes a new contact with the ground. The exploration is done by stochastically selecting the critic and adding noise in the chosen actor. We follow the same Boltzmann exploration scheme as in [50] to select the actor

Algorithm 1: Learning abstract-level policy

```
1: Randomly initialize \theta
 2: Initialize training buffer with tuples from DP
 3: while not done do
 4:
          EXPLORATION:
          for k = 1 \cdots K do
 5:
              \mathbf{s}\sim\mathcal{N}_0
 6:
              while \mathbf{s}.\dot{\theta}_1 \ge 0 do
 7:
                  c \leftarrow Select actor stochastically using Equation Equation 3.3
 8:
 9:
                  \mathbf{a} \leftarrow \Pi_c(\mathbf{s}) + \mathcal{N}_t
                  Apply a and simulate until next impact moment
10:
                  \mathbf{s}' \leftarrow \text{Current state of abstract model}
11:
                  r \leftarrow r(\mathbf{s}, \mathbf{a})
12:
                  Add tuple \tau \leftarrow (\mathbf{s}, \mathbf{a}, \mathbf{s}', r, c) in training buffer
13:
                  \mathbf{s} \leftarrow \mathbf{s}'
14:
              end while
15:
          end for
16:
17:
          UPDATE CRITIC:
          Sample a minibatch m tuples \{\tau_i = (\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i, c_i)\}
18:
          y_i \leftarrow \min(r_i, \gamma \max_j \hat{V}_j(\mathbf{s}'_i)) for each \tau_i
19:
          \theta \leftarrow \theta + \alpha \sum_{i} (y_i - V_{c_i}(\mathbf{s})) \nabla_{\theta} V_{c_i}(\mathbf{s})
20:
21:
          UPDATE ACTOR:
          Sample a minibatch m tuples \{\tau_i = (\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i, c_i)\}
22:
          y_i = \max_i V_i(\mathbf{s})
23:
          y'_i \leftarrow \min(r_i, \gamma \max_i \hat{V}_i(\mathbf{s}'_i))
24:
          if y'_i > y_i then
25:
              \theta \leftarrow \theta + \alpha (\nabla_{\theta} \Pi_{c_i}(\mathbf{s}))^T (\mathbf{a}_i - \Pi_{c_i}(\mathbf{s}))
26:
          end if
27:
```

```
28: end while
```

 $\Pi_i(\mathbf{s})$  based on the probability defined by the predicted output of the critics:

$$\mathcal{P}_i(\mathbf{s}) = \frac{e^{V_i(\mathbf{s}|\theta)/T_t}}{\sum_j e^{V_j(\mathbf{s}|\theta)/T_t}},$$
(3.3)

where  $T_t(=5)$  is the temperature parameter, decreasing linearly to zero in the first 250 iterations. While the actor corresponding to the critic with the highest value is most likely to be chosen, the learning algorithm occasionally explores other actor-critic pairs, essentially trying other possible contacting body parts to be the next contact. Once an actor is selected, we add a zero-mean Gaussian noise to the output of the actor. The covariance of the Gaussian is a user-defined parameter.

After K rollouts are simulated and tuples are collected, the algorithm proceeds to update the critics and actors. In critic update, a minibatch is first sampled from the training buffer with m(=32) tuples,  $\tau_i = (\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i, c_i)$ . We use the temporal difference to update the chosen critic similar to [43]:

$$y_{i} = \min(r_{i}, \gamma \max_{j} \hat{V}_{j}(\mathbf{s}'_{i}))$$

$$\theta \leftarrow \theta + \alpha \sum_{i} (y_{i} - V_{c_{i}}(\mathbf{s})) \nabla_{\theta} V_{c_{i}}(\mathbf{s})$$
(3.4)

where  $\theta$  is updated by following the negative gradient of the loss function  $\sum_i (y_i - V_{c_i}(\mathbf{s}))^2$ with the learning rate  $\alpha (= 0.0001)$ . The discount factor is set to  $\gamma = 0.9$ . Note that we also adopt the idea of target network from [43] to compute the target,  $y_i$ , for the critic update. We denote the target networks as  $\hat{V}(\mathbf{s})$ .

The actor update is based on supervised learning where the policy is optimized to best match the experiences:  $\min_{\theta} \sum_{i} ||\mathbf{a}_{i} - \prod_{c_{i}}(\mathbf{s}_{i})||^{2}$ . We use the positive temporal difference



Figure 3.4: The average reward for 10 test cases.

to decide whether matching a particular tuple is advantageous:

$$y = \max_{j} V_{j}(\mathbf{s}_{i})$$

$$y' = \min(r_{i}, \gamma \max_{j} \hat{V}_{j}(\mathbf{s}'_{i}))$$
if  $y' > y, \quad \theta \leftarrow \theta + \alpha (\nabla_{\theta} \Pi_{c_{i}}(\mathbf{s}))^{T}(\mathbf{a}_{i} - \Pi_{c_{i}}(\mathbf{s})).$ 
(3.5)

# 3.1.3 Results

We validate our policy in both simulation and on physical hardware. The testing platform is a small humanoid, BioloidGP [124] with a height of 34.6 cm, a weight of 1.6 kg, and 16 actuated degrees of freedom. We also compare the results from our policy against those calculated by the dynamic-programming (DP) based method proposed by Ha and Liu [5]. Because DP conducts a full search in the action space online, which is 50 to 400 times slower than our policy, in theory DP should produce better solutions than ours. Thus, the goal of the comparison is to show that our policy produces comparable solutions as DP while enjoying the speed gain by two orders of magnitude.

We implement and train the proposed network architecture using PyCaffe [125] on Ubuntu Linux. For simulation, we use a Python binding [126] of an open source physics engine, DART [127].


Figure 3.5: The histogram of rewards for the 1000 test cases. Our policy outperforms DP in 65% of the tests.

# Learning of Abstract-Level Policy

In our experiment, we construct a network with 8 pairs of actor-critic to represent 8 possible contacting body parts: right toe, right heel, left toe, left heel, knees, elbows, hands, and head. During training, we first generate 5000 tuples from DP to initialize the training buffer. The learning process takes 1000 iterations, approximately 4 hours on 6 cores of 3.3GHz Intel i7 processor. Figure Figure 3.4 shows the average reward of 10 randomly selected test cases over iterations. Once the policy is trained, a single query of the policy network takes approximately 0.8 milliseconds followed by 25 milliseconds of the inverse kinematics routine. The total of 25.8 milliseconds computation time is a drastic improvement from DP which takes 1 to 10 seconds of computation time.

We compare the rewards of the trained policy with those of DP by running 1000 test cases starting from randomly sampled initial states. The results are shown in Figure Figure 3.5, a histogram of rewards for the 1000 tests computed by our policy and by DP. Our policy achieves not only comparable rewards, it actually outperforms DP in 64 % of the test cases. The average reward of our policy is 0.8093, comparing to 0.7784 of DP. In terms of impulse, the average of maximum impulse produced by our policy is 0.2540, which shows a 15 % improvement from 0.2997 produced by DP.

Theoretically, DP, which searches the entire action space for the given initial state, should be the upper bound of the reward our policy can ever achieve. In practice, the discretization of the state and action spaces in DP might result in suboptimal plans. In



Figure 3.6: **Top:** A fall from a two-feet stance due to a 3 N push (Two-feet). **Middle:** A fall from an unbalanced stance due to a 5 N push (Unbalanced). **Bottom:** A fall from a one-foot stance due to a 6 N push (One-foot).

Initial Pose	Push	Algorithm	Contact Sequence	Maximum Impulse
Two-feet	3 N	Unplanned	Torso (0.90)	0.90
		DP	Hands (0.58)	0.58
		Ours	Knees (0.53), Hands (0.20)	0.53
Unbalanced	5 N	Unplanned	Torso (2.50)	2.50
		DP	Hands (0.53)	0.53
		Ours	Hands (0.45)	0.45
One-foot	6 N	Unplanned	Torso (2.10)	2.10
		DP	L-Heel (0.20), Hands (0.45)	0.45
		Ours	L-Heel (0.10), Hands (0.43)	0.43

Table 3.1: Different falling strategies.

contrast, our policy exploits the mixture of actor and critic network to optimize continuous action variables without discretization. The more precise continuous optimization often results in more optimal contact location or timing. In some cases, it also results in different contact sequences being selected (45 out of 641 test cases where our policy outperforms DP).

# Different falling strategies

With different initial conditions, various falling strategies emerge as the solution computed by our policy. Table 3.1 showcases three distinctive falling strategies from the test



Figure 3.7: Comparison of the impulse profiles among the unplanned motion, the motion planned by DP, and the motion planned by our policy.

cases. The table shows the initial pose, the external push, the resulting contact sequence, the impulse due to each contact (the number in the parenthesis), as well as the maximal impulse for each test case. Starting with a two-feet stance, the robot uses knees and hands to stop a fall. If the initial state is leaning forward and unbalanced, the robot directly uses its hands to catch itself. If the robot starts with a one-foot stance, it is easer to use the swing foot followed by the hands to stop a fall. The robot motion sequences can be visualized in Figure 3.15 and the supplementary video. For each case, we compare our policy against DP and a naive controller which simply tracks the initial pose (referred as *Unplanned*). Both our policy and DP significantly reduce the maximum impulse comparing to Unplanned. In the cases where our policy outperforms DP, the improvement can be achieved by different contact timing (One-foot case), better target poses (Unbalanced case), or different contact sequences (Two-feet case, Figure Figure 3.7).

# Hardware Results

Finally, we compare the falling strategy generated by our policy against the unplanned motion on the hardware of BioloidGP. Due to the lack of on-board sensing capability, Bi-



Figure 3.8: Comparison of measured acceleration between motion computed by our policy and unplanned motion. Three trials for each condition are plotted. **Left:** A fall from a two-feet stance due to a 3 N push. **Right:** A fall from an one-foot stance due to a 5 N push.

oloidGP cannot take advantage of the feedback aspect of our policy. Nevertheless, we can still use this platform to demonstrate the falling strategy generated by our policy and compare it against an unplanned motion.

We first match the initial pose of the simulated BioloidGP with the real one and push the simulated BioloidGP from the back by 3 N and 5 N, assuming that the pushes we applied to the robot by hand are about the same. We then apply our policy on the simulated BioloidGP to obtain a sequence of target poses. In the hardware experiment, we program BioloidGP to track these poses once a fall is detected. During the falls, we measure the acceleration of the head using an external IMU. Figure 3.8 shows the results of two different falls. In the first case, the robot is pushed with a force of 3N and is initialized with both the feet on the ground and an upright position, the robot uses its knees first and then the hands to control the fall. The maximal acceleration from our policy is 2.9 G while that from an unplanned motion is 5.7 G, showing a 49% of improvement. In the second case, the robot is pushed with one foot on the ground, the falling strategy for this includes using the left-heel first then the hands to control the fall. The maximal acceleration from our policy is 2.3 G while that from an unplanned motion is 6.4 G, showing a 64% of improvement.

# 3.1.4 Conclusions

We proposed a new policy optimization method to learn the appropriate actions for minimizing the damage of a humanoid fall. Unlike most optimal control problems, the action space of our problem consists of both discrete and continuous variables. To address this issue, our algorithm trains n control policies (actors) and the corresponding value functions (critics) in an actor-critic network. Each actor-critic pair is associated with a candidate contacting body part. When the robot establishes a new contact with the ground, the policy corresponding to the highest value function will be executed while the associated body part will be the next contact. As a result of this mixture of actor-critic architecture, we cast the discrete contact planning into the problem of expert selection, while optimizing the policy in continuous space. We show that our algorithm reliably reduces the maximal impulse of a variety of falls. Comparing to the previous work [5] that employs an expensive dynamic programming method during online execution, our policy can reach better reward and only takes 0.25% to 2% of computation time on average.

One limitation of this work is the assumption that humanoid falls primarily lie on the sagittal plane. This limitation is due to our choice of the simplified model, which reduces computation time but only models planar motions. This assumption can be easily challenged when considering real-world falling scenarios, such as those described in [42, 6]. One possible solution to handling more general falls is to employ a more complex model similar to the inertia-loaded inverted pendulum proposed by [6]. Our algorithm also assumes that the robot is capable of motion that is fast enough to achieve the pose that the network outputs. Sensors to detect contacts are also important to use the trained policy as a feedback controller.

Another possible future work direction is to learn control policies directly in the fullbody joint space, bypassing the need of an abstract model and the restrictions come with it. This allows us to consider more detailed features of the robot during training, such as full body dynamics or precise collision shapes. Given the increasingly more powerful policy



Figure 3.9: An Linear inverted pendulum abstract model used in computing control signals for ankle and hip strategy

learning algorithms for deep reinforcement learning [31, 128], motor skill learning with a large number of variables, as is the case with falling, becomes a feasible option.

# **3.2** Fall prevention for bipedal robots

# 3.2.1 Motivation

In the previous section, we described our approach to generate falling motion using policy optimization method. However, in many situations a fall can be prevented by employing some effecting balance strategies that humans use. Similar to falling, balancing strategies are often not periodic in nature and require special control algorithms to achieve stability.

In Kumar *et al.* [14], we propose a curriculum learning framework to learn push recovery control policies for humanoid robots. Our algorithm aimed to improve existing model-based balance controllers inspired by human balance recovery motion like hip and ankle strategies [4, 21, 22] and stepping strategies [24, 25, 26, 27, 28]. Typically, modelbased control algorithms use an abstract dynamical model (like Linear-inverted pendulum LIPM or the 3D version 3D-LIPM) for computational tractability, however lower dimensional models fail to capture accurate dynamics of the system. Abstract models also neglect arm motions which can be important for balancing. Model-free reinforcement learning, on the other hand, has been shown to work well for high-dimensional systems at the cost of



Figure 3.10: An abstract model to compute stepping distance when a large external push is applied

sample efficiency during training. To take advantage of both these methods, our policy outputs residual control signals in addition to control computed by model-based controllers thereby simplifying learning the control policy.

Figure 3.9 and Figure 3.10 illustrate the abstract dynamical models used by modelbased controllers to compute the control signal. The objective function to train the policy minimizes Center-of-Mass position and velocity changes when reacting to pushes to encourage balancing. The training process also incorporates an adaptive sampling scheme that identifies the magnitude of push that needs to be applied to enable efficient learning. Our key insight here was based on observation that during training, applying pushes of large magnitudes could be detrimental to learning where as pushes of low magnitudes do not help the policy learn much. We maintain an estimate of the push magnitudes the policy is currently capable of handling , called *Region of Attraction* illustrated in **??**, and draw perturbations near its boundary. We show that the policy trained using our approach is capable of handling pushes of larger magnitudes compared to simple DRL baselines while also being able to generate appropriate arm-motions to aid in balance recovery.

# 3.2.2 Adaptive sampling to simplify learning

The control policy takes as input a partial observation vector  $\mathbf{o} = [\mathbf{C}^T, \dot{\mathbf{C}}^T, \mathbf{q}^T]^T$ , where each term represents the COM position C, the COM velocity  $\dot{\mathbf{C}}$ , and the joint positions

q. The outputs are 12 control signals  $\mathbf{u} = [\Delta \tau_h^T, \Delta \tau_a^T, \Delta \mathbf{q}_s^T]^T$  where  $\Delta \tau_h$  and  $\Delta \tau_a$  are offsets to the hip and ankle torques and  $\Delta \mathbf{q}_s$  is the offset to the shoulder target angles. All terms have four entries for the pitch and roll axes in the left and right limbs. The final torques are computed as:

$$au = \mathbf{PD}(\mathbf{ar{q}} + \mathbf{I}_s \Delta \mathbf{ar{q}}_s, \mathbf{q}, \mathbf{\dot{q}}) + \mathbf{I}_h(\mathbf{ar{ au}}_h + \mathbf{\Delta au}_h) + \mathbf{I}_a(\mathbf{ar{ au}}_a + \mathbf{\Delta au}_a)$$
 (3.6)

where  $\bar{\mathbf{q}}$ ,  $\bar{\tau}_h$ ,  $\bar{\tau}_a$  are outputs from the model-based controllers and  $\mathbf{I}_h$ ,  $\mathbf{I}_a$ ,  $\mathbf{I}_s$  matrices map the hip, ankle, and shoulder joints to the full-body joint space. The PD controller PD provides the joint torques to track the modified target position  $\bar{\mathbf{q}} + \mathbf{I}_s \Delta \bar{\mathbf{q}}_s$  for all joints except hips and ankles.

Another important component of RL is the reward function. One of the possible reward functions is a binary success flag that is 1 if and only if the COM height is greater than the certain threshold  $\bar{z}$ . However, training such a binary function is usually time-consuming and impractical. Rather, we choose a simple continuous function that penalizes the COM position and velocity as follows:

$$r(\mathbf{q}, (\dot{\mathbf{q}})) = -w_p |\bar{\mathbf{C}} - \mathbf{C}(\mathbf{q})|^2 - w_d |\dot{\mathbf{C}}(\mathbf{q})|^2$$
(3.7)

where the first term penalizes the positional error of the COM and the second term penalizes the velocity. The terms  $w_p$  and  $w_d$  are the weights for adjusting the scales.

### 3.2.3 Adaptive Sampling of Perturbations

This section describes our scheme for adaptively sampling perturbations. In our implementation, perturbations are parameterized by the directions and magnitudes of external forces. The objective of adaptive sampling is to expedite the learning process by providing more informative trials. For instance, trying to recover from an extremely strong perturbation will not be a useful data point because the robot will quickly fall even with the optimal



Figure 3.11: Polygon representation of a RoA (red line) and the range of perturbations the controller can handle (cyan area).

control policy. However, it is difficult to determine whether a randomly sampled perturbation is too weak or too strong because it depends on the performance of the current policy.

Our key idea is to maintain the RoA information during the learning process and sample perturbations for training using a probabilistic distribution around the boundary of the current RoA. This process requires us to implement the following three functionalities:

- Representing the RoA,
- Sampling perturbations from the given RoA, and
- Updating the RoA.

More precisely, we define the RoA  $\mathcal{R}$  as the range of perturbations for which the controller shows success rates of above 90 %. A trial is declared success if the COM height  $C_z$ at the end of the simulation is greater than a user-defined threshold  $\bar{z}$ .

We represent the RoA as a polygon with M sides where the vertices are defined in the polar coordinate by a set of magnitudes  $\mathcal{R} = \{\Omega_1, \Omega_2, \dots, \Omega_M\}$  at predefined angles  $\{\theta_1, \theta_2, \dots, \theta_M\}$  chosen to uniformly cover 0 to  $2\pi$  radians (Figure Figure 3.11). The assumption underlying this representation is that a controller that works for a particular



Figure 3.12: The function  $f_{\alpha}$  that maps the success rate to the update rate  $\alpha$ .

perturbation is likely to recover from a weaker perturbation in the same direction. Although there may be counterexamples, this assumption allows us to use a simplified representation that works in practice. In some cases, it is more convenient to represent the magnitude as a function of the angle, i.e.  $\Omega = f_{\mathcal{R}}(\theta)$ .

The adaptive sampling method is summarized in Algorithm algorithm 2. For each episode of learning, we first randomly sample  $\theta \in \{\theta_1, \theta_2, \dots, \theta_M\}$ , and then sample  $\Omega \in [k_l f_{\mathcal{R}}(\theta), k_u f_{\mathcal{R}}(\theta)]$  where  $k_l \leq 1$  and  $1 \leq k_u$  are predefined constants. An example of RoA representation and the range of permissible perturbations is shown in Figure Figure 3.11.

During the training process, we update  $\mathcal{R}$  for every K iterations using a simple feedback rule based on the success rates at recent episodes. At the beginning, we initialize  $\mathcal{R}$  with a constant magnitude  $\tilde{\Omega}$  in all directions, i.e.  $\mathcal{R} = \{\tilde{\Omega}, \dots, \tilde{\Omega}\}$ . We collect the success rates in each direction by counting the number of good and bad trials  $n_i^{good}$  and  $n_i^{bad}$  where i is the index of the direction closest to the direction of the sampled perturbation. At every Kiterations, we increase  $\omega_i$  if the success rate  $n_i^{good}/(n_i^{good} + n_i^{bad})$  is greater than 0.9, and shrink it otherwise. The update rate  $\alpha$  is defined by a function  $f_{\alpha}$  of the success rate shown in Figure Figure 3.12.

Algorithm 2: Learning with Adaptive Sampling 1 Initialize  $\mathcal{R}$  with a constant  $\hat{\Omega}$ ; 2 Initialize  $n_1^{good}, n_1^{bad}, \cdots, n_M^{good}, n_M^{bad}$  to 0; 3 Initialize the policy  $\Phi$ ; 4 for each iteration do for each trial do 5  $\theta = random(0, 2\pi);$ 6  $\Omega = random(k_l f_{\mathcal{R}}(\theta), k_u f_{\mathcal{R}}(\theta));$ 7 run a simulation with perturbation  $\theta$ ,  $\Omega$ ; 8 i =index of direction closest to  $\theta$ ; 9 if  $C_z > \bar{z}$  then  $\lfloor n_i^{good} +=1;$ 10 11 else 12  $n_i^{bad}$ +=1; 13 update the policy  $\Phi$ ; 14 if every K-th iteration then 15 for each direction i do 16 17 18 set  $n_1^{good}, n_1^{bad}, \cdots, n_M^{good}, n_M^{bad}$  to 0; 19 **20 return** trained policy  $\Phi$ ;

Catagony	Nome	Postural	Stepping
Category	Name	Controller	Controller
	Perturbation Magnitudes	(1.0, 51.0)	(41.0, 81.0)
Droblom	Perturbation Angles	$(-\pi,\pi)$	$(-\pi/4,\pi/4)$
1 I ODICIII	Position Weight $(w_p)$	1.0	1.0
	Velocity Weight $(w_d)$	0.1	0.5
	Neural Network Structure	32,32	32, 32, 32
	Activation Functions	tanh	tanh
TPPO	Initial Standard Deviation	0.3	0.3
INIO	Batch Size	50000	50000
	Max Iteration	1200	1200
	Step Size	0.01	0.01
	# of Sides $(M)$	8	3
	Initial Magnitude $(\tilde{\Omega})$	5.00	41.0
Adaptive Sampling	Update Frequency $(K)$	4	4
	Sampling Range Low $(k_l)$	0.8	0.7
	Sampling Range High $(k_l)$	1.1	1.1

We use the simulation model of the humanoid robot COMAN [129] to conduct our experiments. The humanoid robot model is about 0.9 meters tall and 31 kg in weight. It has 23 joints (7 in each leg, 4 in each arm, and 3 in the torso) and 31 DoFs including the six underactuated DoFs of the floating base, although we do not use the elbow joints in our controller. The maximum torque of each joint is set to 35 Nm, which is 70 % of the specification sheet values. The proportional and derivative gains for the PD controller are set to 100.0 Nm/rad and 1.0 Nms/rad respectively.

The simulations are conducted with an open-source physics simulation engine Py-DART [130]. It handles contacts and collisions by formulating velocity-based linearcomplementarity problems to guarantee non-penetration and approximated Coulomb friction cone conditions. The simulation and control time steps are set to 0.002 s (500 Hz), which is enough to be executed on the actual COMAN hardware. The computations are conducted on a single core of 3.40GHz Intel i7 processor.

For all cases, perturbations are applied to the torso of the robot for 0.2 seconds. The perturbation direction is defined as the direction of the external force. For instance, a backward perturbation (180°) is a push from the front direction that causes the robot to fall backward.

We use the TRPO [131] implementation of Duan *et al.* [132] to train the control policy whose structures and activation functions are listed in Table Table 3.2. Learning of each control policy takes 12 hours in average. Also refer to Table Table 3.2 for other hyper-parameters for problems and algorithms.

In the following sections, we will compare the following three learning methods:

- Naive RL
- RL with a model-based controller and uniform sampling
- RL with a model-based controller and adaptive sampling

State Weight (Q)	Control Weight (R)	Maximum Backward Perturb.	Maximum Forward Perturb.
diag(100, 100, 1, 1)	diag(0.005, 0.005)	-11 N	11 N
diag(200, 200, 2, 2)	diag(0.005, 0.005)	-16  N	26 N
diag(1000, 1000, 10, 10)	diag(0.005, 0.005)	-21  N	46 N
diag(2000, 2000, 20, 20)	diag(0.005, 0.005)	-16  N	41 N
-20 -20 -20 -20 -20 -20 -20 -20 -20 -20	RL + LQ	R R + Adaptive Sz	nan an
0 200	400 600 8 Iteration	00 1000	1200

Table 3.3: Comparison of Various LQR Parameters

Figure 3.13: The learning curves for postural balance controllers. Note that the learning curve with adaptive sampling (blue) is measured from a different set of perturbations.

The baseline controller for naive RL is tracking with local PD position servos at each joint to maintain a given target pose  $\bar{q}$ . Uniform sampling draws the perturbations from all possible ranges defined in Table Table 3.2.

# 3.2.5 Postural Balance Controller

First, we apply the proposed learning framework to the postural balance controller.

We select the LQR parameters by comparing various feedback matrices K from different state and input weights Q and R (Table Table 3.3). Based on the maximum permissible perturbation, we set Q to diag(1000, 1000, 10, 10) and R to diag(0.005, 0.005) that yields the following feedback gain:

$$\mathbf{K} = \begin{bmatrix} -448 & -1.30 & -45.2 & -0.14 \\ -1.30 & 447 & -0.13 & 47.3 \end{bmatrix}.$$
 (3.8)

First, we compare the learning curves of RL and RL with the LQR controller (Figure Figure 3.13). Within the same iterations, learning with the LQR achieves a reward of



Figure 3.14: RoAs for postural balance controllers.



Figure 3.15: The learned motions with adaptive sampling. **Top:** Postural balancing with the 26.0 N backward perturbation. **Bottom:** Stepping with the 81.0 N forward perturbation.

-39.5 while naive RL achieves -64.8. Therefore, it is not surprising that the controller with the LQR has a larger RoA than that of RL, especially for forward perturbations that the LQR controller itself is already able to handle well (Figure Figure 3.14).

Although learning with the LQR controller is successful, the resulting RoA is not as expanded as we expected. For example, the final controller becomes more vulnerable to backward pushes than the initial controller, even though the total area of the RoA has been increased. We suspect that the final controller tends to bend forward to increase the RoA in all other directions while sacrificing the ability to recover from backward perturbations.

By training with adaptive sampling, we are able to achieve a larger RoA than with uniform sampling. The most improved directions of perturbations are backward  $(180^{\circ})$ 

and left (90°), where the controller can recover from 50 % and 38 % larger perturbations than other controllers, respectively. When we take a look at the motions, we observe more proper arm reactions that are important to recover from side pushes.

Note that the learning curve of the adaptive sampling method cannot be compared directly to those of uniform sampling methods because the sample perturbations are drawn from different distributions. For example, the initial average reward for adaptive sampling is much greater than uniform sampling at the first iteration due to the small  $\tilde{\Omega}$  we chose for the initial RoA estimation. Even within the adaptive sampling method, the sample distributions are different depending on the shape of the estimated RoA. Therefore, the average reward of -16.0 at the 100 th iteration and the 1200 th iteration have different meanings.

The learning curve with adaptive sampling fluctuates a lot at the beginning of the learning process because the initial estimation of the RoA does not match well with the actual RoA. Once they become similar to each other as learning proceeds, the reward is stabilized and shows slow changes. The relatively flat learning curve is due to the fact that all the successful trials have similar rewards mostly between -18.0 and -10.0. We believe this relationship between the reward and RoA expansion can be an interesting future work.

# 3.2.6 Stepping Controller

We also conduct experiments for expanding the RoA of the stepping controller Section **??**. For this set of experiments, we only consider forward perturbations Table 3.2.

We compare the learning curves in Figure Figure 3.16. Learning without a modelbased controller cannot recover balance from most of perturbations, even though it shows improvement in terms of the reward function value. We believe that the stepping behavior is difficult to automatically emerge without providing any prior knowledge, even after many iterations. On the other hand, learning with a model-based controller is able to handle a wide range of perturbations. The trained policy successfully adjusts stepping locations for completely stopping the robot by generating additional torques to the hip and ankle joints.



Figure 3.16: The learning curves for stepping controllers. Note that the learning curve with adaptive sampling (blue) is measured from a different set of perturbations.



Figure 3.17: RoAs for stepping controllers.

It also shows reactive upper body motions to compensate induced angular momentum.

The adaptive sampling expands the RoA for the stepping controller more effectively. Note that the learning curve of adaptive sampling (Figure Figure 3.16) does not have the initial fluctuation seen with the postural controller (Figure Figure 3.13). This is because the initial RoA estimate is larger than the actual RoA, which was smaller in the case of postural control. Interestingly, learning with adaptive sampling almost sequentially expands each direction of RoA: it learns first how to step right, and uses the experience to the front and left steps. It might be due to asymmetric upper body motions of the initial control policy that tends to move the left arm forward. We believe that this can be another way of structuring tasks, which might be an interesting direction for future work.

# 3.2.7 Correlation between RoA Size and Average Reward

One interesting observation from our experiments is that the average reward is not necessarily correlated to the RoA size, which makes learning with uniform sampling more difficult if the objective is maximizing the RoA. This is mainly because the rewards of failed simulation samples have larger variances than successful trials. Therefore, the average reward can increase when the controller slightly improves failed samples while sacrificing the successful samples. This issue can be even worse when there exists many impossible perturbations that cannot be handled by any controllers.

However, it is difficult to identify the exact RoA of the optimal controller before training many different controllers. Therefore, learning with uniform sampling is likely to have many impossible samples that can skew the average reward. On the other hand, adaptive sampling can reduce impossible trials by continuously updating the current RoA.

# 3.2.8 Conclusion

This paper presented a learning framework for enhancing the performance of humanoid balance controllers so that they can recover from a wider range of perturbations. The key idea is to combine reinforcement learning with model-based controllers in order to expedite the learning process. We also proposed a novel adaptive sampling scheme that maintains the RoA during the learning process and samples the perturbations for training around its boundary.

We conducted simulation experiments using two model-based controllers: an LQRbased postural balance controller and an LIP-based stepping controller. We demonstrated that the proposed framework can improve the controller performance within the same number of iterations. Furthermore, the controllers trained with adaptive sampling can accommodate larger perturbations than those with uniform sampling.

We tested the framework with different hyper-parameters for the proposed adaptive sampling method. In our experience, the performance of the learning algorithm is not very sensitive to the hyper-parameters. For instance, the algorithm is able to automatically adjust the estimated RoA even if we choose wrong initial magnitudes ( $\tilde{\Omega}$ ). The update frequency K also did not significantly affect the convergence rate. However, we observed that a very narrow range such as  $k_l=0.9$  and  $k_u=1.0$  caused over-fitting for stepping controllers and cannot handle weaker perturbations around 41 N.

Although we empirically demonstrated the effectiveness of our framework, future work could include more theoretical analysis. For instance, the idea of drawing samples from the RoA boundary is based on intuition. It would be more convincing if we can provide more statistical data that those boundary samples are more useful for learning than others.

Another possible future direction is to apply the proposed framework to other control problems such as locomotion or getting up. In particular, the polygon (or polyhedron) representation of the RoA may not be sufficient for these tasks due to high-dimensional state space. Furthermore, the discontinuity due to collisions may result in discontinuous RoA that cannot be expressed with the current representation. However, we believe that the general idea of providing the learning process with tasks of appropriate difficulty levels can be effective in other problems.

#### **3.3** Expanding motor skills using relay networks

# 3.3.1 Motivation

Often, challenging robotic tasks suffer from sub-optimal behaviour due to finding local optimum. Our insight into developing this algorithm is that, instead of learning one control policy, the task could by simplified by learning multiple policies which work towards a common goal of completing the task. Our approach begins by defining an initial set of states from which achieving the goal is easy. For example, in the classical control task of pendulum swing-up and balance, we start by learning a policy that achieves pole balancing. Then, our algorithm gradually expands by finding states where the first policy fails, and adds a new policy to the graph where the goal is to reach states in which the first policy is

good at. The algorithm gradually expands the set of successful initial states by connecting new policies, one at a time, to the existing graph until the robot can achieve the original complex task. Each policy is trained for a new set of initial states with an objective function that encourages the policy to drive the new states to previously identified successful states.

In our next work [15], we introduced a new technique of building a graph of policies to simplify learning a challenging control task. Similar to hierarchical reinforcement learning (HRL) [58, 133] which decomposes a large-scale Markov Decision Process (MDP) into sub-tasks, the key idea of this work is to build a directed graph of local policies represented by neural networks, which we refer to as *relay neural networks*, illustrated in Figure 3.18. The nodes of the graph are state distributions and the edges are control policies that connect two state distributions.

We apply our algorithm to a set of challenging control problems and show that the policy is capable of solving the tasks in a sample efficient manner compared to baselines such as a single policy or curriculum learning. In addition to classical control problems like Cartpole swingup, Hopper and 2D walker, we show that with a combination of policies, we can control a 3-Dimensional



Figure 3.18: Illustration of a graph of policies in a relay network. The policies in the graph as sequentially executed in order to complete a task.

humanoid character with 23-Degrees of Freedom (DoF) with a goal of walking forward while also balancing.

#### 3.3.2 Method

Our approach to a complex robotic task is to automatically decompose it to a sequence of subtasks, each of which aims to reach a state where the policy of its preceding subtask is able to handle. The original complex task is formulated as a MDP described by a tuple  $\{S, A, r, T, \rho, P\}$ , where S is the state space, A is the action space, r is the reward function, T is the set of termination conditions,  $\rho = N(\mu_{\rho}, \Sigma_{\rho})$  is the initial state distribution, and *P* is the transition probability. Instead of solving for a single policy for the MDP, our algorithm solves for a set of policies and value functions for a sequence of simpler MDP's. A policy,  $\pi : S \times A \mapsto [0, 1]$ , is represented as a Gaussian distribution of action  $\mathbf{a} \in A$  conditioned on a state  $\mathbf{s} \in S$ . The mean of the distribution is represented by a fully connected neural network and the covariance is defined as part of the policy parameters to be optimized. A value function,  $V : S \mapsto R$ , is also represented as a neural network whose parameters are optimized by the policy learning algorithm.

We organize the MDP's and their solutions in a directed graph  $\Gamma$ . A node  $\mathcal{N}_k$  in  $\Gamma$  stores the initial state distribution  $\rho_k$  of the  $k^{th}$  MDP, while a directed edge connects an MDP to its parent MDP and stores the solved policy ( $\pi_k$ ), the value function ( $V_k$ ), and the threshold of value function ( $\bar{V}_k$ ) (details in Section subsection 3.3.4). As the robot expands its skill set to accomplish the original task, a chain of MDP's and policies is developed in  $\Gamma$ . If desired, our algorithm can be extended to explore multiple solutions to achieve the original task. Section subsection 3.3.6 describes how multiple chains can be discovered and merged in  $\Gamma$  to solve multi-modal problems.

# 3.3.3 Learning Relay Networks

The process of learning the relay networks (Algorithm algorithm 3) begins with defining a new initial state distribution  $\rho_0$  which reduces the difficulty of the original MDP (Line Line 3). Although our algorithm requires the user to define  $\rho_0$ , it is typically intuitive to find a  $\rho_0$  which leads to a simpler MDP. For example, we can define  $\rho_0$  as a Gaussian whose mean,  $\mu_{\rho_0}$ , is near the goal state of the problem.

Once  $\rho_0$  is defined, we proceed to solve the first subtask  $\{S, A, r, T, \rho_0, P\}$ , whose objective function is defined as the expected accumulated discounted rewards along a trajectory:

$$J_0 = \mathbb{E}_{\mathbf{s}_{0:t_f}, \mathbf{a}_{0:t_f}} [\sum_{t=0}^{t_f} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t)],$$
(3.9)

where  $\gamma$  is the discount factor,  $s_0$  is the initial state of the trajectory drawn from  $\rho_0$ , and

A	lgori	thm	3:	LearnRel	layN	etwor	ks
---	-------	-----	----	----------	------	-------	----

1: Input: MDP  $\{S, \overline{A}, r, \overline{T}, \rho, P\}$ 

2: Add root node,  $\overline{\mathcal{N}} = \{\emptyset\}$ , to  $\Gamma$ 

3: Define a simpler initial state distribution  $\rho_0$ 

4: Define objective function  $J_0$  according to Equation Equation 3.9

5:  $[\pi_0, V_0] \leftarrow \text{PolicySearch}(\mathcal{S}, \mathcal{A}, \rho_0, J_0, \mathcal{T})$ 

6:  $\bar{V}_0 \leftarrow \text{ComputeThreshold}(\pi_0, V_0, \mathcal{T}, \rho_0, J_0)$ 

7: Add node,  $\mathcal{N}_0 = \{\rho_0\}$ , to  $\Gamma$ 

8: Add edge,  $\mathcal{E} = \{\pi_0, V_0, \overline{V}_0, \}$ , from  $\mathcal{N}_0$  to  $\overline{\mathcal{N}}$ 

9: k = 0

10: while  $\pi_k$  does not succeed from  $\rho$  do

11:  $\mathcal{T}_{k+1} = \mathcal{T} \cup (V_k(\mathbf{s}) > \overline{V}_k)$ 

12:  $\rho_{k+1} \leftarrow$  Compute new initial state distribution using Equation Equation 3.10

```
13: Define objective function J_{k+1} according to Equation Equation 3.11
```

14:  $[\pi_{k+1}, V_{k+1}] \leftarrow \text{PolicySearch}(\mathcal{S}, \mathcal{A}, \rho_{k+1}, J_{k+1}, \mathcal{T}_{k+1})$ 

15:  $\overline{V}_{k+1} \leftarrow \text{ComputeThreshold}(\pi_{k+1}, V_{k+1}, \mathcal{T}, \rho_{k+1}, J_{k+1})$ 

16: Add node,  $\mathcal{N}_{k+1}^i = \{\rho_{k+1}\}$ , to  $\Gamma$ 

17: Add edge,  $\mathcal{E} = \{\pi_{k+1}, V_{k+1}, \overline{V}_{k+1}\}$ , from node  $\mathcal{N}_{k+1}$  to  $\mathcal{N}_k$ 

18:  $k \leftarrow k+1$ 

19: end while

20: return  $\Gamma$ 

 $t_f$  is the terminal time step of the trajectory. We can solve the MDP using PPO/TRPO or A3C/DDPG to obtain the policy  $\pi_0$ , which drives the robot from a small set of initial states from  $\rho_0$  to states where the original task is completed, as well as the value function  $V_0(\mathbf{s})$ , which evaluates the return by following  $\pi_0$  from state s (Line Line 5).

The policy for the subsequent MDP aims to drive the rollouts toward the states which  $\pi_0$  can successfully handle, instead of solving for a policy that directly completes the original task. To determine whether  $\pi_0$  can handle a given state s, one can generate a rollout by following  $\pi_0$  from s and calculate its return. However, this approach can be too slow for online applications. Fortunately, many of the modern policy gradient methods, such as PPO or A3C, produce a value function, which provides an approximated return from s without generating a rollout. Our goal is then to determine a threshold  $\bar{V}_0$  for  $V_0(s)$  above which s is deemed "good" (Line Line 6). The details on how to determine such a threshold are described in Section subsection 3.3.4. We can now create the first node  $\mathcal{N}_0 = \{\rho_0\}$  and add

it to  $\Gamma$ , as well as an edge  $\mathcal{E} = \{\pi_0, V_0, \overline{V}_0\}$  that connects  $\mathcal{N}_0$  to the dummy root node  $\overline{\mathcal{N}}$ (Line Line 7-Line 8).

Starting from  $\mathcal{N}_0$ , the main loop of the algorithm iteratively adds more nodes to  $\Gamma$ by solving subsequent MDP's until the last policy  $\pi_k$ , via relaying to previous policies  $\pi_{k-1}, \cdots \pi_0$ , can generate successful rollouts from the states drawn from the original initial state distribution  $\rho$  (Line Line 10). At each iteration k, we formulate the  $(k + 1)^{th}$  subsequent MDP by redefining  $\mathcal{T}_{k+1}$ ,  $\rho_{k+1}$ , and the objective function  $J_{k+1}$ , while using the shared S, A, and P. First, we define the terminal conditions  $\mathcal{T}_{k+1}$  as the original set of termination conditions ( $\mathcal{T}$ ) augmented with another condition,  $V_k(\mathbf{s}) > \overline{V}_k$  (Line Line 11). Next, unlike  $\rho_0$ , we define the initial state distribution  $\rho_{k+1}$  through an optimization process. The goal of the optimization is to find the mean of the next initial state distribution,  $\mu_{\rho_{k+1}}$ , that leads to unsuccessful rollouts under the current policy  $\pi_k$ , without making the next MDP too difficult to relay. To balance this trade-off, the optimization moves in a direction that reduces the value function of the most recently solved MDP,  $V_k$ , until it reaches the boundary of the "good state zone" defined by  $V_k(s) \ge \overline{V}_k$ . In addition, we would like  $\mu_{
ho_{k+1}}$  to be closer to the mean of the initial state distribution of the original MDP,  $\mu_{
ho}$ . Specifically, we compute  $\mu_{\rho_{k+1}}$  by minimizing the following objective function subject to constraints:

$$\boldsymbol{\mu}_{\rho_k} =_{\mathbf{s}} V_{k-1}(\mathbf{s}) + w \|\mathbf{s} - \boldsymbol{\mu}_{\rho}\|^2$$
  
subject to  $V_{k-1}(\mathbf{s}) \ge \bar{V}_{k-1}$   
 $C(\mathbf{s}) \ge 0$  (3.10)

where  $C(\mathbf{s})$  represents the environment constraints, such as the constraint that enforces collision-free states. Since the value function  $V_k$  in Equation Equation 3.10 is differentiable through back-propagation, we can use any standard gradient-based algorithms to solve this optimization. Procedure for selecting the weighting coefficient  $\mathbf{w}$  is explained in the appendix.

In addition, we define the objective function of the  $(k + 1)^{th}$  MDP as follows:

$$J_{k+1} = \mathbb{E}_{\mathbf{s}_{0:t_f}, \mathbf{a}_{0:t_f}} \left[ \sum_{t=0}^{t_f} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \gamma^{t_f} g(\mathbf{s}_{t_f}) \right],$$
(3.11)  
where  $g(\mathbf{s}_{t_f}) = \begin{cases} V_k(\mathbf{s}_{t_f}), & V_k(\mathbf{s}_{t_f}) > \bar{V}_k \\ 0 & \text{otherwise.} \end{cases}$ 

Besides the accumulated reward, this cost function has an additional term to encourage "relaying". That is, if the rollout is terminated because it enters the subset of S where the policy of the parent node is capable of handling (i.e.  $V_k(\mathbf{s}_{t_f}) > \overline{V}_k$ ), it will receive the accumulated reward by following the parent policy from  $\mathbf{s}_{t_f}$ . This quantity is approximated by  $V_k(\mathbf{s}_{t_f})$  because it recursively adds the accumulated reward earned by each policy along the chain from  $\mathcal{N}_k$  to  $\mathcal{N}_0$ . If a rollout terminates due to other terminal conditions (e.g. falling on the ground for a locomotion task), it will receive no relay reward. Using this objective function, we can learn a policy  $\pi_{k+1}$  that drives a rollout towards states deemed good by the parent's value function, as well as a value function  $V_{k+1}$  that measures the long-term reward from the current state, following the policies along the chain (Line Line 14). Finally, we compute the threshold of the value function (Line Line 15), add a new node,  $\mathcal{N}_{k+1}$ , to  $\Gamma$  (Line Line 16), and add a new edge that connects  $\mathcal{N}_{k+1}$  to  $\mathcal{N}_k$  (Line Line 17).

The weighting parameter  $\alpha$  determines the importance of "relaying" behavior. If  $\alpha$  is set to zero, each MDP will attempt to solve the original task on its own without leveraging previously solved policies. The value of  $\alpha$  in all our experiments is set to 30 and we found that the results are not sensitive to  $\alpha$  value, as long as it is sufficiently large (Section subsection 3.3.10).

#### 3.3.4 Computing Threshold for Value Function

In practice, V(s) provided by the learning algorithm is only an approximation of the

Algorithm 4: ComputeThreshold

1: **Input:**  $\pi, V, \mathcal{T}, \rho = N(\boldsymbol{\mu}_{\rho}, \boldsymbol{\Sigma}_{\rho}), J$ 2: Initialize buffer  $\mathcal{D}$  for training data 3:  $[\mathbf{s}_1, \cdots, \mathbf{s}_M] \leftarrow \text{Sample states from } N(\boldsymbol{\mu}_{\rho}, 1.5\Sigma_{\rho})$ 4:  $[\tau_1, \cdots, \tau_M] \leftarrow$  Generate rollouts by following  $\pi$  and  $\mathcal{T}$  from  $[\mathbf{s}_1, \cdots, \mathbf{s}_M]$ 5: Compute returns for rollouts:  $R_i = J(\tau_i), i \in [1, M]$ 6:  $\bar{R} \leftarrow$  Compute average of returns for rollouts not terminated by  $\mathcal{T}$ 7: for i = 1 : M do if  $R_i > \overline{R}$  then 8: 9: Add  $(V(\mathbf{s}_i), 1)$  in  $\mathcal{D}$ else 10: Add  $(V(\mathbf{s}_i), 0)$  in  $\mathcal{D}$ 11: 12: end if 13: end for 14:  $\bar{V} \leftarrow \text{Classify}(\mathcal{D})$ 15: return  $\bar{V}$ 

true value function. We observe that the scores V(s) assigns to the successful states are relatively higher than the unsuccessful ones, but they are not exactly the same as the true returns. As such, we can use V(s) as a binary classifier to separate "good" states from "bad" ones, but not as a reliable predictor of the true returns.

To use V as a binary classifier of the state space, we first need to select a threshold  $\bar{V}$ . For a given policy, separating successful states from unsuccessful ones can be done as follows. First, we compute the average of true return,  $\bar{R}$ , from a set of sampled rollouts that do not terminate due to failure conditions. Second, we compare the true return of a given state to  $\bar{R}$  to determine whether it is a successful state (successful if the return is above  $\bar{R}$ ). In practice, however, we can only obtain an approximated return of a state via V(s) during policy learning and execution. Our goal is then to find the optimal  $\bar{V}$  such that the separation of approximated returns by  $\bar{V}$  is as close as possible to the separation of true returns by  $\bar{R}$ .

Algorithm algorithm 4 summarizes the procedure to compute  $\bar{V}$ . Given a policy  $\pi$ , an approximated value function V, termination conditions  $\mathcal{T}$ , a Gaussian initial state distribution  $\rho = N(\mu_{\rho}, \Sigma_{\rho})$ , and the objective function of the MDP J, we first draw M states

from an expanded initial state,  $N(\mu_{\rho}, 1.5\Sigma_{\rho})$  (Line Line 3), generate rollouts from these sampled states using  $\pi$  (Line Line 4), and compute the true return of each rollout using J(Line Line 5). Because the initial states are drawn from an inflated distribution, we obtain a mixed set of successful and unsuccessful rollouts. We then compute the average return of successful rollouts  $\bar{R}$  that do not terminate due to the terminal conditions  $\mathcal{T}$  (Line Line 6). Next, we generate the training set where each data point is a pair of the predicted value  $V(\mathbf{s}_i)$  and a binary classification label, "good" or "bad", according to  $\bar{R}$  (i.e.  $R_i > \bar{R}$ means  $\mathbf{s}_i$  is good) (Line Line 7-Line 10). We then train a binary classifier represented as a decision tree to find to find  $\bar{V}$  (Line Line 14).

#### 3.3.5 Applying Relay Networks

Once the graph of relay networks  $\Gamma$  is trained, applying the polices is quite straightforward. For a given initial state s, we select a node c whose V(s) has the highest value among all nodes. We execute the current policy  $\pi_c$  until it reaches a state where the value of the parent node is greater than its threshold  $(V_{p(c)}(s) > \overline{V}_{p(c)})$ , where p(c) indicates the index of the parent node of c. At that point, the parent control policy takes over and the process repeats until we reach the root policy. Alternatively, instead of always switching to the parent policy, we can switch to another policy whose V(s) has the highest value.

# 3.3.6 Extending to Multiple Strategies

Optionally, our algorithm can be extended to discover multiple solutions to the original MDP. Take the problem of cartpole swing-up and balance as an example. In Algorithm algorithm 3, if we choose the mean of  $\rho_0$  to be slightly off to the left from the balanced goal position, we will learn a chain of relay networks that often swings to the left and approaches the balanced position from the left. If we run Algorithm algorithm 3 again with the mean of  $\rho_0$  leaning toward right, we will end up learning a different chain of polices that tends to swing to the right. For a problem with multi-modal solutions, we extend

Algorithm algorithm 3 to solve for a directed graph with multiple chains and describe an automatic method to merge the current chain into an existing one to improve sample efficiency. Specifically, after the current node  $\mathcal{N}_k$  is added to  $\Gamma$  and the next initial state distribution  $\rho_{k+1}$  is proposed (Line Line 12 in Algorithm algorithm 3), we compare  $\rho_{k+1}$ against the initial state distribution stored in every node on the existing chains (excluding the current chain). If there exists a node  $\tilde{\mathcal{N}}$  with a similar initial state distribution, we merge the current chain into  $\tilde{\mathcal{N}}$  by learning a policy (and a value function) that relays to  $\mathcal{N}_k$  from the initial state distribution of  $\tilde{\mathcal{N}}$ , essentially adding an edge from  $\tilde{\mathcal{N}}$  to  $\mathcal{N}_k$  and terminating the current chain. Since now  $\tilde{\mathcal{N}}$  has two parents, it can choose which of the two policies to execute based on the value function or by chance. Either path will lead to completion of the task, but will do so using different strategies.

# 3.3.7 Results

We evaluate our algorithm on motor skill control problems in simulated environments. We use DART physics engine [134] to create five learning environments similar to Cartpole, Hopper, 2D Walker, and Humanoid environments in Open-AI Gym [135]. To demonstrate the advantages of relay networks, our tasks are designed to be more difficult than those in Open-AI Gym. Implementation details can be found in the supplementary document. We compare our algorithm to three baselines:

- A single policy (ONE): ONE is a policy trained to maximize the objective function of the original task from the initial state distribution *ρ*. For fair comparison, we train ONE with the same number of samples used to train the entire relay networks graph. ONE also has the same number of neurons as the sum of neurons used by all relay policies.
- No relay (NR): NR validates the importance of relaying which amounts to the second term of the objective function in Equation Equation 3.11 and the terminal condition, V<sub>k</sub>(s) > V<sub>k</sub>. NR removes these two treatments, but otherwise identical to relay

networks. To ensure fairness, we use the same network architectures and the same amount of training samples as those used for relay networks. Due to the absence of the relay reward and the terminal condition  $V_k$  (s >  $\bar{V}_k$ , each policy in NR attempts to achieve the original task on its own without relaying. During execution, we evaluate every value function at the current state and execute the policy that corresponds to the highest value.

Curriculum learning (CL): We compare with curriculum learning which trains a single policy with increasingly more difficult curriculum. We use the initial state distributions computed by Algorithm algorithm 3 to define different curriculum. That is, we train a policy to solve a sequence of MDP's defined as {S, A, r, T, ρ<sub>k</sub>, P}, k ∈ [0, K], where K is the index of the last node on the chain. When training the next MDP, we use previously solved π and V to "warm-start" the learning.

# 3.3.8 Tasks

We will briefly describe each task in this section. Please see Appendix B in the supplementary document for detailed description of the state space, action space, reward function, termination conditions, and initial state distribution for each problem.

- **Cartpole:** Combining the classic cartpole balance problem with the pendulum swingup problem, this example trains a cartpole to swing up and balance at the upright position. The mean of initial state distribution,  $\mu_{\rho}$ , is a state in which the pole points straight down and the cart is stationary. Our algorithm learns three relay policies to solve the problem.
- Hopper: This example trains a 2D one-legged hopper to get up from a supine position and hop forward as fast as possible. We use the same hopper model described in Open AI Gym. The main difference is that μ<sub>ρ</sub> is a state in which the hopper lies flat on the ground with zero velocity, making the problem more challenging than the

one described in OpenAI Gym. Our algorithm learns three relay policies to solve the problem.

- 2D walker with initial push: The goal of this example is to train a 2D walker to overcome an initial backward push and walk forward as fast as possible. We use the same 2D walker environment from Open AI Gym, but modify μ<sub>ρ</sub> to have a negative horizontal velocity. Our algorithm learns two relay policies to solve the problem.
- 2D walker from supine position: We train the same 2D walker to get up from a supine position and walk as fast as it can. μ<sub>ρ</sub> is a state in which the walker lies flat on the ground with zero velocity. Our algorithm learns three relay policies to solve the problem.
- Humanoid walks: This example differs from the rest in that the subtasks are manually designed and the environment constraints are modified during training. We train a 3D humanoid to walk forward by first training the policy on the sagittal plane and then training in the full 3D space. As a result, the first policy is capable of walking forward while the second policy tries to brings the humanoid back to the sagittal plane when it starts to deviate in the lateral direction. For this example, we allow the policy to switch to non-parent node. This is necessary because while walking forward the humanoid deviates from the sagittal plane many times.

#### 3.3.9 Baselines Comparisons

We compare two versions of our algorithm to the three baselines mentioned above. The first version (AUTO) is exactly the one described in Algorithm algorithm 3. The second version (MANUAL) requires the user to determine the subtasks and the initial state distributions associated with them. While AUTO presents a cleaner algorithm with less user intervention, MANUAL offers the flexibility to break down the problem in a specific way to incorporate domain knowledge about the problem.



Figure 3.19: Testing curve comparisons.

Figure Figure 3.19 shows the *testing curves* during task training. The learning curves are not informative for comparison because the initial state distributions and/or objective functions vary as the training progresses for AUTO, MANUAL, NR, and CL. The testing curves, on the other hand, always computes the average return on the original MDP. That is, the average objective value (Equation Equation 3.9) of rollouts drawn from the original initial state distribution  $\rho$ . Figure Figure 3.19 indicates that while both AUTO and MAN-UAL can reach higher returns than the baselines, AUTO is in general more sample efficient than MANUAL. Further, training the policy to relay to the "good states" is important as demonstrated by the comparison between AUTO and NR. The results of CL vary task by task, indicating that relying learning from a warm-started policy is not necessarily helpful.

#### 3.3.10 Analyses

**Relay reward:** One important parameter in our algorithm is the weight for relay reward, i.e.  $\alpha$  in Equation Equation 3.11. Figure 3.20(a) shows that the policy performance is not sensitive to  $\alpha$  as long as it is sufficiently large.

Accuracy of value function: Our algorithm relies on the value function to make accurate binary prediction. To evaluate the accuracy of the value function, we generate 100 rollouts using a learned policy and label them negative if they are terminated by the termination conditions  $\mathcal{T}$ . Otherwise, we label them positive. We then predict the rollouts positive if they satisfy the condition,  $V(\mathbf{s}) > \overline{V}$ . Otherwise, they are negative. Figure 3.20(c) shows the confusion matrix of the prediction. In practice, we run an additional regression on the value function after the policy training is completed to further improve the consistency be-



Figure 3.20: (a) The experiment with  $\alpha$ . (b) Comparison of ONE with different numbers of neurons. (c) Confusion matrix of value function binary classifier (d) Confusion matrix after additional regression.

tween the value function and the final policy. This additional step can further improve the accuracy of the value function as a binary predictor (Figure 3.20(d)).

# 3.3.11 Conclusion

We propose a technique to learn a robust policy capable of controlling a wide range of state space by breaking down a complex task to simpler subtasks. Our algorithm has a few limitations. The value function is approximated based on the visited states during training. For a state that is far away from the visited states, the value function can be very inaccurate. Thus, the initial state distribution of the child node cannot be too far from the parent's initial state distribution. In addition, as mentioned in the introduction, the relay networks are built on locally optimal policies, resulting globally suboptimal solutions to the original task. The theoretical bounds of the optimality of relay networks can be an interesting future direction.

# CHAPTER 4 FALL PREVENTION USING ASSISTIVE DEVICES

#### 4.1 Motivation

More than three million older adults every year in the United States are treated for fall injuries. In 2015, the medical costs for falls amounted to more than \$50 billion. Compounding to the direct injuries, fall-related accidents have long-lasting impact because falling once doubles one's chances of falling again. Even with successful recovery, many older adults develop fear of falling, which may make them reduce their everyday activities. When a person is less active, their health condition plummets which increases their chances of falling again.

Designing a control policy to prevent falls on an existing wearable robotic device has multiple challenges. First, the control policy needs to run in real-time with limited sensing and actuation capabilities dictated by the walking device. Second, a large dataset of human falling motions is difficult to acquire and unavailable to public to date, which imposes fundamental obstacles to learning-based approaches. Lastly and perhaps most importantly, the development and evaluation of the fall-prevention policy depends on intimate interaction with human users. The challenge of modeling realistic human behaviors in simulation is daunting, but the risk of testing on real humans is even greater.

We tackle these issues by taking the approach of model-free reinforcement learning (RL) in simulation to train a fall-prevention policy that operates on the walking device in real-time, as well as to model the human locomotion under disturbances. The model-free RL is particularly appealing for learning a fall-prevention policy because the problem involves non-differentiable dynamics and lacks existing examples to imitate. In addition, demonstrated by recent work in learning policies for human motor skills [136, 137], the

model-free RL provides a simple and automatic approach to solving under-actuated control problems with contacts, as is the case of human locomotion. To ensure the validity of these models, we compare the key characteristics of human gait under disturbances to those reported in the biomechanics literature [17, 69].

Specifically, we propose a framework to automate the process of developing a **fall predictor** and a **recovery policy** on an assistive walking device, by only utilizing the on-board sensors and actuators. When the fall predictor predicts that a fall is imminent based on the current state of the user, the recovery policy will be activated to prevent the fall and deactivated when the stable gait cycle is recovered. The core component of this work is a **human walking policy** that is robust to a moderate level of perturbations. We use this human walking policy to provide training data for the fall predictor, as well as to teach the recovery policy how to best modify the person's gait to prevent falling.

Our evaluation shows that the human policy generates walking sequences similar to the real-world human walking data both with and without perturbation. We also show quantitative evaluation on the stability of the recovery policy against various perturbation. In addition, our method provides a quantitative way to evaluate the design choices of assistive walking device. We analyze and compare the performance of six different configurations of sensors and actuators, enabling the engineers to make informed design decisions which account for the control capability prior to manufacturing process.

# 4.2 Method

We propose a framework to automate the process of augmenting an assistive walking device with the capability of fall prevention. Our method is built on three components: a human walking policy, a fall predictor, and a recovery policy. We formulate the problem of learning human walking and recovery policies as Markov Decision Processes (MDPs),  $(S, A, T, r, p_0, \gamma)$ , where S is the state space, A is the action space, T is the transition function, r is the reward function,  $p_0$  is the initial state distribution and  $\gamma$  is a discount factor. We take the approach of model-free reinforcement learning to find a policy  $\pi$ , such that it maximizes the accumulated reward:

$$J(\pi) = \mathbb{E}_{\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_T} \sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t),$$

where  $\mathbf{s}_0 \sim p_0$ ,  $\mathbf{a}_t \sim \pi(\mathbf{s}_t)$  and  $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t)$ .

We denote the human walking policy as  $\pi_h(\mathbf{a}_h|\mathbf{s}_h)$  and the recovery policy as  $\pi_e(\mathbf{a}_e|\mathbf{s}_e)$ , where  $\mathbf{s}_h$ ,  $\mathbf{a}_h$ ,  $\mathbf{s}_e$ , and  $\mathbf{a}_e$ , represent the corresponding states and actions, respectively. Our method can be applied to assitive walking devices with any sensors or actuators, though we assume that the observable state  $\mathbf{s}_e$  of the walking device is a subset of the full human state  $\mathbf{s}_h$  due to sensor limitations. Since our method is intended to augment an assistive walking device, we also assume that the user who wears the device is capable of walking. Under such an assumption, our method only needs to model normal human gait instead of various pathological gaits. Our framework can be applied to any kind of external perturbation which causes a fall. We evaluate our algorithm with pushes to the pelvis, although falls in the real-world are typically not caused by pushes but rather by slips or trips, data for validating our human policy is more readily available for pushing, such as [69] and [138]. we consider this as a first-step towards a more general recovery policy.



Figure 4.1: Left : We model a 29-Degree of Freedom(DoF) humanoid and the 2-DoF exoskeleton in PyDart. Right : Assistive device design used in our experiments.

# 4.2.1 Human Walking Policy

We take the model-free reinforcement learning approach to develop a human locomotion policy  $\pi_h(\mathbf{a}_h|\mathbf{s}_h)$ . To achieve natural walking behaviors, we train a policy that imitates the human walking reference motion similar to Peng *et al.* [136]. The human 3D model (agent) consists of 23 actuated joints with a floating base as shown in Figure Figure 4.1. This gives rise to a 53 dimensional state space  $\mathbf{s}_h = [\mathbf{q}, \dot{\mathbf{q}}, \mathbf{v}_{com}, \boldsymbol{\omega}_{com}, \phi]$ , including joint positions, joint velocities, linear and angular velocities of the center of mass (COM), and a phase variable that indicates the target frame in the motion clip. We model the intrinsic sensing delay of a human musculoskeletal system by adding a latency of 40 milliseconds to the state vector before it is fed into the policy. The action determines the target joint angles  $\mathbf{q}_t^{target}$  of the proportional-derivative (PD) controllers, deviating from the joint angles in the reference motion:

$$\mathbf{q}_t^{target} = \hat{\mathbf{q}}_t(\phi) + \mathbf{a}_t, \tag{4.1}$$

where  $\hat{\mathbf{q}}_t(\phi)$  is the corresponding joint position in the reference motion at the given phase  $\phi$ . Our reward function is designed to imitate the reference motion:

$$r_h(\mathbf{s}_h, \mathbf{a}_h) = w_q(\mathbf{q} - \hat{\mathbf{q}}(\phi)) + w_c(\mathbf{c} - \hat{\mathbf{c}}(\phi)) + w_e(\mathbf{e} - \hat{\mathbf{e}}(\phi)) - w_\tau ||\boldsymbol{\tau}||^2, \quad (4.2)$$

where  $\hat{\mathbf{q}}$ ,  $\hat{\mathbf{c}}$ , and  $\hat{\mathbf{e}}$  are the desired joint positions, COM positions, and end-effector positions from the reference motion data, respectively. The reward function also penalizes the magnitude of torque  $\boldsymbol{\tau}$ . We use the same weight  $w_q = 5.0$ ,  $w_c = 2.0$ ,  $w_e = 0.5$ , and  $w_{\tau} = 0.005$ for all experiments. We also use early termination of the rollouts, if the agent's pelvis drops below a certain height or if the base rotates about any axis beyond a threshold, we end the rollout and re-initialize the state.

Although the above formulation can produce control policies that reject small disturbances near the target trajectory, they often fail to recover from perturbations with larger magnitude, such as those encountered during locomotion. It is critical to ensure that our human walking policy can withstand the same level of perturbation as a capable real person, so that we can establish a fair baseline to measure the increased stability due to our recovery policy.

Therefore, we exert random forces to the agent during policy training. Each random force has a magnitude uniformly sampled from [0, 800] N and a direction uniformly sampled from  $[-\pi/2, \pi/2]$ , applied for 50 milliseconds on the agent's pelvis in parallel to the ground. The maximum force magnitude induces a velocity change of roughly 0.6m/sec. This magnitude of change in velocity is comparable to experiments found in literature such as [69],[138] and [70]. We also randomize the time when the force is applied within a gait cycle. Training in such a stochastic environment is crucial for reproducing the human ability to recover from a larger disturbance during locomotion. We represent a human policy as a multi-layered perceptron (MLP) neural network with two hidden layers of 128 neurons each. The formulated MDP is trained with Proximal Policy Optimization (PPO) [PPO].

# 4.2.2 Fall Predictor

Being able to predict a fall before it happens gives the recovery policy critical time to alter the outcome in the near future. We take a data-driven approach to train a classifier capable of predicting the probability of the fall in the next 40 milliseconds. Collecting the training data from the real world is challenging because induced human falls can be unrealistic and dangerous/tedious to instrument. As such, we propose to train such a classifier using only simulated human motion. Our key idea is to automatically label the future outcome of a state by leveraging the trained human policy  $\pi_h$ . We randomly sample a set of states  $s_h$  and add random perturbations to them. By following the policy  $\pi_h$  from each of the sampled states, we simulate a rollout to determine whether the state leads to successful recovery or falling. We then label the corresponding state observed by the walking device, ( $s_e$ , 1) if succeeds, or ( $s_e$ , 0) if fails. We collect about 50000 training samples. Note that the input of the training data corresponds to the state of the walking device, not the full state of human, as the classifier will only have access to the information available to the onboard sensors.

We train a support vector machine (SVM) classifier with radial basis function kernel to predict if a state  $s_e$  has a high chance of leading to a fall or not. We perform a six-fold validation test on the dataset and the classifier achieves an accuracy above 94%.

# 4.2.3 Recovery Policy

The recovery policy aims to utilize the onboard actuators of the assistive walking device to stabilize the gait such that the agent can continue to walk uninterruptedly. The recovery policy  $\pi_e$  is trained to provide optimal assistance to the human walking policy when a fall is detected. The state of  $\pi_e$  is defined as  $\mathbf{s}_e = [\dot{\boldsymbol{\omega}}, \boldsymbol{\omega}, \mathbf{q}_{hip}, \dot{\mathbf{q}}_{hip}]$ , which comprises of angular acceleration, angular velocity, and hip joint angle position and velocity, amounting to 10 dimensional state space. We envision to measure these quantities through an Inertial Measurement Unit(IMU) and a motor encoder at the hip. IMU outputs angular velocity and we would compute angular acceleration using finite difference method. Although, IMU provides us with linear acceleration as well, in practice we observed it has very noisy readings, hence we decided to omit it in our state. The action space consists of torques at two hip joints. The reward function maximizes the quality of the gait while minimizing the impact of disturbance:

$$r_e(\mathbf{s}_h, \mathbf{a}_e) = r_{walk}(\mathbf{s}_h) - w_1 \|\mathbf{v}_{com}\| - w_2 \|\boldsymbol{\omega}_{com}\| - w_3 \|\mathbf{a}_e\|,$$
(4.3)

where  $r_{walk}$  evaluates walking performance using Equation Equation 5.2 except for the last term, and  $\mathbf{v}_{com}$  and  $\boldsymbol{\omega}_{com}$  are the global linear and angular velocities of the pelvis. We use the same weight  $w_1 = 2.0$ ,  $w_2 = 1.2$  and  $w_3 = 0.001$  for all our experiments. Note that the input to the reward function includes the full human state  $\mathbf{s}_h$ . While the input to the recovery policy  $\pi_e$  should be restricted by the onboard sensing capability of the assistive
walking device, the input to the reward function can take advantage of the full state of the simulated world, since the reward function is only needed at training time. The policy is represented as a MLP neural network with two hidden layers of 64 neurons each and trained with PPO.

#### 4.2.4 Results

We validate the proposed framework using the open-source physics engine DART [139]. Our human agent is modeled as an articulated rigid body system with 29 degrees of freedom (dofs) including the six dofs for the floating base. The body segments and the mass distribution are determined based on a 50th percentile adult male in North America. We select the prototype of our assistive walking device as the testbed.Similar prototypes are described in [**ExoDesign**, 140, 141]. It has two cable-driven actuators at hip joints, which can exert about 200 Nm at maximum. However, we limit the torque capacity to 30, beyond this value, the torque saturates. Sensors, such as Inertial Measurement Units (IMU) and hip joint motor encoders, are added to the device. We also introduce a sensing delay of 40 to 50 ms. We modeled the interaction between the device and human by adding positional constraints on the thigh and anchor points. For all experiments, the simulation time step is set to 0.002s.

We design experiments to systematically validate the learned human behaviors and effectiveness of the recovery policy. Particularly, our goal is to answer the following questions:

- 1. How does the motion generated by the learned human policy compare to data in the biomechanics literature?
- 2. Does the recovery policy improve the robustness of the gaits to external pushes?
- 3. How does the effectiveness of the assistive walking device change with design choices?



Figure 4.2: Comparison between hip and knee joint angles during walking generated by the policy and human data [17].

#### 4.2.5 Comparison of Policy and Human Recovery Behaviors

We first validate the steady walking behavior of the human policy by comparing it to the data collected from human-subject experiments. Figure Figure 4.2 shows that the hip and knee joint angles generated by the walking policy well match the data reported in Winter *et al.* [17]. We also compare the "torque loop" between the gait generated by our learned policy and the gait recorded from the real world [17]. A torque loop is a plot that shows the relation between the joint degree of freedom and the torque it exerts, frequently used in the biomechanics literature as a metric to quantify human gait. Although the torque loops in Figure **??** are not identical, both trajectories form loops during a single gait cycle indicating energy being added and removed during the cycle. We also notice that the torque range and the joint angle range are similar.

In addition, we compare adjusted footstep locations due to external perturbations with the studies reported by Wang *et al.* [69]. Their findings strongly indicate that the COM dynamics is crucial in predicting the step placement after disturbance that leads to a balanced state. They introduced a model to predict the changes in location of the foot placement of a normal gait as a function of the COM velocity. Figure Figure 4.3 illustrates the foot placements of our model and the model of Wang *et al.* against four pushes with different



Figure 4.3: (a) Comparison of torque loops of a typical trajectory generated by our policy and human data reported by [17] at the hip of stance leg during a gait cycle. The green dots indicate the start and the black dots indicate 50% of the gait cycle. The arrows show the progression of the gait from 0% to 100%. (b) Comparison of the forward foot step locations predicted by the policy and by the model reported by Wang *et al.* [69] as a function of the COM velocity.

magnitudes in the sagittal plane. For all scenarios, the displacement error is below 4 cm.

# 4.2.6 Effectiveness of Recovery Policy



Figure 4.4: Four different timing of the left leg swing phase during which we test the performance of the assistive device. First is at 10% of the phase and then subsequently 30%, 60% and 90% of the left swing leg.

We test the performance of the learned recovery policy in the simulated environment with external pushes. As a performance criterion, we define the *stability region* as a range of external pushes from which the policy can return to the steady gait without falling. For better 2D visualization, we fix the pushes to be parallel to the plane, applied on the same location with the same timing and duration (50 milliseconds). All the experiments in this section use the default sensors and actuators provided by the prototype of the walking



Figure 4.5: Stability region with and without the use of a recovery policy. A larger area shows increased robustness to an external push in both magnitude and direction.

device: an IMU, hip joint motor encoders, and hip actuators that control the flexion and extension of the hip.

Figure Figure 4.5 compares the stability region with and without the learned recovery policy. The area of stability region is expanded by 35% when the recovery policy is used. Note that the stability region has very small coverage on the negative side of y-axis which corresponds to the rightward forces. This is because we push the agent when the swing leg is the left one, making it difficult to counteract the rightward pushes. Figure Figure 4.7 shows one example of recovery motion.

The timing of the push in a gait cycle has a great impact on fall prevention. We test our recovery policy with perturbation applied at four different phases during the swing phase (Figure Figure 4.4). We found that the stability region is the largest when the push is applied at 30% of the swing phase and the smallest at 90% (Figure Figure 4.6, Top). This indicates that the perturbation occurs right before heel strike is more difficult to recover than the one occurs in early swing phase possibly due to the lack of time to adjust the foot location. The difference in the stability region is approximately 28%. The bottom of Figure Figure 4.6 shows the impact of the perturbation timing on COM velocity over four gait cycles. The results echo the previous finding as it shows that the agent fails to return to the steady state when the perturbation occurs later in the swing phase.



Figure 4.6: Comparison of recovery performance when perturbation is applied at four different phases. **Top:** Comparison of stability region. **Bottom:** Comparison of COM velocity across five gait cycles. Perturbation is applied during the gait cycle 'p'. The increasing velocity after perturbation indicates that our policy is least effective at recovering when the perturbation occurs later in the swing phase.



Figure 4.7: **Top:** Successful gait with an assistive device. **Bottom:** Unsuccessful gait without an assistive device. Torques are set to zero.



Figure 4.8: Average torques at the hip joints from 50 trials with various perturbations. The shaded regions represent the 3-sigma bounds. **Red:** Joint torques exerted by the human and the recovery policy. **Blue:** Human joint torques without a recovery policy. **Green:** Torques produced by a recovery policy.

We also compare the generated torques with and without the recovery policy when perturbation is applied. Figure Figure 4.8 shows the torques at the hip joint over the entire gait cycle (not just swing phase). We collect 50 trajectory for each scenario by applying random forces ranging from 200N to 800N at the fixed timing of 30% of the gait cycle. The results show that hip torques exerted by the human together with the recovery policy do not change the overall torque profile significantly, suggesting that the recovery policy makes minor modification to the torque profile across the remaining gait cycle, instead of generating a large impulse to stop the fall. We also show that the torque exerted by the recovery policy never exceeds the actuation limits of the device.

## 4.2.7 Evaluation of Different Design Choices

Our method can be used to inform the selection of sensors and actuators when designing a walking device with the capability of fall prevention. We test two versions of actuators: the 2D hip device can actuate the hip joints only in the sagittal plane while the 3D device also allows actuation in the frontal plane. We also consider three different configurations of sensors: an inertial measurement unit (IMU) that provides the COM velocity and ac-



Figure 4.9: Stability region for six policies trained with three sensor configurations and two actuator configurations.

celeration, a motor encoder that measures hip joint angles, and the combination of IMU and motor encoder. In total, we train six different recovery policies with three sensory inputs and two different actuation capabilities. For each sensor configuration, we train a fall predictor using only sensors available to that configuration.

Figure Figure 4.9 shows the stability region for each of the six design configurations. The results indicate that 3D actuation expands the stability region in all directions significantly comparing to 2D actuation, even when the external force lies on the sagittal plane. We also found that the IMU sensor plays a more important role than the motor encoder, which suggests that COM information is more critical than the hip joint angle in informing the action for recovery. The recovery policy performs the best when combining the IMU and the joint encoder, as expected.

# 4.2.8 Conclusion

We presented an approach to automate the process of augmenting an assistive walking device with ability to prevent falls. Our method has three key components : A human walking policy, fall predictor and a recovery policy. In a simulated environment we showed that an assistive device can indeed help recover balance from a wider range of external perturbations. We introduced *stability region* as a quantitative metric to show the benefit

of using a recovery policy. In addition to this, *stability region* can also be used to analyze different design choices for an assistive device. We evaluated six different sensor and actuator configurations.

In this work, we only evaluated the effectiveness of using a recovery policy for an external push. It would be interesting to extend our work to other kinds of disturbances such as tripping and slipping. Another future direction we would like to take is deploying our recovery policy on the real-world assistive device. This would need additional efforts to make sure that our recovery policy also can adjust for the differences in body structure of users.

# CHAPTER 5 ERROR-AWARE POLICY LEARNING

#### 5.1 Motivation

Humans exhibit remarkably large variations in gait characteristics during steady-state walking as well as push recovery. Due to this, assistive device controllers tuned for one individual often do not work well when tested on another person. Typically, fine tuning a controller for each person has been the go-to approach thus far, but this process can be a very time consuming ,tedious and unsafe for the user.

In this work, we introduce a novel approach to tackle sim-to-real problems in which the environment dynamics has high variance and is highly unobservable. While our approach is motivated by physical assistive robotic applications, the method can be applied to other tasks in which many dynamic parameters are challenging to model. We propose to train a policy explicitly aware of the effect of unobservable factors during training, called an Error-Aware policy (EAP). Akin to the high-level idea of meta learning, we divide the dynamical environments into training and validation sets and "emulate" a reality gap in simulation. Instead of estimating the model parameters that give rise to the emulated reality gap, we train a function that predicts the deviation (i.e. error) of future states due to the emulated reality gaps. Conditioned on the error predictions, the error-aware policies (EAPs) can learn to overcome the reality gap, in addition to mastering the task.

The main application in this work is to learn an error-aware policy for assistive device control, such as a hip-exoskeleton that helps the user to recover balance during locomotion. From biomechanical data of human gait, we model multiple virtual human walking agents, each varying in physical characteristics as well as parameters that affect the dynamics such as joint damping, torque limits, ground friction, and sensing and actuation delay. We then train a single policy on this group of human agents and show that that the learned EAP works effectively when tested on a different human agent without needing additional data. We extend the prior work, [16], that trained a control policy for push-recovery assistive device for just one simulated human agent, and develop an algorithm that enables the learned policy to transfer to other human agents with unseen biomechanical characteristics.

We evaluate our approach on assistive wearable device by quantifying the stability and gait characteristics generated by an unseen human agent wearing the device with the trained EAP. We present a comprehensive study of the benefits of our approach over prior zero-shot methods such as universal policy (UP) and domain randomization (DR). We also provide results on some standard RL environments, such as Cartpole, Hopper, 2D walker and a quadrupedal robot.

#### 5.2 Method

We present a method to achieve the zero-shot transfer of control policies in partially observable dynamical environments. We consider robotic systems and environments with unobservable or unmeasurable model parameters, which make building accurate simulation models difficult.

We present a novel policy architecture, an Error-Aware Policy (EAP), that is explicitly aware of errors induced by unobservable dynamics parameters and self-corrects its actions according to the errors. An EAP takes the current state, observable dynamic parameters, and predicted errors as inputs and generates corrected actions. We learn an additional errorprediction function that outputs the expected error. Both the error-aware policy and the error-prediction function, are iteratively learned using model-free reinforcement learning and supervised learning.



Figure 5.1: Overview of An Error-aware Policy (EAP). An EAP takes the "expected" future state error as an additional input. The expected error is predicted based on the current state s, observable parameters  $\mu$ , and an uncorrected action a that assumes zero error.

## 5.2.1 Problem Formulation

We formulate the problem as Partially Observable Markov Decision Processes (PoMDPs),  $(S, O, A, P, R, \rho_0, \gamma)$ , where S is the state space, O is the observation space, A is the action space, P is the transition function, R is the reward function,  $\rho_0$  is the initial state distribution and  $\gamma$  is a discount factor. In our formulation, we make a clear distinction between observable model parameters  $\mu$  and unobservable parameters  $\nu$  of the agent and environment. Observable quantities are parameters that can be easily measured such as masses or link lengths, whereas unobserved quantities are challenging to estimate, such as circuit dynamics or backlash. Therefore, both  $\mu$  and  $\nu$  affect the transition function  $P(\mathbf{s}'|\mathbf{a}, \mathbf{s}, \mu, \nu)$ . Since we can configure our simulator with both  $\mu$  and  $\nu$ , we can randomly sample  $\mu$  and  $\nu$ and create a list of K different environments  $\mathbf{D} = \{(\mu_0, \nu_0), (\mu_1, \nu_1), \cdots, (\mu_K, \nu_K)\}$ , but it is hard to obtain  $\nu$  at testing time. In this case, the transition function will be abbreviated as  $P(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \mu)$ .

Instead of estimating the values of unobserved quantities, we capture the effect of these parameters by defining a metric called a *state-error*. When transferring from one environment to another, the action a applied at a given state s will produce different next states due to the differences in both  $\mu$  and  $\nu$ , in other words, a state-error.

We hypothesize that a policy which is explicitly aware of the state-error would be able to make better decisions by self-correcting its action. We call this an error-aware policy  $\pi(\mathbf{a}|\mathbf{s}, \boldsymbol{\mu}, \mathbf{e})$  (EAP), which takes in observable parameters  $\boldsymbol{\mu}$  as well as the "expected" future state error in a new environment e as input (Figure Figure 5.1).

We present a novel training methodology using model-free reinforcement learning that involves learning two functions: an error-aware policy and an error prediction function. First, we learn an error-aware policy that takes the output of error prediction function Eas an input and has the ability to generalize to novel environments in a zero-shot manner. Simultaneously, we learn an error-prediction function, which takes as inputs the state s, an uncorrected action a and observable parameters  $\mu$ , and outputs the expected state error e when a policy trained in one environment is deployed to a different one  $E : (s, a, \mu) \mapsto \mathbb{R}^n$ . We will discuss more details of training in the following sections.

Algorithm 5: Train an Error Aware Policy.
1: Input: Environments $\mathbf{D} = \{(\boldsymbol{\mu}_0, \boldsymbol{\nu}_0), \cdots, (\boldsymbol{\mu}_K, \boldsymbol{\nu}_K)\}$
2: Pre-train $\pi(\mathbf{a} \mathbf{s}, \boldsymbol{\mu}_0, \mathbf{e} = 0)$ for $P(\mathbf{s}' \mathbf{s}, \mathbf{a}, \boldsymbol{\mu}_0)$ reference environment with $\mathbf{e} = 0$
3: while not done do
4: Sample an environment with $(\boldsymbol{\mu}, \boldsymbol{\nu})$ from <b>D</b>
5: <b>for</b> each policy update iteration <b>do</b>
6: Initialize buffer $\mathbf{B} = \{\}$
7: Update an error function $E$ using Algorithm algorithm 7
8: $\mathbf{B}$ = Generate rollouts using Algorithm algorithm 6
9: Update policy $\pi$ using <b>B</b> with PPO.
10: end for
11: end while
12: return $\pi(\mathbf{a},   \mathbf{s}, \boldsymbol{\mu}, \mathbf{e})$

# 5.2.2 Training an Error-aware Policy

**Training Procedure.** The training process of an error-aware policy is summarized in Algorithm algorithm 5. Assume that we have an oracle error function  $E(\mathbf{s}, \mathbf{a}, \boldsymbol{\mu})$  that outputs the expected state error in a novel environment, which will be explained in the following section. First, the policy is pre-trained to achieve the desired behavior only in the reference environment  $(\boldsymbol{\mu}_0, \boldsymbol{\nu}_0)$  assuming there is no state error,  $\pi(\mathbf{a}|\mathbf{s}, \boldsymbol{\mu}_0, \mathbf{e} = 0)$ . Once the policy is trained in the reference environment, we sample dynamics parameters  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\nu}_i$  (i > 0) uniformly from the data set **D** and evaluate the EAP in this new environment. The

policy parameters are updated using a model-free reinforcement learning algorithm, Proximal Policy Optimization [52]. Sampling new testing environments and updating policy parameters are repeated until the convergence.

# Algorithm 6: Generate Rollouts

1: Input: Observable dynamics parameters $\mu$ , Transition function P, Current policy $\pi$
and error function $E$ , Replay buffer <b>B</b>
2: Sample state s from initial state distribution $\rho_0$
3: while not done do
4: $\mathbf{a} \sim \pi(\mathbf{a} \mathbf{s}, oldsymbol{\mu}, \mathbf{e} = 0)$ // original action
5: $\mathbf{e}=E(\mathbf{s},\mathbf{a},oldsymbol{\mu})$ // predicted error
6: $\hat{\mathbf{a}} \sim \pi(\mathbf{a} \mathbf{s}, oldsymbol{\mu}, \mathbf{e})$ // error-aware action
7: $\mathbf{s}' \sim P(\mathbf{s}'   \mathbf{s}, \hat{\mathbf{a}}, oldsymbol{\mu})$
8: $r = R(\mathbf{s}, \hat{\mathbf{a}})$
9: $B = B \cup \{(\mathbf{s}, \hat{\mathbf{a}}, r, \mathbf{s}', \boldsymbol{\mu})\}$
10: $\mathbf{s} = \mathbf{s}'$
11: end while
12: return B

**Rollout Generation.** A roll-out generation procedure is described in Algorithm algorithm 6. Given a state s in this environment  $(\mu, \nu)$ , we query an action from policy  $\pi$  as if the policy is being deployed in the reference environment with e = 0. This action a is fed into the error function E which predicts the expected state error in this environment, then the state error is passed into the error-aware policy to query a corrected action  $\hat{a}$  which will be applied to the actual system. The task reward R(s, a) guides the policy optimization to find the best "corrected" action that maximizes the reward.

### 5.2.3 Training an Error Function

In reality, we do not have an oracle error function that can predict the next state due to the lack of unobservable parameters  $\nu$ . To this end, we will learn this function simultaneously with EAP, by splitting the dataset **D** into the training and validation sets. Similar to training methodology followed in meta-learning algorithms, we repeatedly apply the trained policy into sampled environments from the validation set. Because our nominal behavior is pre-

Algor	ithm 7: Train an Error Prediction Function.
1:	<b>Input:</b> Reference environment with $\mu_0$
2:	<b>Input:</b> Target environment with $\mu$
3:	Input: Replay Buffer B
4:	Input: Dataset Z
5:	<b>Input:</b> Error Horizon T
6:	while not done do
7:	Sample the initial state $s_0^0$ from $B$
8:	$\mathbf{s}^0 = \mathbf{s}^0_0$
9:	for $t = 0 : T - 1$ do
10:	<pre>// Simulation in Reference Env</pre>
11:	$\mathbf{a}_0^t \sim \pi(\mathbf{a} \mathbf{s}_0^t, \boldsymbol{\mu}_0, \mathbf{e}=0)$
12:	$\mathbf{s}_0^{t+1} \sim P(\mathbf{s}_0^t, \mathbf{a}_0^t, oldsymbol{\mu}_0)$
13:	<pre>// Simulation in Validation Env</pre>
14:	$\mathbf{a}^t \sim \pi(\mathbf{a} \mathbf{s}^t, \boldsymbol{\mu}, \mathbf{e} = 0)$
15:	$\mathbf{s}^{t+1} \sim P(\mathbf{s}^t, \mathbf{a}^t, oldsymbol{\mu})$
16:	end for
17:	$\mathbf{Z} = \mathbf{Z} \cup \{(\mathbf{s^0}, \mathbf{a^0}, \mathbf{s}^T, \mathbf{s}^T_0, oldsymbol{\mu})\}$
18:	end while
19:	minimize the $L(\phi)$ in Eq. Equation 5.1 using <b>Z</b> .
20:	return $\phi$

trained in the reference environment  $(\mu_0, \nu_0)$ , we compute the errors by measuring the differences in the reference environment  $(\mu_0, \nu_0)$  and the validation environment  $(\mu, \nu)$ :  $e = (\bar{\mathbf{s}}' - \mathbf{s}') \in \mathbb{R}^n$ , generated by two dynamic models  $P(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \mu_0)$  and  $\bar{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \mu)$ .

Horizon of Error Prediction. In practice, we found that the error accumulated during one step is often not sufficient to provide useful information to the EAP. To overcome this challenge, we take the state in the collected trajectory and further simulate it for T steps in both the reference environment  $P(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \boldsymbol{\mu}_0)$  and the validation environment  $\overline{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \boldsymbol{\mu})$ . We provide analysis on the effect of horizon length from T = 1 to T = 8 in the Section ??. Loss Function. Since the differences between the two dynamical environments reflects the reality gap caused by unobservable parameters, the error prediction function E enables us to learn the effect of the unobserved parameters captured through the state error. We train our error prediction function E to learn this "emulated" sim-to-real gap by minimizing the



Figure 5.2: Left : A full state error representation input into the policy vs **Right** : Projected error representation as an input to the policy

following loss:

$$L(\boldsymbol{\phi}) = \sum_{(\mathbf{s}^0, \mathbf{a}^0, \mathbf{s}^T, \mathbf{s}_0^T, \boldsymbol{\mu}) \in \mathbf{Z}} ||E(\mathbf{s}^0, \mathbf{a}^0, \boldsymbol{\mu}) - (\mathbf{s}_0^T - \mathbf{s}^T)||^2,$$
(5.1)

where Z is the collected dataset and  $\phi$  is the parameters or the neural net representing E. Algorithm algorithm 7 summarizes the training procedure.

**Reduced Representations.** We experiment with two different representations of the error input to the policy. First, we input the full state error  $\mathbf{e} = \mathbf{s}_0^T - \mathbf{s}^T$  (with the same dimension as the state) approximated by a MLP neural network, into the policy. Second, we use a network architecture with an information bottle neck, as illustrated in Figure Figure 5.2, and input the latent representation  $\mathbf{e}_p$  into the policy. The same loss function *L* is used to train both the functions.

## 5.3 Results

We design experiments to validate the performance of error-aware policies. We aim to answer the following research questions.

- 1. Does an EAP show better zero-shot transfer on unseen environments compared to the baseline algorithms?
- 2. How does the choice of hyperparameters affect the performance of an EAP?

Task	Observable Params. $\mu$	Unobservable Params. $\nu$	Net. Arch.	Err. Dim. $ \mathbf{e}_p $
Assitive Walking	mass, height, leg length, and foot length	joint damping, max torques, PD gains and delay	(64, 32)	6
Aliengo	PD gains, link masses	sensor delay, joint damping, ground friction	(64, 32)	6
Cartpole	pole length, pole mass, cart mass	joint damping, joint friction	(32, 16)	2
Hopper	thight mass, foot mass, shin length	joint damping, ground friction	(32, 16)	4
Walker 2D	link masses, shin length	sensing delay, joint damping, ground friction	(64, 32)	5

Table 5.1:	Tasks	and Net	work A	rchitectures
------------	-------	---------	--------	--------------

Table 5.2: Ranges of variation for observable parameters during training and testing in the assistive walking task.

Observable Params. $\mu$			
Parameter	Training range	Testing range	
Mass	[45,76] kg	[55,95] kg	
Height	[143,182] cm	[155,197] cm	
Leg-length	[70,88] cm	[80,95] cm	
foot length	[21,24] cm	[24,26] cm	

# 5.3.1 Baseline Algorithms

We compare our method with two baselines commonly used for sim-to-real policy transfer, Domain Randomization (DR)[82, 85] and Universal Policy (UP) [95]. DR aims to learn a more robust policy for zero-shot transfer, by training with randomly sampled dynamics parameters (in our case, both  $\mu$  and  $\nu$ ). UP extends DR by taking dynamics parameters as additional input. UP often transfer to target environments better than DR, but it explicitly requires to know dynamics parameters, where  $\nu$  is assumed to be unobservable in our scenario. We did not compare EAPs against meta-learning algorithms [105, 107, 106], which require additional samples from the validation environment.

# 5.3.2 Tasks

We evaluate the performance of error-aware policies on five different tasks. The first task is about push-recovery of an assistive walking device for simulated humans, inspired by the work of Kumar et al [16]. The second task is locomotion of a quadrupedal robot, Aliengo Explorer[142]. The rest three tasks are CartPole, Hopper, and Walker2D, which are from the OpenAI benchmark suite [143].

Unobservable Params. $\nu$			
Parameter	Training range	Testing range	
Joint damping	[0.3,0.6]	[0.5,0.8]	
Max torques	[120,180]	[155,200]	
PD gains (P,D)	[(500,25),(750,50)]	[(650,30),(800,50)]	
Delay	[30,60] ms	[45,70] ms	

Table 5.3: Ranges of variation for unobservable parameters during training and testing in the assistive walking task.



Figure 5.3: Five different test subjects for the assistive walking experiment with varying height, mass, leg length and foot length from the biomechanical gait dataset [144].

# Assistive walking device for push recovery

In this task, the goal is to learn a policy for an assistive wearable device (i.e. exoskeleton) to help a human recover balance after an external push is applied. (inset figure). We use a hip exoskeleton that applies torques in 2-degrees of freedom at each hip joint. Our algorithm begins by training 15 human agents using public biomechanical gait data [144] to walk in a steady-state gait cycle, similar to the approach presented in [136]. The 15 agents vary in mass, height, leg length, and foot length according to the biomechanical data used to train their corresponding policies, which formulate the four-dimensional observable parameters  $\mu$  (Figure Figure 5.3). We also vary each human agent's joint damping, maximum joint torques, PD gains, and sensory delay as the four dimensional unobservable parameters  $\nu$ . We split the 15 human agents into 10 for the training set and 5 as the testing set.

Human Behavior Modeling. First, we capture the human behavior by training a humanonly walking policy  $\pi_h$  that mimics the reference motion which frames are denoted as  $\bar{\mathbf{q}}$ . Each human model has 23 actuated joints along with a floating base. The state space has 53 dimensions,  $\mathbf{s}_h = [\mathbf{q}, \dot{\mathbf{q}}, \mathbf{v}_{com}, \boldsymbol{\omega}_{com}, \psi]$ , which represent joint positions, joint velocities, linear and angular velocities of the center of mass, and a phase variable  $\psi$  that indicates the target frame in the reference biomechanical gait cycle. The action  $\mathbf{a}$  is defined as the offset to the reference biomechanical joint trajectory  $\bar{\mathbf{q}}(\psi)$ , which results in the target angles:  $\mathbf{q}^{target} = \bar{\mathbf{q}} + \mathbf{a}$ . The reward function encourages mimicking the reference motions from public biomechanical data:

$$R_{human}(\mathbf{s}_h, \mathbf{a}_h) = w_q(\mathbf{q} - \bar{\mathbf{q}}) + w_v(\dot{\mathbf{q}} - \dot{\bar{\mathbf{q}}}) + w_c(\mathbf{c} - \bar{\mathbf{c}}) + w_p(\mathbf{p} - \bar{\mathbf{p}}) - w_\tau ||\boldsymbol{\tau}||^2, \quad (5.2)$$

where the terms include the reference joint positions  $\bar{\mathbf{q}}$ , joint velocities  $\bar{\mathbf{q}}$ , end-effector locations  $\bar{\mathbf{p}}$ , contact flags  $\mathbf{c}$ , and the joint torques  $\tau$ . During training, we exert random forces to the agent during policy training. Each random force has a magnitude uniformly sampled from [0, 800] N and a direction uniformly sampled from  $[-\pi/2, \pi/2]$ , applied for 50 milliseconds on the agent's pelvis in parallel to the ground. The maximum force magnitude induces a velocity change of roughly 0.6 - 0.8 m/sec. This magnitude of change in velocity is comparable to experiments found in biomechanics literature such as [69],[138] and [70]. We also randomize the time when the force is applied within a gait cycle. The forces are applied once at a randomly chosen time in each trajectory rollout. Similar to [16], we enforce joint torque constraints and introduce sensing delays during training to prevent the human agent to adapt to external disturbance really well.

**MDP Formulation.** Once the human agents are trained, we begin learning the pushrecovery EAP for the assistive device. The objective is to stabilize the human gait from external perturbations. The 17 dimensional state of robot is defined as  $\mathbf{s}_e = [\boldsymbol{\omega}, \boldsymbol{\alpha}, \ddot{\mathbf{x}}, \mathbf{q}_{hip}, \dot{\mathbf{q}}_{hip}]$ , which comprises angular velocity, orientation, linear acceleration, hip joint positions, and hip joint velocity. The four dimensional action  $\mathbf{a}_e$  consists of torques at two hip joints. The reward function maximizes the quality of the gait while minimizing the impact of an external push.

$$R_{exo}(\mathbf{s}_{h}, \mathbf{s}_{e}, \mathbf{a}_{e}) = R_{human}(\mathbf{s}_{h}) - w_{1} \|\mathbf{v}_{com}\| - w_{2} \|\boldsymbol{\omega}_{com}\| - w_{3} \|\mathbf{a}_{e}\|,$$
 (5.3)

where  $R_{human}$  is defined in equation Equation 5.2, and  $\mathbf{v}_{com}$  and  $\boldsymbol{\omega}_{com}$  are the global linear and angular velocities of the pelvis. The last term penalizes the torque usage. We use the same weight  $w_1 = 2.0$ ,  $w_2 = 1.2$  and  $w_3 = 0.001$  for all our experiments.

# Quadrupedal Locomotion

In our second task, we learn a control policy that generates a walking motion for a quadrupedal robot, Aliengo Explorer [142]. For this task, the 17 observable parameters ( $\mu$ ) are PD gains of the joints, link and root masses and the 10 unobservable parameters  $\nu$  include sensing delay, joint damping of thigh and knee joints and ground friction. The 39-dimensional state space consists of torso position and orientation and corresponding velocities, joint position and velocities, foot contact variable that indicates when each foot should be in contact with the ground, while the 12-dimensional action space consists of joint velocity targets which is fed into a PD controller that outputs torques to each joint.

The reward function is designed to track the target motion that walks at 0.8 m/s:

$$r(\mathbf{s}, \mathbf{a}) = w_1 e^{-k_1 * (\mathbf{q} - \bar{\mathbf{q}})} + w_2 e^{-k_2 * (\dot{\mathbf{q}} - \bar{\dot{\mathbf{q}}})} + w_3 \min(\dot{x}, 0.8) + \sum_{i=1}^4 ||c_i - \bar{c}_i||^2.$$
(5.4)

In this equation, the first term encourages to track the desired joint positions, the second term is to track the desired joint velocities, the third term is for matching the forward velocity  $\dot{x}$  to a target velocity of 0.8 m/s. and the four term tracks the predefined contact flags. We use the same weight  $k_1 = 35, w_1 = 0.75, k_2 = 2, w_2 = 0.20$ , and  $w_3 = 1.0$  for all experiments.

# **OpenAI** Environments

We test our method on three OpenAI environments: CartPole, Hopper and Walker2D. While using the same state spaces, action spaces, and the reward functions described in the benchmark [143], we additionally define observable and unobservable dynamics parameters as follows:

- Cartpole. Observable parameters μ ∈ R<sup>3</sup> includes the length of the pole, the mass of the pole, and the mass of cart. Unobservable parameters ν ∈ R<sup>3</sup> include the damping at the rotational joint, the friction at the rotational joint, and the friction at the translational joint.
- Hopper. Observable parameters μ ∈ R<sup>3</sup> include the mass of the thigh and foot and the length of the shin bodynode. Unobservable parameters ν ∈ R<sup>3</sup> include joint damping of shin and foot joints and ground friction.
- Walker 2D. Observable parameters μ ∈ R<sup>6</sup> include the masses of thigh and foot for both legs, the mass of pelvis, and the length of shin. Unobservable parameters ν ∈ R<sup>4</sup> include joint damping of foot joints, the delay in observation, and ground friction.

### 5.3.3 Zero-shot Transfer with EAPs

In this section, we compare the zero-shot transfer of error-aware policies against two other baseline algorithms, Domain Randomization (DR) and Universal Policies (UP).

**Learning Curves.** First, we compare the learning curves of the EAP, DR, and UP approaches on four selected tasks in Figure Figure 5.4. We set the same ranges of the observable and unobservable parameters for all three algorithms. In our experience, EAPs learn faster than DR and UP for three tasks, the Hopper, Walker2D, and assistive walking tasks, while showing comparable performance for the quadrupedal locomotion task. Note that,



Figure 5.4: Learning curves for four tasks. The number of samples for EAP include the ones generated for training an error function.

to make the comparison fair to baselines, we also include the samples for training error functions (Algorithm algorithm 7) when we evaluate the performance of EAPs. We do not include the experiment on the CartPole environment for brevity but the EAP outperforms the baselines as well

**Zero-shot Transfer.** Then we evaluate the learned policies on unseen validation environments, where their dynamics parameters  $\mu$  and  $\nu$  are sampled from the outside of the training range. We conduct the experiments for the CartPole, Hopper, Walker2D and quadrupedal locomotion tasks and compare the *normalized* average returns (the average return divided by the maximum return). The results are plotted in Figure Figure 5.5, which indicate that EAP outperforms DR by 60% to 116% and UP by 12% to 77%. Note that UP may perform well for the real-world transfer due to the lack of the unobservable parameters. We also observe that UP is consistently better than DR by being aware of the dynamics parameters,  $\mu$  and  $\nu$ , which meets our expectation.



Figure 5.5: Comparison of EAP and baselines DR and UP. The error bars represent the variation in the average return of the policy in the target environment when trained with 4 different seeds.



Figure 5.6: Average stability region in five test subjects. The results indicate the better zero-shot transfer of EAP over DR and UP.

For evaluating the zero-shot transfer for the assistive walking task, we define an additional metric "stability region", which depicts the ranges of maximum perturbations in all directions that can be handled by the human with the EAP-controlled exoskeleton. We train policies for 10 training human subjects and test the learned policies for 5 new human subjects. Figure Figure 5.6 compares the average performance of EAP with DR and UP. The larger area of stability region indicates that EAP significantly outperforms two baselines. The ranges of variation for observable and unobservable parameters during training and testing phases are included in tables Table 5.2 and Table 5.3.

#### 5.3.4 Ablation study

We further analyze the performance of EAPs by conducting a set of ablation studies. We studied four categories of parameters: choices of observable parameters, reference dynamics, error prediction horizons, and error representations.



Figure 5.7: Ablation study with choosing different observable parameters as  $\mu$ . The result indicates that our approach (EAP) shows more reliable zero-shot transfers for all different scenarios.



Figure 5.8: Ablation study with different reference dynamics. The results indicate that our algorithm is robust against the choice of different references.

Choice of Observable and Unobservable Parameters. We check the robustness of EAPs by testing with different choices of observable and unobservable parameters. We randomly split the parameters into  $\mu$  and  $\nu$  and test three different splits. Figure Figure 5.7 shows the stability regions for all three algorithms for three different scenarios. In all cases, EAPs are more robust than the baseline algorithms.

Choice of Reference Dynamics. In this study, we analyze the effect of choosing three different reference dynamics  $P(\mathbf{s}'|\mathbf{a}, \mathbf{s}, \boldsymbol{\mu}_0, \boldsymbol{\nu}_0)$  on the performance of EAP. We randomly choose three different human agents as the reference dynamics and follow the learning procedure of EAPs to train three different policies. These policies are then deployed on the same test subjects along with UP and DR policies. Figure Figure 5.8 shows that all the EAPs outperforming the baselines by having larger stability regions, although EAPs have



Figure 5.9: Ablation study with different parameter setting for EAP training.

slightly larger variances.

Horizon of Error Prediction. As we motivated in Section subsection 5.2.3, one step error might be too subtle to inform the learning of EAPs and we may need T step expansion to enlarge them. We studied the effect of the error prediction horizon T in Algorithm algorithm 7 by varying its value from T = 1 to T = 8 for the assistive walking task. Figure Figure 5.9 shows the normalized average return over T gradually changes over the different values of T and peaks at T = 5. Therefore, we set T = 5 for all the experiments.

The error representation. We also compare the effect of the error representation. Figure Figure 5.9 also plots the normalized average returns of the unprojected errors (blue) and projected errors (orange), where projected errors show slightly better performance for all the different T values.

# 5.4 Hardware experiments

Although the primary inspiration for developing this algorithm was to enable transfer of control policies for assistive devices, the proposed zero-shot transfer approach is applicable to any robot whose dynamics are partial observable due to challenges in measuring certain parameters. To validate the effectiveness of EAP algorithm on a real robotic system, we use the popular quadrupedal robot A1 from Unitree [145] robotics, shown in Figure 5.10, as our testbed. A common, well established, approach while working with real robotic systems



Figure 5.10: Unitree's A1 quardupedal robot.

is to first perform a thorough system identification of the robot's kinematic and dynamic properties by collecting real world data. However, standard system identification methods fail when a robot is tested in an environment from which no prior data was available. For example, when data is collected on regular ground but the robot is tested on a ground with different contact parameters such as a soft foam mat, the resulting variation in dynamics is not captured by the identified parameters. As a result, control policies trained in simulation fail when tested on a novel environment in the real world. Using EAP, we demonstrate that a policy trained using our approach is able to overcome differences in the environments it is being tested in. We take the following approach to illustrate this :

- 1. We first perform a system identification of A1 robot parameters using data collected over carpeted ground.
- Then we define two tasks for the quadrupedal robot In-place walking and walking forward. For each of these tasks, we define a set of observable and unobservable parameters and train EAP along with policies using baseline algorithms - UP and DR.
- 3. We compare the performance of each of these policies on three novel environments from which no data for system identification was collected. First environment uses a soft foam mat as ground. Second, we add an additional unknown mass to the robot's

torso. Third, we test the robot with additional mass using soft foam mat as ground.

#### 5.4.1 Data collection, simulation environment and system identification

#### Data collection

The A1 robot is controlled using a position command which is converted into torques using a simple PD-controller. It is equipped with joint angle sensors, Inertial-Measurement Unit (IMU) and a foot contact sensor. Using the joint angle sensors we compute both joint angle and joint velocities. We can also compute the robots torso's orientation and rate of change of orientation using the IMU sensor.

The data collection process on the real A1 robot involves generating two kinds of data:

- We first suspend the robot in air, thus removing ground contact forces being applied to the robot, and apply a variety of sinusoidal and step commands to the robot. During this process, we make sure that the robots torso remains fixed as much as possible. We collect joint position and velocity responses from the robot.
- 2. Next, we script simple sinusoidal commands with different frequencies and amplitudes that the robot legs can follow while being in contact with the ground (Note, we stick to simple commands because it is challenging to script commands for the robot in contact with the ground while ensuring safety). We collect joint position, velocities , the torso's orientation , angular velocity (computed using simple finite-difference method).

We collect 60 trajectories on the robot each of length 10 seconds. The sensors are sampled at 30 Hz, which gives us a total of 60,000 data points. These form dataset D, each trajectory  $\tau$  is consists of the command sent to the robot a and the state of the robot s (which includes the joint positions, velocities, torso position and orientation) for length of time T.

$$D = \{( au_1, ), ( au_2)...( au_n)\}$$
 $au_i = \{(s_1, a_1), (s_2, a_2)...(s_T, a_T)\}$ 

#### Simulation environment and system identification

We use DART physics engines [134] to simulate the robot. In total, the robot has 18 degrees of freedom, 12 controllable joints of the legs and 6 free degrees of freedom of the torso. The simulation runs at 500 Hz, however we use a control frequency of 20 Hz. In simulation, we enforce a torque limit of 30 N/m based on the torque capabilities of the real robot. In our experiments with the real robot, we measured a sensor delay of 20-30 ms, we incorporate this sensor delay into our simulation framework as well. The feet of A1 robot are made from deformable rubber material which acts like a spring-damper when there is contact, to model this interaction we use hunt-crossley contact model [146, 147]. Hunt-crossley model comprises of parameters such as stiffness and friction which allows us to control the softness of the contact. To identify a simulation model from real world data, we choose 8 dynamical parameters : The proportional and derivative gains of the low-level controller  $K_p, K_d$  on the robot. DART also allows us to write custom lowlevel controllers (such as PD controllers), which enables users to design these controllers to compute the torques. Since, the PD controller gains affect the dynamic response of the robot joint legs, we include these parameters as optimization variables. We also include joint damping  $\beta$ , contact stiffness  $\kappa$ , ground friction  $\gamma$ , mass of torso  $m_{torso}$ , mass of thigh link  $m_{thigh}$  and mass of the calf link  $m_{calf}$ . This makes the vector containing these variables  $\eta \in \mathbb{R}^8$ .

$$\eta = \{K_p, K_d, \beta, \kappa, \gamma, m_{torso}, m_{thigh}, m_{calf}\}$$
(5.5)

For optimization, we use Covariance Matrix Adaptation (CMA-ES) [148] method to

Task	Observable Params. $\mu$	Unobservable Params. $\nu$	Net. Arch.	Err. Dim. $ \mathbf{e}_p $
In-Place Walking	PD-gains, link masses	joint damping, contact params , delay	(64,64,64)	6
Walking forward	PD gains, link masses	sensor delay, joint damping, contact params	(64,64,64)	6
Walking forward w unknown mass	PD-gains	joint damping, delay, torso mass	(64,64,64)	6

Table 5.4:	Tasks and Net	work Architectu	res on the real robot

find the best set of parameter values in simulation that matches real world trajectories by minimizing the cost function described in Equation 5.6.

$$J = \frac{1}{N} \min_{\eta} \sum_{n=1}^{N} \sum_{t=0}^{T} ||s_t' - s_t||^2$$
(5.6)

Where  $s_t$  is the states collected in real world and  $s'_t$  are the states observed in simulation. The cost function minimizes the average error generated by N state trajectories each of T time steps in simulation when compared to real world data when the same actions are applied to both. Once the optimized parameters are obtained, we use this as the base environment in simulation.

## 5.5 In-place walking

We define an in-place walking task for a quadrupedal robot in which the center of mass of the robot remains in the same place while the legs follow a walking motion. Inspired by the lower level control framework presented in [149], we first define a end-effector reference trajectory for each leg and our action space is defined as delta end-effector cartesian positions from the pre-defined reference trajectory. We then use analytical inverse kinematics to compute the corresponding joint angles for each leg. For in-place walking, since there is no movement in the x-y plane of the robot, the reference trajectory only consists of motion in the z-axis which is perpendicular to the ground. The reference trajectory is defined in the equation below, where A is the amplitude, f is the frequency and t is the time. When the cartesian z is below zero or the ground plane the value just remains zero.

$$z = \begin{cases} Asin(2\pi ft), & \text{if } z \ge 0\\ 0, & \text{otherwise} \end{cases}$$

The action space is 12-dimensional  $(R^{12})$  which comprises of residual cartesian position for each leg end-effector  $(R^3)$ . So, the output of the neural network policy is added to the reference trajectory to compute the target leg positions. The state space consists of the angular velocity of torso  $R^3$  computed using the gyroscope. The angular position  $R^3$  of torso, joint angles  $R^{12}$  and binary foot contact information for each leg which indicates if the foot is in contact or not  $R^4$ . An additional phase variable  $\phi(t)$  is included as part of the state to indicate the phase of the gait cycle the robot is in. In total, the state vector is 23 dimensional  $R^{23}$ . The reward function , described in Equation 5.7 encourages the policy to generate actions that penalizes center of mass velocity of the robot while following the reference trajectory  $\bar{q}$  of the end-effector of each leg.

$$r(\mathbf{s}, \mathbf{a}) = w_1 e^{-k_1 * (\mathbf{q} - \bar{\mathbf{q}})} + w_2 e^{-k_2 * (\dot{\mathbf{c}})} - w_3 ||a||^2$$
(5.7)

Where  $w_1 = 0.5$ ,  $w_2 = 0.5$ ,  $k_1 = 20$ ,  $k_2 = 10$  and  $w_3 = 1e^{-3}$  During training, we consider the PD-gains, the mass of the robot links as observable parameters  $\mu$ . In our preliminary experiments, we found that for simple tasks, these quantities can be identified by standard system identification procedures. However, quantities like joint damping, sensor delays and ground contact parameters are considered unobservable parameters  $\nu$ . We train EAP, UP and DR with these parameters as the set of observable and unobservable quantities.

The neural network policy is represented as a multi-layered perceptron (MLP) with three hidden layers of 64 neurons each. The error function is trained with a horizon length of five T = 5 and the projected error dimension is six  $R^6$ .

#### 5.6 Walking forward

In this task, we train a policy to walk forward at a certain velocity while maintaining the yaw orientation. The framework for training the policy is similar to in-place walking task. A pre-defined reference trajectory is defined for each leg end-effector and a neural network policy outputs delta cartesian positions for each leg. The reference trajectories are defined as follows :

$$z = \begin{cases} Asin(2\pi ft - \frac{\pi}{2}), & \text{if } z \ge 0\\ 0, & \text{otherwise} \end{cases}$$

$$x = Asin(2\pi ft), y = 0$$

The reward function encourage forward walking at a velocity of v = 0.6m/s while ensuring that the end-effectors follow the pre-defined reference trajectory. The reward function also encourages the policy to walk in a straight path by penalizing the yaw orientation  $\psi$ .

$$r(\mathbf{s}, \mathbf{a}) = w_1 e^{-k_1 * (\mathbf{q} - \bar{\mathbf{q}})} + w_2 e^{-k_2 * \psi} + \min(\dot{x}, 0.6) - w_3 ||a||^2$$
(5.8)

Where  $w_1 = 0.5, k_1 = 20, w_2 = 0.6, k_2 = 10$  and  $w_3 = 1e^{-3}$ . The state space, action space, observed, unobserved parameters and test environment remains the same as in-place walking task. The parameters of the policy and the error function remains the same as the previous task as well.

#### 5.7 Evaluation on real robot

We test policies on five environments : A training environment and four test environments.

1. Base training environment : The same environment from which data for system iden-

tification was collected (robot moving on carpeted ground).

- 2. Foam mat environment (T1): The robot is deployed on a soft foam mat with different contact parameters than over ground.
- 3. Additional mass environment (T2): Additional unknown mass is added to the robot's torso.
- 4. Foam mat with additional mass environment (T3): Robot is tested with additional mass and deployed over a foam mat.

Owing to the difficulty in computing the rewards in the real world. We use a surrogate metric to evaluate the policies in the real world. For in-place walking task, we collect the joint angles when the policy is being deployed and compare how well the joint angles track the reference trajectory. This quantity is computed using the term  $w_1 e^{-k_1*(\mathbf{q}-\bar{\mathbf{q}})}$  in equation Equation 5.7. For walking forward task, we just use the straight line distance travelled by the robot as the surrogate metric, this is measured using a simple measuring tape.

## Performance in base training environment

In the base training environment we notice that for both in-place walking as well as walking forward task all methods (EAP,DR and UP) show similar performance. The results are illustrated in Figure 5.12a and Figure 5.12b where the task performance for the base environment is similar. This corroborates with findings of prior work [150] which suggested that for simple tasks, carefully randomizing parameters is sufficient for sim-2-real transfer for quadrupedal robots. Since, the policies are being tested in the environment from which data for system identification was collected, the policies are able to transfer well.

# Performance in test environments

In our first test environment (T1), we change the surface on which the robot performs the task, with the aim of changing the contact parameter values. The differences in the contact



Figure 5.11: Comparison of contact profile generated by ground and soft foam mat

profile generated by the two surfaces is illustrated in figure Figure 5.11. We notice that for in-place walking task, all methods perform similarly Figure 5.12a. However, in the walking forward task Figure 5.12b, EAP beats both the baselines DR and UP by walking forward for a longer straight-line distance on the mat with an average of 2.3m whereas both DR and UP move less than 2m. We notice that both DR and UP tend to change the yaw-orientation of the robot, which is not a desired behaviour.

In test environment 2 (T2), we attach a mass of 7.5 *lbs* (unknown to the policy parameters) to the robot's torso. In this environment, the ability for EAP to adapt is best illustrated, beating both baselines DR and UP in straight line distance travelled on carpeted ground. With EAP the robot moves an average distance of 1.7m, UP and DR manage close 1m. Results illustrated in Figure 5.12b.

The performance of the policies in test environment 3 (T3) are similar to that of T2 with EAP managing to travel the furthest distance compared to baselines DR and UP.

# 5.8 Conclusion

We presented a novel approach to train an error-aware policy (EAP) that transfers effectively to unseen target environments in a zero-shot manner. Our method learns an EAP for an assistive wearable device to help a human recover balance after an external push is applied. We show that a single trained EAP is able to assist different human agents with



Figure 5.12: Evaluation of real world results.

unseen biomechanical characteristics. We also validate our approach by comparing EAP to common baselines like Universal Policy and Domain randomization to show our hypothesis that a policy which explicitly takes future state error as input can enable better decision making. Our approach outperforms the baselines in all the tasks. We also evaluated the performance of our algorithm through a series of ablation studies that sheds some light on the importance of parameters such as error horizon length, error representation, choice of observable parameters and choice of reference dynamics. We find that EAP is not sensitive to either the choice of observable parameters or the reference dynamics, and outperforms the baselines with variations in these quantities as well. Further, we validated the ability of EAP to transfer to real world scenario using A1 quadrupedal robot. The preliminary results presented for simple tasks such as in-place walking and walking forward shows promise of applying the proposed approach to real world scenarios.

Our work has a few limitations. At the core, our algorithm relies on the error function to make predictions of the expected state errors. The accuracy of this prediction can be improved by better function approximators such as recurrent neural networks (RNN) that takes a history of states as input, we leave this for future work.

# CHAPTER 6 CONCLUSION AND FUTURE WORK

In conclusion, we presented of a set of learning based algorithms that address the important challenge of safety in bipedal locomotion. Our contributions are along two closely related directions 1) Fall prevention and safe falling for bipedal robots and 2) Fall prevention for humans during walking using assistive devices. With bipedal robots, we showed the effectiveness of reinforcement learning (RL) based approaches to learn control policies that demonstrate a wide range of falling and balancing strategies. Our algorithmic contributions also included improving efficiency of learning by incorporating abstract dynamical models as priors, curriculum learning, imitation learning and a novel method of building a graph of policies into the standard RL framework. In the chapters 4 and 5 of the thesis, we shifted focus towards an important healthcare related topic of fall prevention for humans using assistive devices. To this end, our contributions include using imitation learning approach to create virtual human walking agents whose bio-mechanical gait characteristics are similar to real world humans. Using these simulated agents we showed that, assistive devices could indeed help prevent falls when pushed whereas human agents without assistants fail. Finally, since humans exhibit a wide range of variations in gait characteristics, we developed a novel algorithm to enable zero-shot generalization of fall prevention policies when tested on new human subjects. Although inspired by assistive devices, the algorithm is applicable for any robotic system and we validated the proposed approach on a real world quadrupedal A1 robot.

# 6.1 Safe locomotion for bipedal robots

In chapter 3, we developed an algorithm to generate safe falling motion for a bipedal robot. Unlike steady state walking, which is periodic in nature, optimizing for a falling motion is challenging because it involves reasoning about discrete contact sequence planning (which body part should hit the ground and in what sequence) as well a continuous joint motion to minimize impact upon ground contact. We showed that we could solve this problem using an offline policy learning algorithm. Our approach significantly improved on the computation time required by prior dynamic programming methods while also ensuring that a wide variety of falling strategies naturally emerge from the algorithm. Towards the same goal of safety for bipedal locomotion, inspired by strategies humans use to recover balance as a response to external pushes, in our second work we presented a training methodology for a control policy which learns residual control signals to those generated by traditional model-based controllers to maintain balance. By learning residual control signals, we not just simplify learning the task because the model-based controllers act as strong priors, but also improve upon methods that use either just model-based or naive RL approaches. Our algorithm also included a novel sampling method that adaptively adjusts the difficulty in training samples to encourage efficient learning. By doing so, we achieved better performance at a lower sample cost compared to baseline naive RL approach.

We also addressed the common issue of sample inefficiency while solving complex task such as locomotion with model free reinforcement learning algorithms. RL has shown promising results in learning complex motor skills in high dimensional systems, however the control policies typically perform well only from a small set of initial states and take millions of samples to complete the task. In the final section of chapter 3, we take a divide and conquer approach to solve this by breaking down a complex task into multiple sub-tasks and learn a directed-graph of control policies to achieve the same task but by consuming less samples. The edges of our directed graph contains policies and the nodes contain state distributions, so a policy can take the robotic system from the one set of states to another set of goal states via rolling out policy/policies. Starting from the first policy that attempts to achieve the task from a small set of initial states, the algorithm automatically discovers the next subtask with increasingly more difficult initial states until the last subtask matches

the initial state distribution of the original task. We showed that our approach takes lesser samples than common baselines such as naive RL, curriculum learning, manually generated sub-tasks while ensuring that the task can be completed from a wider range of initial states.

#### Future work

The limitation with our work on push recovery, be it either fall prevention or safe falling, is that the control policies work best when the robot's initial center of mass velocity is close to zero when pushed. However, in real world scenarios the ability to recover when the robot's motion is more dynamic, such as walking or running, would be crucial to ensure safety under all circumstances. Developing algorithms that can ensure safety for a wider range of scenarios would be an interesting extension to our work.

## 6.2 Fall prevention using assistive devices

In chapter 4, we shifted our focus towards learning fall prevention control policies for assistive devices. We developed an approach to automate the process of augmenting an assistive walking device with the ability to prevent falls. Our method has three key components : A human walking policy, fall predictor and a recovery policy. In a simulated environment we showed that an assistive device can indeed help recover balance from a wider range of external perturbations. We introduced stability region as a quantitative metric to show the benefit of using a recovery policy, larger area of stability region would imply better recovery motion. In addition to this, stability region can also be used to analyze different sensor and actuator design choices for an assistive device. We evaluated the stability region of six different policies using various sensor and actuator configurations to shed some light on appropriate design choices to help develop effective assistive devices.

A significant bottleneck in making assistive devices more ubiquitous in society is due to the fact that humans exhibit remarkable variations in gait characteristics, making design of control policies that generalize for a large population extremely challenging. A common
approach has been to tune controllers for each individual user, however this can be tedious and a time consuming process. We addressed this issue by viewing it through the lens of transfer learning problem. We developed a novel algorithm that enables a policy to transfer in a zero-shot manner for partially observed dynamical systems like humans. Our work introduces an Error-aware policy (EAP) which explicitly takes as input a predicted state error generated in the target environment to produce a corrected action that successfully completes the task. We show that a single trained EAP is able to assist different human agents with unseen biomechanical characteristics. We also validate our approach by comparing EAP to common baselines like Universal Policy (UP) and Domain randomization (DR) to show our hypothesis that a policy which explicitly takes future state error as input can enable better decision making. Our approach outperforms the baselines in all the tasks. We also evaluated the performance of our algorithm on a real world quadrupedal robot, which strengthens the case for applicability of EAP on a real system. We show that for tasks such as walking forward, EAP is able to outperform baselines such as DR and UP when tested on novel environments with unique ground contact characteristics.

## Future work

A key component of our work with assistive devices is virtual human walking agent. Simulated agents that generate bio-mechanically accurate motion [151, 152, 153] can play a crucial role in not just improving our understanding of human motion but also help in developing better assistive and rehabilitation strategies for people with disability. Algorithms that can generate such human-like motion for not just periodic steady state walking gaits, but for more dynamic motion such as running, push-recovery, slip/ trip recovery can be a very exciting research direction. Incorporating musculo-tendon human models could also be a contributing factor in closing the gap between real world and simulated human agents. The major bottleneck for such applications is the shortage of real-world data for validation purposes, recent advancements in machine learning methods to learn from limited data as well as a collaborative efforts between researchers working in bio-mechanics and machine learning can help overcome some of these challenges.

In this work, our primary contributions have been on developing control policies for robots, however in the spirit of the famous phrase "form follows function", co-optimization of robot design along with policy learning can be a powerful approach to create functional robots in the future. In chapter 4, we provided some preliminary results by shedding light on the best sensor and actuator design choices for assistive devices based on the policy's performance. Extending this to optimize for robot structural design can be an interesting direction.

Another closely related topic to our work with assistive devices is modelling humanrobot interaction. In our current, work we assume that the human perfectly adapts to the forces applied by the exoskeleton. However, in reality this does not necessarily hold true. Better predictive modelling of the interaction between the robot and the human is a promising research direction. Modelling human intention while collaborating with a robot could also be key towards wide-spread use of these robots. Contributions in this area can have a lasting impact not just with assistive devices but also numerous other applications such as human-robot collaboration in a factory setting, self-driving cars and robots that are built for home environments.

As highlighted in Chapter 5, sim-2-real transfer of policies has gained a lot of interest in the research community. However, the primary focus from researchers has been through the lens of algorithmic contributions, an exciting future direction will be to develop methods that improves simulation environments by bringing them closer to reality. This could be done through leveraging machine learning methods to improve contact models, focusing efforts on differentiable simulators to enable end-to-end learning or by improving computational efficiency of more traditional methods such as finite-element methods (FEM) to improve accuracy of simulation, especially for deformable bodies.

## REFERENCES

- [1] K. E. Adolph, W. G. Cole, M. Komati, J. S. Garciaguirre, D. Badaly, J. M. Lingeman, G. L. Chan, and R. B. Sotsky, "How do you learn to walk? thousands of steps and dozens of falls per day," *Psychological science*, vol. 23, no. 11, pp. 1387–1394, 2012.
- [2] Z. Wang and H. Gu, "A review of locomotion mechanisms of urban search and rescue robot," *Industrial Robot: An International Journal*, 2007.
- [3] H. L. Lee, Y. Chen, B. Gillai, and S. Rammohan, "Technological disruption and innovation in last-mile delivery," *Value Chain Innovation Initiative*, 2016.
- [4] B. Stephens, "Integral control of humanoid balance," *IEEE International Conference on Intelligent Robots and Systems*, 2007.
- [5] S. Ha and C. K. Liu, "Multiple Contact Planning for Minimizing Damage of Humanoid Falls," *IEEE IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015.
- [6] A. Goswami, S.-k. Yun, U. Nagarajan, S.-H. Lee, K. Yin, and S. Kalyanakrishnan, "Direction-changing fall control of humanoid robots: theory and experiments," *Autonomous Robots*, vol. 36, no. 3, pp. 199–223, Jul. 2014.
- [7] C. McGreavy, K. Yuan, D. Gordon, K. Tan, W. J. Wolfslag, S. Vijayakumar, and Z. Li, "Unified push recovery fundamentals: Inspiration from human study," in 2020 IEEE International Conference on Robotics and Automation (ICRA), 2020, pp. 10876–10882.
- [8] W.-b. I. S. Query and R. S. (WISQARS), "Centers for disease control and prevention, national center for injury prevention and control," 2013.
- [9] C. S. Florence, G. Bergen, A. Atherly, E. Burns, J. Stevens, and C. Drake, "Medical costs of fatal and nonfatal falls in older adults," *Journal of the American Geriatrics Society*, vol. 66, no. 4, pp. 693–698, 2018.
- [10] D. A. Sterling, J. A. O'connor, and J. Bonadies, "Geriatric falls: Injury severity is high and disproportionate to mechanism," *Journal of Trauma and Acute Care Surgery*, vol. 50, no. 1, pp. 116–119, 2001.
- [11] B. H. Alexander, F. P. Rivara, and M. E. Wolf, "The cost and frequency of hospitalization for fall-related injuries in older adults.," *American journal of public health*, vol. 82, no. 7, pp. 1020–1023, 1992.

- [12] T. E. Jager, H. B. Weiss, J. H. Coben, and P. E. Pepe, "Traumatic brain injuries evaluated in us emergency departments, 1992-1994," *Academic Emergency Medicine*, vol. 7, no. 2, pp. 134–140, 2000.
- [13] V. C. Kumar, S. Ha, and C. K. Liu, "Learning a unified control policy for safe falling," *IEEE/RSJ International Conference on Intelligent Robots and Systems* (*IROS*), 2017.
- [14] V. C. V. Kumar, S. Ha, and K. Yamane, "Improving model-based balance controllers using reinforcement learning and adaptive sampling," in 2018 IEEE International Conference on Robotics and Automation (ICRA), 2018, pp. 7541–7547.
- [15] V. C. Kumar, S. Ha, and C. Liu, "Expanding motor skills using relay networks," in *Proceedings of The 2nd Conference on Robot Learning*, A. Billard, A. Dragan, J. Peters, and J. Morimoto, Eds., ser. Proceedings of Machine Learning Research, vol. 87, PMLR, 29–31 Oct 2018, pp. 744–756.
- [16] V. C. V. Kumar, S. Ha, G. Sawicki, and C. K. Liu, *Learning a control policy for fall prevention on an assistive walking device*, 2019. arXiv: 1909.10488 [cs.RO].
- [17] D. A. Winter, *Biomechanics and motor control of human gait: normal, elderly and pathological.* Waterloo Biomechanics, 1991.
- [18] L. M. Nashner and G. McCollum, "The organization of human postural movements: A formal basis and experimental synthesis," *Behavioral and brain sciences*, vol. 8, no. 1, pp. 135–150, 1985.
- [19] M. E. Gordon, "An analysis of the biomechanics and muscular synergies of human standing," Ph.D. dissertation, Stanford University, 1991.
- [20] A. D. Kuo and F. E. Zajac, "Human standing posture: Multi-joint movement strategies based on biomechanical constraints," *Progress in brain research*, vol. 97, pp. 349– 358, 1993.
- [21] B. Stephens, "Push recovery control for force-controlled humanoid robots," Ph.D. dissertation, 2011.
- [22] A. Macchietto, V. Zordan, and C. R. Shelton, "Momentum control for balance," *ACM Transactions on graphics (TOG)*, vol. 28, 2009.
- [23] A. D. Kuo, "An optimal control model for analyzing human postural balance," *IEEE transactions on biomedical engineering*, vol. 42, no. 1, pp. 87–101, 1995.
- [24] Z. Aftab, T. Robert, and P.-B. Wieber, "Ankle, hip and stepping strategies for humanoid balance recovery with a single model predictive control scheme," in *Hu*-

manoid Robots (Humanoids), 2012 12th IEEE-RAS International Conference on, IEEE, 2012, pp. 159–164.

- [25] N. Perrin, N. Tsagarakis, and D. G. Caldwell, "Compliant attitude control and stepping strategy for balance recovery with the humanoid COMAN," *IEEE International Conference on Intelligent Robots and Systems*, pp. 4145–4151, 2013.
- [26] T. Komura, H. Leung, S. Kudoh, and J. Kuffner, "A feedback controller for biped humanoids that can counteract large perturbations during gait," *Proceedings - IEEE International Conference on Robotics and Automation*, vol. 2005, pp. 1989–1995, 2005.
- [27] J. Pratt, J. Carff, S. Drakunov, and A. Goswami, "Capture point: A step toward humanoid push recovery," in *Humanoid Robots, 2006 6th IEEE-RAS International Conference on*, IEEE, 2006, pp. 200–207.
- [28] J. Pratt, T. Koolen, T. de Boer, J. Rebula, S. Cotton, J. Carff, M. Johnson, and P. Neuhaus, "Capturability-based analysis and control of legged locomotion, Part 2: Application to M2V2, a lower-body humanoid," *The International Journal of Robotics Research*, vol. 31, no. 10, pp. 1117–1133, 2012.
- [29] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [30] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [31] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust region policy optimization," *CoRR*, *abs/1502.05477*, 2015.
- [32] G. Ma, Q. Huang, Z. Yu, X. Chen, K. Hashimoto, A. Takanishi, and Y.-H. Liu, "Bio-inspired falling motion control for a biped humanoid robot," in 2014 IEEE-RAS International Conference on Humanoid Robots, IEEE, 2014, pp. 850–855.
- [33] E. T. Hsiao and S. N. Robinovitch, "Common protective movements govern unexpected falls from standing height," *Journal of biomechanics*, vol. 31, no. 1, pp. 1–9, 1997.
- [34] S. N. Robinovitch, R. Brumer, and J. Maurer, "Effect of the "squat protective response" on impact velocity during backward falls," *Journal of biomechanics*, vol. 37, no. 9, pp. 1329–1337, 2004.

- [35] J.-S. Tan, J. J. Eng, S. N. Robinovitch, and B. Warnick, "Wrist impact velocities are smaller in forward falls than backward falls from standing," *Journal of biomechanics*, vol. 39, no. 10, pp. 1804–1811, 2006.
- [36] S. N. Robinovitch, J. Chiu, R. Sandler, and Q. Liu, "Impact severity in self-initiated sits and falls associates with center-of-gravity excursion during descent," *Journal* of biomechanics, vol. 33, no. 7, pp. 863–870, 2000.
- [37] K. Fujiwara, F. Kanehiro, S. Kajita, K. Kaneko, K. Yokoi, and H. Hirukawa, "UKEMI: falling motion control to minimize damage to biped humanoid robot," *IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 3, no. October, 2002.
- [38] J. Ruiz-del-solar, S. Member, J. Moya, and I. Parra-tsunekawa, "Fall Detection and Management in Biped Humanoid Robots," *Management*, vol. 12, no. April, pp. 3323–3328, 2010.
- [39] K. Fujiwara, S. Kajita, K. Harada, K. Kaneko, M. Morisawa, F. Kanehiro, S. Nakaoka, and H. Hirukawa, "Towards an optimal falling motion for a humanoid robot," *Proceedings of the 2006 6th IEEE-RAS International Conference on Humanoid Robots, HUMANOIDS*, pp. 524–529, 2006.
- [40] —, "An Optimal planning of falling motions of a humanoid robot," *IEEE International Conference on Intelligent Robots and Systems*, no. Table I, pp. 456–462, 2007.
- [41] J. Wang, E. Whitman, and M. Stilman, "Whole-body trajectory optimization for humanoid falling," ... *Control Conference (ACC)*, ..., pp. 4837–4842, 2012.
- [42] S. K. Yun and A. Goswami, "Tripod fall: Concept and experiments of a novel approach to humanoid robot fall damage reduction," *Proceedings IEEE International Conference on Robotics and Automation*, pp. 2799–2805, 2014.
- [43] V. Mnih, K. Kavukcuoglu, D. Silver, A. a. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015. arXiv: 1312.5602.
- [44] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, pp. 1–14, 2015. arXiv: arXiv:1509.02971v1.

- [45] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic Policy Gradient Algorithms," *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- [46] R. S. Sutton and A. G. Barto, *Reinforcement learning : an introduction*. 1988, ISBN: 9780262193986. arXiv: 1603.02199.
- [47] H. Van Hasselt and M. A. Wiering, "Reinforcement Learning in Continuous Action Spaces," no. Adprl, pp. 272–279, 2007.
- [48] H. Van Hasselt, "Reinforcement Learning in Continuous State and Action Spaces," *Reinforcement Learning*, pp. 207–251, 2012.
- [49] M. Hausknecht and P. Stone, "Deep Reinforcement Learning in Parameterized Action Space," *arXiv*, pp. 1–12, 2016. arXiv: 1511.04143.
- [50] X. B. Peng, G. Berseth, and M. van de Panne, "Terrain-Adaptive Locomotion Skills using Deep Reinforcement Learning," ACM Transactions on Graphics, vol. 35, no. 4, pp. 1–10, 2016.
- [51] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International Conference on Machine Learning*, 2016, pp. 1928–1937.
- [52] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [53] R. S. Sutton, D. Precup, and S. Singh, "Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning," *Artificial intelligence*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [54] C. Daniel, G. Neumann, and J. Peters, "Hierarchical relative entropy policy search," in *Artificial Intelligence and Statistics*, 2012, pp. 273–281.
- [55] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum, "Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation," in Advances in neural information processing systems, 2016, pp. 3675–3683.
- [56] N. Heess, G. Wayne, Y. Tassa, T. Lillicrap, M. Riedmiller, and D. Silver, "Learning and transfer of modulated locomotor controllers," *arXiv preprint arXiv:1610.05182*, 2016.
- [57] X. B. Peng, G. Berseth, K. Yin, and M. Van De Panne, "Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning," *ACM Transactions on Graphics (TOG)*, vol. 36, 2017.

- [58] T. G. Dietterich, "Hierarchical reinforcement learning with the maxq value function decomposition," *J. Artif. Int. Res.*, vol. 13, no. 1, pp. 227–303, Nov. 2000.
- [59] A. Bai, F. Wu, and X. Chen, "Online planning for large mdps with maxq decomposition," in *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems Volume 3*, ser. AAMAS '12, Valencia, Spain: International Foundation for Autonomous Agents and Multiagent Systems, 2012, pp. 1215–1216, ISBN: 0-9817381-3-3, 978-0-9817381-3-0.
- [60] K. Gräve and S. Behnke, "Bayesian exploration and interactive demonstration in continuous state maxq-learning," in 2014 IEEE International Conference on Robotics and Automation, ICRA 2014, Hong Kong, China, May 31 - June 7, 2014, 2014, pp. 3323–3330.
- [61] R. Tedrake, "Lqr-trees: Feedback motion planning on sparse randomized trees," 2009.
- [62] M. A. Borno, M. V. D. Panne, and E. Fiume, "Domain of attraction expansion for physics-based character control," ACM Transactions on Graphics (TOG), vol. 36, no. 2, p. 17, 2017.
- [63] G. Konidaris and A. G. Barto, "Skill discovery in continuous reinforcement learning domains using skill chaining," in Advances in neural information processing systems, 2009, pp. 1015–1023.
- [64] S. Coros, P. Beaudoin, and M. Van de Panne, "Robust task-based control policies for physics-based characters," in ACM Transactions on Graphics (TOG), ACM, vol. 28, 2009, p. 170.
- [65] L. Liu and J. Hodgins, "Learning to schedule control fragments for physics-based characters using deep q-learning," ACM Transactions on Graphics (TOG), vol. 36, no. 3, p. 29, 2017.
- [66] S. Kakade and J. Langford, "Approximately optimal approximate reinforcement learning," in *ICML*, vol. 2, 2002, pp. 267–274.
- [67] I. Popov, N. Heess, T. Lillicrap, R. Hafner, G. Barth-Maron, M. Vecerik, T. Lampe, Y. Tassa, T. Erez, and M. Riedmiller, "Data-efficient deep reinforcement learning for dexterous manipulation," *arXiv preprint arXiv:1704.03073*, 2017.
- [68] C. Florensa, D. Held, M. Wulfmeier, M. Zhang, and P. Abbeel, "Reverse curriculum generation for reinforcement learning," in *Proceedings of the 1st Annual Conference on Robot Learning*, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., ser. Proceedings of Machine Learning Research, vol. 78, PMLR, 13–15 Nov 2017, pp. 482–495.

- [69] Y. Wang and M. Srinivasan, "Stepping in the direction of the fall: The next foot placement can be predicted from current upper body state in steady-state walking," *Biology Letters*, vol. 10, no. 9, 2014.
- [70] A. Hof, S. Vermerris, and W. Gjaltema, "Balance responses to lateral perturbations in human treadmill walking," *Journal of Experimental Biology*, vol. 213, no. 15, pp. 2655–2664, 2010.
- [71] V. Joshi and M. Srinivasan, "A controller for walking derived from how humans recover from perturbations," 2019.
- [72] F.Antoine, S. Gil, D. Christopher, G. Joris, J. Ilse, and F. Groote, "Rapid predictive simulations with complex musculoskeletal models suggest that diverse healthy and pathological human gaits can emerge from similar control strategies," 2019.
- [73] S. Wang, L. Wang, C. Meijneke, E. van Asseldonk, T. Hoellinger, G. Cheron, Y. Ivanenko, V. La Scaleia, F. Sylos-Labini, M. Molinari, F. Tamburella, I. Pisotta, F. Thorsteinsson, M. Ilzkovitz, J. Gancet, Y. Nevatia, R. Hauffe, F. Zanow, and H. van der Kooij, "Design and control of the mindwalker exoskeleton," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 23, no. 2, pp. 277–286, Mar. 2015.
- [74] S. Jezernik, G. Colombo, and M. Morari, "Automatic gait-pattern adaptation algorithms for rehabilitation with a 4-dof robotic orthosis," *IEEE Transactions on Robotics and Automation*, vol. 20, no. 3, pp. 574–582, Jun. 2004.
- [75] J. A. Blaya and H. Herr, "Adaptive control of a variable-impedance ankle-foot orthosis to assist drop-foot gait," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 12, no. 1, pp. 24–31, Mar. 2004.
- [76] A. Duschau-Wicke, J. von Zitzewitz, A. Caprez, L. Lunenburger, and R. Riener, "Path control: A method for patient-cooperative robot-aided gait rehabilitation," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 18, no. 1, pp. 38–48, Feb. 2010.
- [77] H. Kazerooni, R. Steger, and L. Huang, "Hybrid control of the berkeley lower extremity exoskeleton (bleex)," *The International Journal of Robotics Research*, vol. 25, no. 5-6, pp. 561–573, 2006. eprint: https://doi.org/10.1177/0278364906065505.
- [78] H. Vallery, J. Veneman, E. van Asseldonk, R. Ekkelenkamp, M. Buss, and H. van Der Kooij, "Compliant actuation of rehabilitation robots," *IEEE Robotics Automation Magazine*, vol. 15, no. 3, pp. 60–69, Sep. 2008.
- [79] M. Hamaya, T. Matsubara, T. Noda, T. Teramae, and J. Morimoto, "Learning assistive strategies for exoskeleton robots from user-robot physical interaction," *Pattern*

*Recognition Letters*, vol. 99, pp. 67–76, 2017, User Profiling and Behavior Adaptation for Human-Robot Interaction.

- [80] G. Bingjing, H. Jianhai, L. Xiangpan, and Y. Lin, "Human–robot interactive control based on reinforcement learning for gait rehabilitation training robot," *International Journal of Advanced Robotic Systems*, vol. 16, no. 2, p. 1 729 881 419 839 584, 2019. eprint: https://doi.org/10.1177/1729881419839584.
- [81] V. C. V. Kumar, S. Ha, G. Sawicki, and C. K. Liu, "Learning a control policy for fall prevention on an assistive walking device," in 2020 IEEE International Conference on Robotics and Automation (ICRA), 2020, pp. 4833–4840.
- [82] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba, *Learning dexterous in-hand manipulation*, 2018. arXiv: 1808.00177 [cs.LG].
- [83] 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," *CoRR*, vol. abs/1703.06907, 2017. arXiv: 1703.06907.
- [84] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," *CoRR*, vol. abs/1703.02702, 2017. arXiv: 1703.02702.
- [85] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," *CoRR*, vol. abs/1710.06537, 2017. arXiv: 1710.06537.
- [86] A. Rajeswaran, S. Ghotra, B. Ravindran, and S. Levine, *Epopt: Learning robust neural network policies using model ensembles*, 2017. arXiv: 1610.01283 [cs.LG].
- [87] F. Muratore, C. Eilers, M. Gienger, and J. Peters, *Data-efficient domain randomization with bayesian optimization*, 2021. arXiv: 2003.02471 [cs.LG].
- [88] B. Mehta, M. Diaz, F. Golemo, C. J. Pal, and L. Paull, Active domain randomization, 2019. arXiv: 1904.04762 [cs.LG].
- [89] F. Ramos, R. C. Possas, and D. Fox, *Bayessim: Adaptive domain randomization via probabilistic inference for robotics simulators*, 2019. arXiv: 1906.01728 [cs.RO].
- [90] 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," *CoRR*, vol. abs/1804.10332, 2018. arXiv: 1804.10332.

- [91] 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, no. 26, 2019. eprint: https://robotics.sciencemag.org/content/4/26/eaau5872. full.pdf.
- [92] Z. Xie, P. Clary, J. Dao, P. Morais, J. Hurst, and M. van de Panne, "Learning locomotion skills for cassie: Iterative design and sim-to-real," in *Proceedings of Machine Learning Research*, L. P. Kaelbling, D. Kragic, and K. Sugiura, Eds., vol. 100, PMLR, 30 Oct–01 Nov 2020, pp. 317–329.
- [93] M. Jegorova, J. Smith, M. Mistry, and T. Hospedales, "Adversarial generation of informative trajectories for dynamics system identification," *arXiv preprint arXiv:2003.01190*, 2020.
- [94] Y. Jiang, T. Zhang, D. Ho, Y. Bai, C. K. Liu, S. Levine, and J. Tan, *Simgan: Hybrid simulator identification for domain adaptation via adversarial reinforcement learning*, 2021. arXiv: 2101.06005 [cs.R0].
- [95] W. Yu, J. Tan, C. K. Liu, and G. Turk, "Preparing for the unknown: Learning a universal policy with online system identification," *arXiv preprint arXiv:1702.02453*, 2017.
- [96] W. Zhou, L. Pinto, and A. Gupta, *Environment probing interaction policies*, 2019. arXiv: 1907.11740 [cs.RO].
- [97] Y. Song, A. Mavalankar, W. Sun, and S. Gao, *Provably efficient model-based policy adaptation*, 2020. arXiv: 2006.08051 [cs.LG].
- [98] S. Desai, H. Karnan, J. P. Hanna, G. Warnell, and P. Stone, *Stochastic grounded* action transformation for robot learning in simulation, 2020. arXiv: 2008.01281 [cs.R0].
- [99] Y. Yang, K. Caluwaerts, A. Iscen, T. Zhang, J. Tan, and V. Sindhwani, "Data efficient reinforcement learning for legged robots," *CoRR*, vol. abs/1907.03613, 2019. arXiv: 1907.03613.
- [100] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox, *Closing the sim-to-real loop: Adapting simulation randomization with real world experience*, 2018. arXiv: 1810.05687 [cs.RO].
- [101] W. Yu, V. C. Kumar, G. Turk, and C. K. Liu, *Sim-to-real transfer for biped loco-motion*, 2019. arXiv: 1903.01390 [cs.R0].

- [102] X. B. Peng, E. Coumans, T. Zhang, T.-W. E. Lee, J. Tan, and S. Levine, "Learning agile robotic locomotion skills by imitating animals," in *Robotics: Science and Systems*, Jul. 2020.
- [103] I. Exarchos, Y. Jiang, W. Yu, and C. K. Liu, *Policy transfer via kinematic domain randomization and adaptation*, 2020. arXiv: 2011.01891 [cs.RO].
- [104] 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, no. 2, pp. 2950–2957, 2020.
- [105] S. Belkhale, R. Li, G. Kahn, R. McAllister, R. Calandra, and S. Levine, Modelbased meta-reinforcement learning for flight with suspended payloads, 2020. arXiv: 2004.11345 [cs.RO].
- [106] I. Clavera, A. Nagabandi, R. S. Fearing, P. Abbeel, S. Levine, and C. Finn, "Learning to adapt: Meta-learning for model-based control," *CoRR*, vol. abs/1803.11347, 2018. arXiv: 1803.11347.
- [107] C. Finn, P. Abbeel, and S. Levine, *Model-agnostic meta-learning for fast adaptation of deep networks*, 2017. arXiv: 1703.03400 [cs.LG].
- [108] Z. Xu, C. Tang, and M. Tomizuka, "Zero-shot deep reinforcement learning driving policy transfer for autonomous vehicles based on robust control," in 2018 21st International Conference on Intelligent Transportation Systems (ITSC), 2018, pp. 2865– 2871.
- [109] D. D. Fan, J. Nguyen, R. Thakker, N. Alatur, A.-a. Agha-mohammadi, and E. A. Theodorou, *Bayesian learning-based adaptive control for safety critical systems*, 2019. arXiv: 1910.02325 [eess.SY].
- [110] A. Gahlawat, P. Zhao, A. Patterson, N. Hovakimyan, and E. Theodorou, "L1-gp: L1 adaptive control with bayesian learning," in *Proceedings of Machine Learning Research*, A. M. Bayen, A. Jadbabaie, G. Pappas, P. A. Parrilo, B. Recht, C. Tomlin, and M. Zeilinger, Eds., vol. 120, The Cloud: PMLR, Oct. 2020, pp. 826–837.
- [111] J. Zhang, P. Fiers, K. A. Witte, R. W. Jackson, K. L. Poggensee, C. G. Atkeson, and S. H. Collins, "Human-in-the-loop optimization of exoskeleton assistance during walking," *Science*, vol. 356, no. 6344, pp. 1280–1284, 2017. eprint: https://science. sciencemag.org/content/356/6344/1280.full.pdf.
- [112] R. W. Jackson and S. H. Collins, "Heuristic-based ankle exoskeleton control for coadaptive assistance of human locomotion," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 27, no. 10, pp. 2059–2069, 2019.

- [113] Z. Peng, R. Luo, R. Huang, J. Hu, K. Shi, H. Cheng, and B. K. Ghosh, "Data-driven reinforcement learning for walking assistance control of a lower limb exoskeleton with hemiplegic patients," in 2020 IEEE International Conference on Robotics and Automation (ICRA), IEEE, 2020, pp. 9065–9071.
- [114] Z. Huang, J. Liu, Z. Li, and C. Su, "Adaptive impedance control of robotic exoskeletons using reinforcement learning," 2016 International Conference on Advanced Robotics and Mechatronics (ICARM), pp. 243–248, 2016.
- [115] Y. Yuan, Z. Li, T. Zhao, and D. Gan, "Dmp-based motion generation for a walking exoskeleton robot using reinforcement learning," *IEEE Transactions on Industrial Electronics*, vol. 67, no. 5, pp. 3830–3839, 2020.
- [116] K. Yin, K. Loken, and M. Van de Panne, "Simbicon: Simple biped locomotion control," *ACM Transactions on Graphics (TOG)*, vol. 26, no. 3, 105–es, 2007.
- [117] J. K. Hodgins, W. L. Wooten, D. C. Brogan, and J. F. O'Brien, "Animating human athletics," in *Proceedings of the 22nd annual conference on Computer graphics and interactive techniques*, 1995, pp. 71–78.
- [118] S. Jain, Y. Ye, and C. K. Liu, "Optimization-based interactive motion synthesis," *ACM Transactions on Graphics (TOG)*, vol. 28, no. 1, pp. 1–12, 2009.
- [119] S. Coros, P. Beaudoin, and M. Van de Panne, "Generalized biped walking control," *ACM Transactions On Graphics (TOG)*, vol. 29, no. 4, pp. 1–9, 2010.
- [120] M. da Silva, Y. Abe, and J. Popović, "Interactive simulation of stylized human locomotion," in *ACM SIGGRAPH 2008 papers*, 2008, pp. 1–10.
- [121] Y. Ye and C. K. Liu, "Optimal feedback control for character animation using an abstract model," in *ACM SIGGRAPH 2010 papers*, 2010, pp. 1–9.
- [122] U. Muico, Y. Lee, J. Popović, and Z. Popović, "Contact-aware nonlinear control of dynamic characters," in ACM SIGGRAPH 2009 papers, 2009, pp. 1–9.
- [123] K. W. Sok, M. Kim, and J. Lee, "Simulating biped behaviors from human motion data," in *ACM SIGGRAPH 2007 papers*, 2007, 107–es.
- [124] Bioloidgp, http://en.robotis.com/.
- [125] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [126] PyDart, A python binding of dynamic animation and robotics toolkit, http://pydart2.readthedocs.io.

- [127] DART, Dynamic animation and robotics toolkit, http://dartsim.github.io/.
- [128] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *arXiv preprint arXiv:1506.02438*, 2015.
- [129] N. G. Tsagarakis, S. Morfey, G. M. Cerda, L. Zhibin, and D. G. Caldwell, "Compliant humanoid coman: Optimal joint stiffness tuning for modal frequency control," in *Robotics and Automation (ICRA)*, 2013 IEEE International Conference on, IEEE, 2013, pp. 673–678.
- [130] PyDART, A python binding of dynamic animation and robotics toolkit, http://pydart2.readthedocs.
- [131] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," pp. 1889–1897, 2015.
- [132] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," in *International Conference on Machine Learning*, 2016.
- [133] A. G. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning," *Discrete Event Dynamic Systems*, vol. 13, no. 1-2, pp. 41–77, Jan. 2003.
- [134] DART: Dynamic Animation and Robotics Toolkit.
- [135] Openai gym.
- [136] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne, "Deepmimic: Exampleguided deep reinforcement learning of physics-based character skills," ACM Transactions on Graphics (Proc. SIGGRAPH 2018), 2018.
- [137] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics (Proc. SIGGRAPH 2018)*, vol. 37, no. 4, 2018.
- [138] D. Martelli, V. Vashista, S. Micera, and S. K. Agrawal, "Direction-dependent adaptation of dynamic gait stability following waist-pull perturbations," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 24, no. 12, pp. 1304– 1313, Dec. 2016.
- [139] J. Lee, M. X. Grey, S. Ha, T. Kunz, S. Jain, Y. Ye, S. S. Srinivasa, M. Stilman, and C. K. Liu, "Dart: Dynamic animation and robotics toolkit," *The Journal of Open Source Software*, vol. 3, no. 22, p. 500, 2018.
- [140] J. M. Caputo and S. H. Collins, "A Universal Ankle–Foot Prosthesis Emulator for Human Locomotion Experiments," *Journal of Biomechanical Engineering*, vol. 136,

no. 3, Feb. 2014, 035002. eprint: https://asmedigitalcollection.asme.org/biomechanical/article-pdf/136/3/035002/3025302/bio $_136_03_035002$ .pdf.

- [141] K. A. Witte, J. Zhang, R. W. Jackson, and S. H. Collins, "Design of two lightweight, high-bandwidth torque-controlled ankle exoskeletons," in 2015 IEEE International Conference on Robotics and Automation (ICRA), May 2015, pp. 1223–1228.
- [142] Unitree's aliengo quadrupedal robot.
- [143] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, *Openai gym*, 2016.
- [144] C. Fukuchi, F. Reginaldo, and M. Duarte, "A public dataset of overground and treadmill walking kinematics and kinetics in healthy individuals," *Peer journal*, 2018.
- [145] Unitree's al quadrupedal robot.
- [146] K. H. Hunt and F. R. E. Crossley, "Coefficient of restitution interpreted as damping in vibroimpact," 1975.
- [147] A. Seth, J. L. Hicks, T. K. Uchida, A. Habib, C. L. Dembia, J. J. Dunne, C. F. Ong, M. S. DeMers, A. Rajagopal, M. Millard, *et al.*, "Opensim: Simulating musculoskeletal dynamics and neuromuscular control to study human and animal movement," *PLoS computational biology*, vol. 14, no. 7, e1006223, 2018.
- [148] N. Hansen, "The CMA evolution strategy: A tutorial," *CoRR*, vol. abs/1604.00772, 2016. arXiv: 1604.00772.
- [149] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science Robotics*, vol. 5, no. 47, 2020. eprint: https://robotics.sciencemag.org/content/5/47/eabc5986.full.pdf.
- [150] Z. Xie, X. Da, M. van de Panne, B. Babich, and A. Garg, "Dynamics randomization revisited: A case study for quadrupedal locomotion," *arXiv preprint arXiv:2011.02404*, 2020.
- [151] A. Falisse, L. Pitto, H. Kainz, H. Hoang, M. Wesseling, S. Van Rossom, E. Papageorgiou, L. Bar-On, A. Hallemans, K. Desloovere, *et al.*, "Physics-based simulations to predict the differential effects of motor control and musculoskeletal deficits on gait dysfunction in cerebral palsy: A retrospective case study," *Frontiers in human neuroscience*, vol. 14, p. 40, 2020.

- [152] G. F. Santos, E. Jakubowitz, N. Pronost, T. Bonis, and C. Hurschler, "Predictive simulation of post-stroke gait and effects of functional electrical stimulation: A case report," 2021.
- [153] J. Weng, E. Hashemi, and A. Arami, "Natural walking with musculoskeletal models using deep reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 4156–4162, 2021.

VITA

Visak Kumar was born in Bangalore, India in 1991. He graduated from National Public School in 2009 and attended M.S.Ramaiah Institute of Technology from 2009 to 2013 where he majored in Mechanical Engineering. After completing undergraduate degree, Visak worked at Indian Institute of Science, Bangalore, under supervision of Prof. Ananthasuresh from 2013 - 2014, as a research assistant where he worked on developing a micro-scale mechanical pump device for drug delivery.

Visak obtained his Masters degree from University of Washington, Seattle in 2016. After obtaining his masters degree, he started PhD program in robotics at Georgia Tech, under supervision of Dr. Karen Liu working on developing learning algorithms for robot control. In 2017, Visak worked as a research intern at Disney Research, Pittsburgh under supervision of Dr. Katsu Yamane and Dr. Sehoon Ha. In 2019, Visak interned at Nvidia Research, Seattle under supervision of Dr. Stan Birchfield and Dr. Jonathan Tremblay working on dextrous manipulation. He returned as an intern to Nvidia in summer of 2020.