

Reinforcement learning based control of an underactuated double pendulum system

Master's Thesis Nr. 0183

Scientific Thesis for Acquiring the Master of Science Degree
in the study program Mechatronic and Robotics at the School of Engineering
and Design of the Technical University of Munich.

Thesis Advisor Laboratory for Product Development and Lightweight Design
Prof. Dr. Markus Zimmermann

Supervisor Laboratory for Product Development and Lightweight Design
Akhil Sathuluri, Felix Wiebe, Prof.Dr. Shivesh Kumar
Prof.Dr. Frank Kirchner (Second corrector)

Submitted by Chi Zhang
Karl Köglspurger Straße 9, 80939, München
Matriculation number: 03735807
chi97.zhang@mytum.de

Submitted on Garching, 15.11.2023

Declaration

I assure that I have written this work autonomously and with the aid of no other than the sources and additives indicated.

Garching, 15.11.2023

Chi Zhang

Project Definition

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Background

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Acknowledgement

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Project Note

Master's Thesis	Nr. 0183
Supervisor	Akhil Sathuluri, Felix Wiebe, Prof.Dr. Shivesh
Kumar	
Partners in industry/research	DFKI GmbH, Robotics Innovation Center
Time period	15.05.2023 - 15.11.2023

The dissertation project of Akhil Sathuluri, Felix Wiebe, Prof.Dr. Shivesh Kumar set the context for the work presented. My supervisor Akhil Sathuluri, Felix Wiebe, Prof.Dr. Shivesh Kumar mentored me during the compilation of the work and gave continuous input. We exchanged and coordinated approaches and results monthly.

An accurate elaboration, a comprehensible and complete documentation of all steps and applied methods, and a good collaboration with industrial partners are of particular importance.

Publication

I consent to the laboratory and its staff members using content from my thesis for publications, project reports, lectures, seminars, dissertations and postdoctoral lecture qualifications.

The work remains a property of the Laboratory for Product Development and Lightweight Design.

Garching, 15.11.2023

Chi Zhang

Akhil Sathuluri, Felix Wiebe, Prof.Dr. Shivesh Kumar

Contents

1	Introduction	1
1.1	Motivation.....	1
1.1.1	Trajectory planning and tracking.....	2
1.1.2	Reinforcement Learning Based Control	4
1.2	Problem setup	6
1.3	Contribution	9
1.4	Content	9
2	State of the art	11
2.1	Theory	11
2.1.1	Variations of double pendulum	11
2.1.2	Dynamics of underactuated double pendulum system.....	13
2.2	Related work	16
3	Methodology	21
3.1	Soft actor critic	21
3.2	Linear quadratic regulator	24
3.3	Combining SAC and LQR with region of attraction.....	25
3.4	Reward shaping	27
3.5	Introduction to leaderboard metrics	30
3.5.1	Performance Leaderboard in Simulation and Real system	30
3.5.2	Simulation Robustness Leaderboard.....	32
4	Agent training and experiments in simulation environment	35
4.1	Training setup.....	35
4.2	Training phase	37
4.3	Simulation results	40
4.3.1	Pendubot simulation in ideal environment.....	40
4.3.2	Acrobot simulation in ideal environment	41
5	Experiments on real hardware	43
5.1	Hardware setup	43
5.2	System identification	49
5.3	Sim2Real problem.....	50
5.3.1	Validation with noisy simulation	51
5.3.2	Training with domain randomization	53

5.4	Real hardware results.....	54
5.4.1	Pendubot results	55
5.4.2	acrobot results	56
6	Discussion and Future Work	57
6.1	Interpretation of simulation leaderboard.....	57
6.2	Interpretation of robust leaderboard.....	58
6.3	Interpretation of real system leaderboard.....	59
6.4	Conclusion and future work.....	60
	Bibliography	61
	Appendix	65
	Appendix A An appendix.....	66

1 Introduction

Nonlinear systems, as their name suggests, do not exhibit linear relationships between inputs and outputs. This inherent non-linearity means the system's response to changes in input can be complex and often unpredictable. In the real world, nearly all systems display some degree of nonlinearity. This nonlinearity can manifest in a variety of phenomena. For example, in systems with multiple inputs and outputs, the interdependencies between variables can become intricate, leading to coupling issues. Another frequent challenge is chaotic behavior, where even minor alterations in initial conditions can lead to vastly different outcomes. Addressing these nonlinearities is crucial when managing control tasks, especially in real-world systems.

Achieving precise control over nonlinear systems has long been a primary focus in the field of control theory. While linear systems, often represented by linear differential equations, can typically be solved quickly and analytically, nonlinear equations representing nonlinear dynamics usually lack closed-form solutions. This necessitates the use of approximations and numerical methods. Efficiently and accurately executing these methods presents a central challenge in nonlinear control. Throughout history, control engineers have devised a wide range of strategies to manage these complex systems. The rise of robotics in recent decades has introduced new methods specifically tailored to address the challenges presented by nonlinearities.

1.1 Motivation

Robots are purposefully engineered as programmable mechanical structures, enabling them to perform a variety of tasks, either autonomously or under partial human oversight. These tasks include mobility, manipulation, and active interaction with their surroundings.

In the field of modern robotics, the mechanical systems are highly complex and nonlinear, posing significant challenges for precise and effective control. However, with the advancement of modern control methodologies and the increasing power of artificial intelligence, numerous innovative control approaches are emerging each year. Many products in modern robotics have achieved remarkable success, some well-known instances include quadruped robotics[9], autonomous vehicles[33], quadcopters[27], and humanoid robots[32].



(a) Quadruped



(b) Quadcopter



(c) Humanoid



(d) Autonomous vehicle

Figure 1.1: Four successful examples for controlling dynamic and nonlinear systems in the field of modern robotics

1.1.1 Trajectory planning and tracking

In the traditional control of complex systems such as robotics, which exhibit significant nonlinearity, it is crucial to employ dependable nonlinear control strategies to enhance motion

capabilities. Typically, the procedure for planning intricate movements within these systems adheres to a two-phase method[7], encompassing trajectory planning and trajectory tracking. This structured approach ensures precision and stability in robotic system control.

Trajectory planning[15][16] involves calculating a smooth and feasible path for the robot to follow, aiming for a specific target position or operational point. A widely-used method in this stage is trajectory optimization[6], which seeks to minimize a cost function accounting for travel time, energy consumption, and motion smoothness. Techniques such as gradient descent or genetic algorithms are often utilized for optimization, with adherence to constraints like maximum velocities and accelerations being a crucial aspect of the process.

Upon successful trajectory planning, the next step is trajectory tracking, which entails the use of control algorithms to guide the robot along the predetermined path. A feedback control approach is typically employed, continually monitoring the robot's position and adjusting control inputs as necessary. However, in real-world systems, external disturbances, uncertainties, and system limitations may lead to significant deviations from the planned trajectory. Therefore, ensuring accurate state estimation and implementing robust control strategies are imperative in feedback control to address these challenges.

In industrial robot control[17], the practice typically involves dividing the control process into distinct phases of offline planning and execution. This segmentation is primarily due to the non-critical requirement for real-time responsiveness in such applications. Take, for instance, the planning phase, where an algorithm such as A* is employed to determine an optimal route based on a predefined task, navigating around obstacles and minimizing travel distance, which results in a set of discrete waypoints[13]. Following this, cubic interpolation techniques are applied to these waypoints, crafting a smooth and continuous trajectory that the robot can realistically follow, ensuring both fluid motion and compliance with the robot's kinematic constraints[8]. The process subsequently moves to the execution phase, where a robust and precise control strategy is implemented. This strategy is crucial as it guarantees the robot's meticulous adherence to the pre-established trajectory, maintaining both accuracy and reliability throughout the operation[28].

Yet, in dynamic domains like automotive and flight control, the demand for real-time responsiveness takes center stage. For example, Model Predictive Control (MPC)[22] stands as a prime illustration of this critical requirement, seamlessly integrating online trajectory planning and execution into a unified framework. It initiates the process by forecasting a series of future control actions as part of the planning phase. Subsequently, during execution, these control inputs are refined to minimize any discrepancies between the real-time system state and the intended trajectory, ensuring precise alignment. MPC excels in its ability to simultaneously create, optimize, and follow trajectories, all while dynamically adjusting in real-time to accom-

moderate disturbances and uncertainties. This capability is vital for maintaining precise control and adaptability in fast-paced and complex environments.

The traditional trajectory generation and tracking approach, while effective for numerous systems, has its own set of limitations.

- **Limited Adaptability:** Trajectory planning typically relies on predefined paths or trajectories, limiting adaptability to unforeseen changes or dynamic environments. If the environment changes significantly, the planned trajectory may no longer be optimal or even feasible.
- **Difficulty with Nonlinear Systems:** Trajectory planning struggles with highly nonlinear systems where the dynamics are hard to model accurately. Linearizing the system for planning purposes may lead to suboptimal or infeasible trajectories.
- **High Computational Demands:** Some trajectory planning algorithms can be computationally intensive, especially for high-dimensional or complex robotic systems. This computational demand becomes a drawback, particularly in real-time or time-critical applications.

Recognizing the limitations inherent in traditional trajectory generation and tracking methodologies is essential for fostering the development of more advanced and efficient strategies in trajectory planning and control. This insight becomes particularly crucial in contexts that require rapid responsiveness and a capacity to skillfully navigate unpredictable changes or dynamic environments.

1.1.2 Reinforcement Learning Based Control

Transitioning from a focus on predefined trajectories, reinforcement learning (RL)-based control presents an alternative framework.

Reinforcement learning (RL) is a subset of machine learning centered on agents learning optimal behavior through trial-and-error interactions with their environment. Essentially, the agent makes sequential decisions, observes the outcomes of its actions, and receives feedback in the form of rewards or penalties. This feedback serves as a guiding mechanism, enabling the agent to evaluate the outcomes of its actions and adjust its policy accordingly to maximize cumulative rewards over time.

The mathematical framework that describes RL is the Markov Decision Process (MDP), which provides a structured model for the decision-making problem. In an MDP, the agent's current state, available actions, potential next states, and the rewards associated with state-action transitions are all clearly defined. This structure ensures that the future state of the system depends only on the current state and the action taken, fulfilling the Markov property.

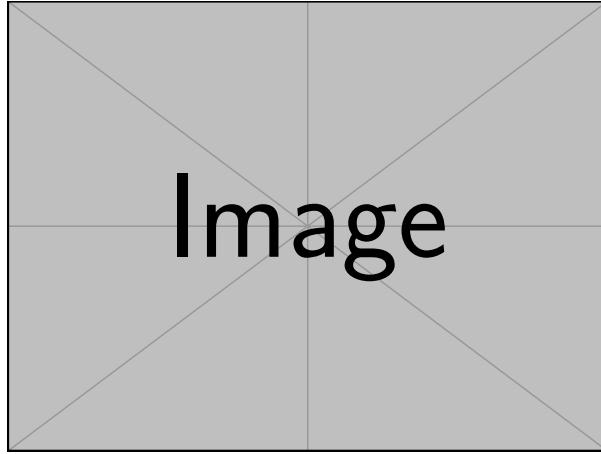


Figure 1.2: Markov chain

The ultimate objective in RL is to find an optimal policy, a mapping from states to actions, that dictates the best action to take in each state to achieve maximum long-term rewards. The learning process is driven by the tuple (s, a, r, s') , representing the current state, action taken, received reward, and subsequent state. Over time, through exploration of the state-action space and exploitation of acquired knowledge, the agent refines its policy, converging towards optimal behavior.

A key challenge in RL is balancing exploration and exploitation. Exploration involves trying out different actions to discover their potential outcomes, essential for acquiring new knowledge. Exploitation, on the other hand, entails leveraging existing knowledge to make optimal decisions. Striking the right balance between these two strategies is crucial for the agent to effectively navigate its environment and learn from received rewards.

In the field of robotic control, the integration of reinforcement learning has become increasingly popular due to the numerous distinct advantages this combined approach offers:

- **Adaptability and Flexibility:** RL allows systems to adapt and improve in dynamic environments. Its control policy evolves through accumulating new experiences and

knowledge from interactions with the environment, making it highly adaptable to varying circumstances.

- **Reduced Dependency on Accurate Models:** In contrast to traditional control methods that depend on exact mathematical models, RL thrives without the necessity of a predetermined model, particularly when learning from real-world data. This feature is immensely beneficial in situations where the system dynamics are either intricate or not well understood, as RL refines its performance through direct interactions with the environment.
- **Effective Handling of Nonlinearities and High-Dimensional Systems:** RL excels at managing nonlinear systems and complex control tasks using neural network-based function approximation. This enables it to navigate high-dimensional input spaces and control intricate relationships between states and actions in complex scenarios.

Reinforcement learning distinguishes itself as a straightforward control method with high adaptability, setting itself apart from traditional trajectory generation and tracking techniques. While conventional methods frequently struggle with challenges such as limited adaptability, nonlinearity in systems, and intensive computational demands, reinforcement learning addresses these issues through direct interactions with the environment. It excels in unpredictable conditions and in managing complex, high-dimensional systems. This is largely attributed to its reduced dependency on accurately predefined dynamic models, ensuring robustness, effectiveness, and versatility across a wide range of applications.

1.2 Problem setup

In the realm of nonlinear systems, underactuated systems[26] present a particularly challenging class. These systems are characterized by having fewer control inputs than degrees of freedom, or their control inputs are constrained in some way. This characteristic makes them significantly more challenging to control when compared to fully actuated systems. What's intriguing is that a majority of robots and even natural organisms fall into the category of underactuated systems. Consequently, the study of underactuated mechanical systems' control holds universal relevance.

The concept of underactuated dynamics can be briefly introduced as follows. According to Newton's second law ($F = ma$), the dynamics that govern any mechanical system can be mathematically expressed as shown in Equation 1.1:

$$\ddot{q} = f(q, \dot{q}, u, t) \quad (1.1)$$

In this equation, \ddot{q} represents acceleration, which is the second derivative of the variable q , typically representing the position of the system. The function $f(q, \dot{q}, u, t)$ describes how \ddot{q} depends on various parameters, including the state variables q and \dot{q} , the control inputs u , and the time t .

The state of the system, denoted as s and represented as $[q, \dot{q}]^T$, consists of two vectors: q , which signifies positions, and \dot{q} , which signifies velocities.

When dealing with control-related tasks, we can express the second-order differential equation in the following manner:

$$\ddot{q} = f_1(q, \dot{q}, t) + f_2(q, \dot{q}, t)u \quad (1.2)$$

In this equation, $f_1(q, \dot{q}, t)$ corresponds to one part of the function that affects \ddot{q} , while $f_2(q, \dot{q}, t)$ represents another part that interacts with the control input u .

In the context of a controlled dynamical system, as described by Equation 1.2, we evaluate the condition of underactuation at specific states $[q, \dot{q}]^T$ at time t . Underactuation can be identified through two distinct scenarios:

- **Case 1:**

The system is classified as underactuated at a particular state if the rank of the matrix $[f_2(q, \dot{q}, t)]$ is less than the dimension of q . This condition is expressed as:

$$\text{rank}[f_2(q, \dot{q}, t)] < \dim[q] \quad (1.3)$$

- **Case 2:**

Alternatively, underactuation may also arise even when f_2 is full rank, provided there are additional constraints on the control inputs. For instance, if constraints such as $|\mathbf{u}| \leq 1$

limit the control inputs, the system's controllability can still be restricted, resulting in underactuation.

These principles are explained in the work by Russ Tedrake[38].

A well-discussed example of underactuated control is found in the double pendulum—a simple setup consisting of two links connected by two rotational joints. These joints include the shoulder joint, which is directly connected to the world frame, and the elbow joint, situated between the two links. The end effector is located at the tip of the second link. Active control is achieved by attaching motors to the shoulder and elbow joints. In the domain of underactuated control, if the shoulder joint is actuated, the setup is referred to as a pendubot (see Figure 1.3a). Conversely, if the elbow joint is actuated, it is known as an acrobot (see Figure 1.3b).

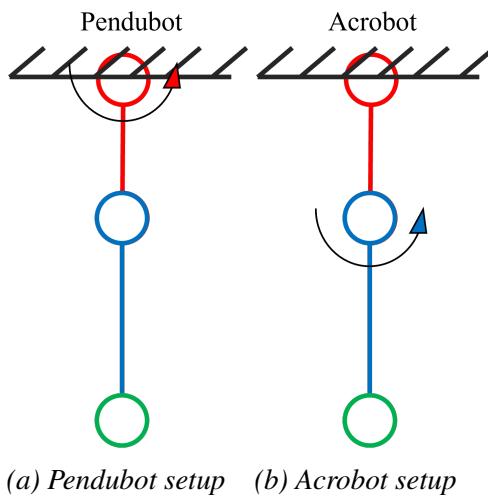


Figure 1.3: Two variations of underactuated double pendulum

Despite its simple configuration, the system exhibits highly nonlinear and chaotic behavior[34]. The double pendulum setup presents two classic tasks: swing-up and stabilization around the highest point. Research on swing-up and stabilization of the double pendulum can be traced back to the 1990s[43], and it continues to be a crucial testbed for validating the effectiveness of newly designed control algorithms[42][44][3].

Our project's objective is to develop a reinforcement learning-based control method suitable for underactuated control of the double pendulum system, specifically addressing swing-up and stabilization tasks. To evaluate the efficacy of this control method, we conduct both simulations and real system experiments.

1.3 Contribution

In this paper, the main contribution is as follows: An effective control strategy has been developed to achieve two key objectives with the double pendulum. The first task involves swinging the double pendulum from its lowest point to its highest point. The second task is to ensure stability at the highest point for an extended time period.

To address the swing-up task, a well-known model-free reinforcement learning algorithm called Soft Actor-Critic (SAC) was utilized[18]. This algorithm enabled the training of a policy capable of reaching the Region of Attraction (RoA) of a continuous-time linear quadratic regulator (LQR) controller[24]. Once the RoA is reached by the system, a seamless transition to the LQR controller is made to maintain stability around the highest point.

1.4 Content

This thesis comprises six chapters, namely: introduction, state of the art, methodology, agent training and experiments in a simulation environment, experiments on real hardware, discussion, and future work. At the end of Chapter One, the content of the following chapters is outlined below.

- **Chapter 2: State of the Art:**
 - This chapter begins by providing an explanation of the fundamental theories related to double pendulum dynamics. It is followed by an overview of recent advancements in the field of robotic control, with a specific focus on learning-based control methods.
- **Chapter 3: Methodology:**
 - This chapter delves into the methodology, covering fundamental aspects of reinforcement learning, with a specific focus on the SAC algorithm. It explains the reward function used for training, the training procedure, and introduces the concept of the LQR controller and the composition of the combined controller.

- Additionally, we introduce the evaluation metrics used to assess the performance and robustness of the newly designed control strategies in both simulated environments and real-world experiments.
- **Chapter 4: Agent training and experiments in simulation environment:**
 - In this chapter, we present the results obtained from simulations, showcasing the performance and behavior of the designed control strategy.
- **Chapter 5: Experiments on real hardware:**
 - This chapter reports the outcomes of experiments conducted on the hardware, providing insights into our approach to addressing the sim2real transfer problem.
- **Chapter 6: Discussion and Future Work:**
 - The final chapter engages in a discussion about the obtained results and explores potential future research and development directions.

2 State of the art

In this chapter, we will explore the state of the art. In the first section, we will first discuss the important variations of double pendulum that is widely used by the research society. Subsequently, we will delve into the basic dynamics of a double pendulum system. In the second section, we will review some of the most renowned works related to learning-based control in the field of robotics.

2.1 Theory

In this section, we will first introduce the various variations of double pendulum setups commonly used in research society. Subsequently, we will explore the fundamental principles governing double pendulum dynamics.

2.1.1 Variations of double pendulum

Broadly speaking, a double pendulum system is a mechanical structure consisting of two pendulum arms or masses suspended in such a way that they can swing freely and independently from each other. These pendulum arms are typically connected in a series, where the motion of the second pendulum is influenced by the motion of the first pendulum.

There are several variations of the double pendulum setup, and one of the most famous ones is the Double Inverted Pendulum on a Cart (DIPC)[10]. The DIPC system is a modification of the pendubot setup, as it is actuated solely at the shoulder joint. The key distinction lies in the actuation mechanism, which employs a prismatic joint instead of a revolute joint.

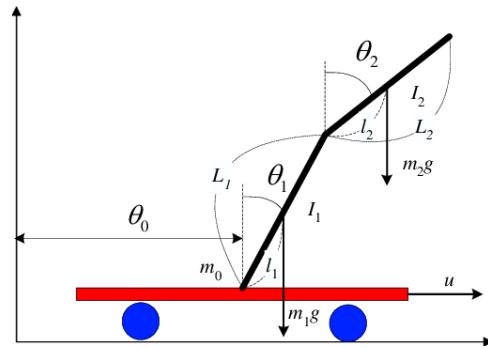


Figure 2.1: Double inverted pendulum on a cart[10]

Another significant variation is the Furuta pendulum[12], initially developed at the Tokyo Institute of Technology by Furuta and his colleagues[14]. This system comprises a driven arm that rotates horizontally and a pendulum attached to this arm, allowing free vertical rotation. Due to the influence of gravitational forces and the coupling stemming from Coriolis and centripetal forces, the Furuta pendulum is characterized by its underactuation and highly nonlinear behavior.

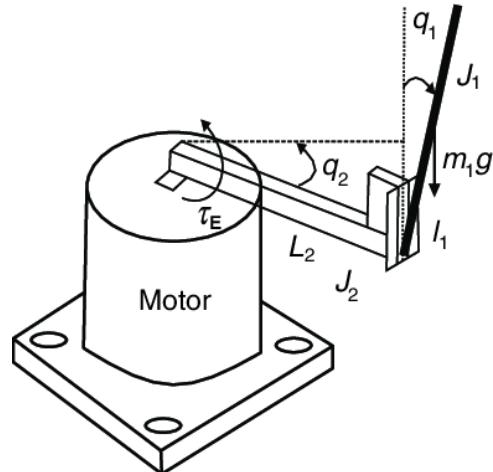


Figure 2.2: Furuta pendulum[2]

The double pendulum utilized in this thesis represents a third variation. It consists of two links connected in series by revolute joints. Unlike the Furuta pendulum, both links of the double pendulum move in the same plane within three-dimensional space. These two links are

connected to the world frame via revolute joints as well. In contrast to the DIPC setup, the actuation capability is solely limited by the torque that the joints can generate, unrestricted by the length of a prismatic rail. The dynamics of the double pendulum system will be discussed in the following section.

2.1.2 Dynamics of underactuated double pendulum system

As shown in Figure 2.3, we model the dynamics of the double pendulum with 15 parameters which include 8 link parameters namely masses (m_1, m_2) , lengths (l_1, l_2) , center of masses (r_1, r_2) , inertias (I_1, I_2) for the two links, and 6 actuator parameters namely motor inertia I_r , gear ratio g_r , coulomb friction (c_{f1}, c_{f2}) , viscous friction (b_1, b_2) for the two joints and gravity.

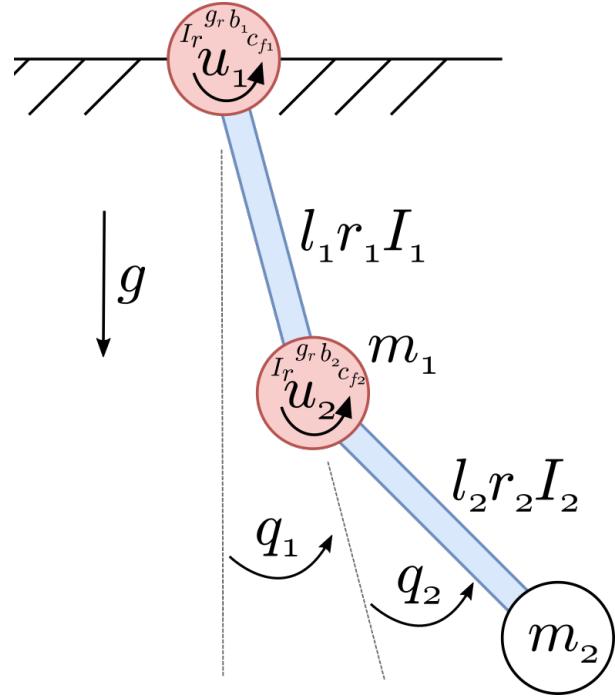


Figure 2.3: Double pendulum dynamics

The generalized coordinates $\mathbf{q} = [q_1, q_2]^T$ are the joint angles measured from the free hanging position. The state vector of the systems contains the position coordinates and their time

derivatives: $\mathbf{x} = [\mathbf{q}, \dot{\mathbf{q}}]^T$. The torque applied by the actuators are $\mathbf{u} = [u_1, u_2]$. The equation of motion for the dynamics of a dynamical system can be derived following the blow steps:

Step 1. Define the Lagrangian (L):

The Lagrangian (L) is defined as the difference between the kinetic energy (T) and the potential energy (U) of the system:

$$L = T - U \quad (2.1)$$

Step 2. Express the Kinetic Energy (T):

The kinetic energy (T) of the double pendulum is the sum of the kinetic energies of both links. The kinetic energy for a link is given by:

$$T = \frac{1}{2}m_1(\dot{x}_1^2 + \dot{y}_1^2) + \frac{1}{2}m_2(\dot{x}_2^2 + \dot{y}_2^2) \quad (2.2)$$

where m_1 and m_2 are the masses of the links, (x_1, y_1) and (x_2, y_2) are their positions, and $\dot{x}_1, \dot{y}_1, \dot{x}_2, \dot{y}_2$ are their respective velocities.

Step 3. Express the Potential Energy (U):

The potential energy (U) of the double pendulum is the sum of the potential energies of both links. The potential energy for a link in a gravitational field is given by:

$$U = m_1gy_1 + m_2gy_2 \quad (2.3)$$

where g is the acceleration due to gravity.

Step 4. Formulate the Lagrange's Equation:

Use Lagrange's equation to derive the equations of motion for the generalized coordinates x_1, y_1, x_2, y_2 .

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_i} \right) - \frac{\partial L}{\partial q_i} = 0 \quad (2.4)$$

Step 5. Solve the Equations of Motion:

Solve the obtained set of second-order differential equations to determine the equations of motion for the system. The system dynamics with friction is:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} = Du + G(q) - F(\dot{q}) \quad (2.5)$$

Because the state vector is $\mathbf{x} = [\mathbf{q}, \dot{\mathbf{q}}]^T$, the equation of motion can also be expressed as:

$$\begin{aligned} \dot{\mathbf{x}} &= f(x, u) \\ &= \begin{bmatrix} \dot{q} \\ M^{-1}(Du - C(q, \dot{q})\dot{q} + G(q) - F(\dot{q})) \end{bmatrix} \end{aligned} \quad (2.6)$$

Consider the forward kinematics of double pendulum system, the coordinate of the joint between the first link and the second link is $P_1 = (x_1, y_1)$, the coordinate of the end effector is $P_2 = (x_2, y_2)$.

$$\begin{cases} x_1 = l_1 \sin(q_1) \\ y_1 = -l_1 \cos(q_1) \end{cases} \quad (2.7)$$

$$\begin{cases} x_2 = l_1 \sin(q_1) + l_2 \sin(q_1 + q_2) \\ y_2 = -l_1 \cos(q_1) - l_2 \cos(q_1 + q_2) \end{cases} \quad (2.8)$$

Put Equation 2.7 and 2.8 into 2.2, 2.3, 2.4, 2.5, we can get the mass matrix (with $s_1 = \sin(q_1)$, $c_1 = \cos(q_1)$, ...)

$$\mathbf{M} = \begin{bmatrix} I_1 + I_2 + l_1^2 m_2 + 2l_1 m_2 r_2 c_2 + g_r^2 I_r + I_r & I_2 + l_1 m_2 r_2 c_2 - g_r I_r \\ I_2 + l_1 m_2 r_2 c_2 - g_r I_r & I_2 + g_r^2 I_r \end{bmatrix} \quad (2.9)$$

The Coriolis matrix:

$$\mathbf{C} = \begin{bmatrix} -2\dot{q}_2 l_1 m_2 r_2 \sin(q_2) & -\dot{q}_2 l_1 m_2 r_2 \sin(q_2) \\ \dot{q}_1 l_1 m_2 r_2 \sin(q_2) & 0 \end{bmatrix}, \quad (2.10)$$

The gravity vector:

$$\mathbf{G} = \begin{bmatrix} -g m_1 r_1 \sin(q_1) - g m_2 (l_1 \sin(q_1) + r_2 \sin(q_{1+2})) \\ -g m_2 r_2 \sin(q_{1+2}) \end{bmatrix}, \quad (2.11)$$

The friction vector:

$$\mathbf{F} = \begin{bmatrix} b_1 \dot{q}_1 + c_{f1} \arctan(100 \dot{q}_1) \\ b_2 \dot{q}_2 + c_{f2} \arctan(100 \dot{q}_2) \end{bmatrix} \quad (2.12)$$

(the $\arctan(100 \dot{q}_i)$ function is used to approximate the discrete step function for the coulomb friction)

And the actuator selection matrix \mathbf{D} :

$$\mathbf{D}_{full} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{D}_{pendu} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{D}_{acro} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad (2.13)$$

for the fully actuated system, the pendubot or the acrobot.

2.2 Related work

This thesis focuses on learning-based robotic motion control, a field that has garnered increasing attention in recent years, with significant contributions from various research institutes.

Today, much of the research in this field is conducted within simulation environments due to their cost-effectiveness and the ability to facilitate rapid iteration. One noteworthy project from

2018 is the DeepMimic project[30], undertaken by researchers at the University of California, Berkeley. This work resides at the intersection of deep reinforcement learning, imitation learning, and robotics.

The DeepMimic project utilizes physics-based simulations to successfully replicate the diverse range of behaviors exhibited in the real world by 3D characters. These characters include real-world examples such as humanoid and Atlas robotics, as well as fictional characters like T-Rexes and dragons. Instead of relying on manually designed controllers, the project employs deep reinforcement learning methods to generalize to new skills and situations, often without human intervention.

In the training process, the agent is provided with reference data recorded by motion capture actors or keyframed animations. Through imitation learning, the agent is guided to achieve specific predefined goals. The central contribution of this project lies in its framework, which combines goal-directed reinforcement learning with reference data generated by humans. This combination enables the imitation of a wide variety of motion skills.

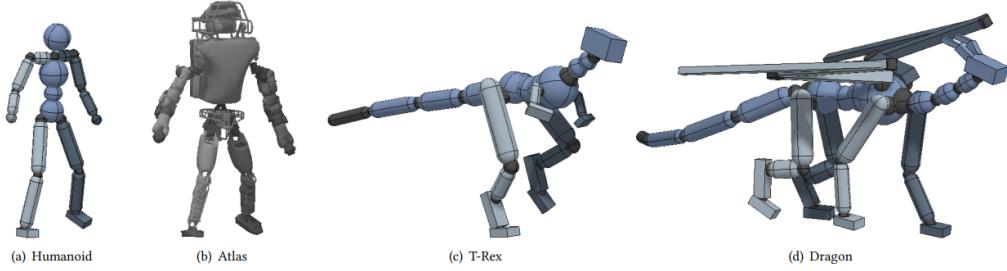


Figure 2.4: 3D characters in Deepmimic project[30]

In the context of reinforcement learning, a distinction exists between model-free and model-based approaches. Model-free reinforcement learning(MFRL) does not require information about the transition model, whereas model-based reinforcement learning(MBRL) leverages the transition model to make decisions based on a prior known model or one learned from interactions.

The model-free approach is notably characterized by its sample inefficiency. Successful training often takes many hours, if not days or weeks.

An intriguing project that relies solely on model-free deep reinforcement learning is the Learning-to-Walk-in-20 Minutes project[35] conducted by researchers from UC Berkeley. They employ

an algorithmic framework closely related to DroQ [19], an extension of the SAC algorithm [18] incorporating dropout [36] and layer normalization [5]. Remarkably, their training is conducted directly on the real system. They demonstrate that current deep RL methods can effectively teach quadrupedal locomotion in under 20 minutes, a stark contrast to previous research conducted by Kumar et al. [23], which employed the same hardware but required 1.2×10^9 samples to train a robust controller for locomotion. This corresponds to roughly 4.5 months' worth of cumulative experience.

Model-based reinforcement learning is recognized for its integration of an environment model with trial-and-error learning. One notable advantage is the potential for higher sample efficiency compared to the model-free approach. This work, conducted by the University of Toronto[41], provides a comprehensive comparison by benchmarking model-based reinforcement learning.

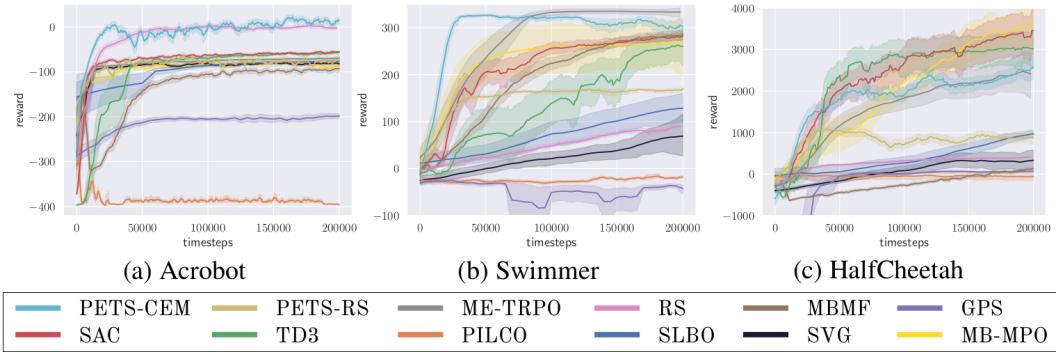


Figure 2.5: Benchmarking model-based reinforcement learning[41]

The team has collected 11 Model-Based Reinforcement Learning (MBRL) algorithms and 4 Model-Free Reinforcement Learning (MFRL) algorithms across 18 benchmarking environments specifically designed for MBRL in simulation. These benchmarking environments are based on the standard OpenAI Gym[11], ranging from simple 2D tasks like cart-pole to complex setups like humanoid. The benchmarking process is further extended by introducing noise into the environment, including disturbances in observations and actions.

The team has discovered that while 1 million time-step training is common for MFRL algorithms, many MBRL algorithms converge much earlier, often within 200k time-steps. Within the field of MBRL, when it comes to the evaluation of sample efficiency, asymptotic performance, and robustness, there is no clear and consistent best MBRL algorithm. This leaves ample opportunities for future research to leverage the strengths of different approaches.

Another notable challenge in reinforcement learning-based robotic control is the sim-to-real gap problem. Since agents are typically trained in carefully designed simulated worlds, these simulations can sometimes be idealized or oversimplified to some extent. Consequently, the optimal policy derived from the simulation often fails to account for the uncertainties of the real world, leading to failures when executing intended tasks. There are several approaches to bridging the sim-to-real gap, and the following work presents an interesting approach.

A research team from ETH Zurich has made significant progress in addressing the sim-to-real interface challenge [20]. Their methodology involves training a control model for a quadruped robot within a simulation environment. By utilizing a neural network and leveraging data collected from the real robot, they approximated the dynamics model of the physical robot. This approach has facilitated the accurate implementation of the control policy derived from the virtual environment onto the real robot. These exemplary works demonstrate the promising applications of learning-based approaches in various aspects of robot control.

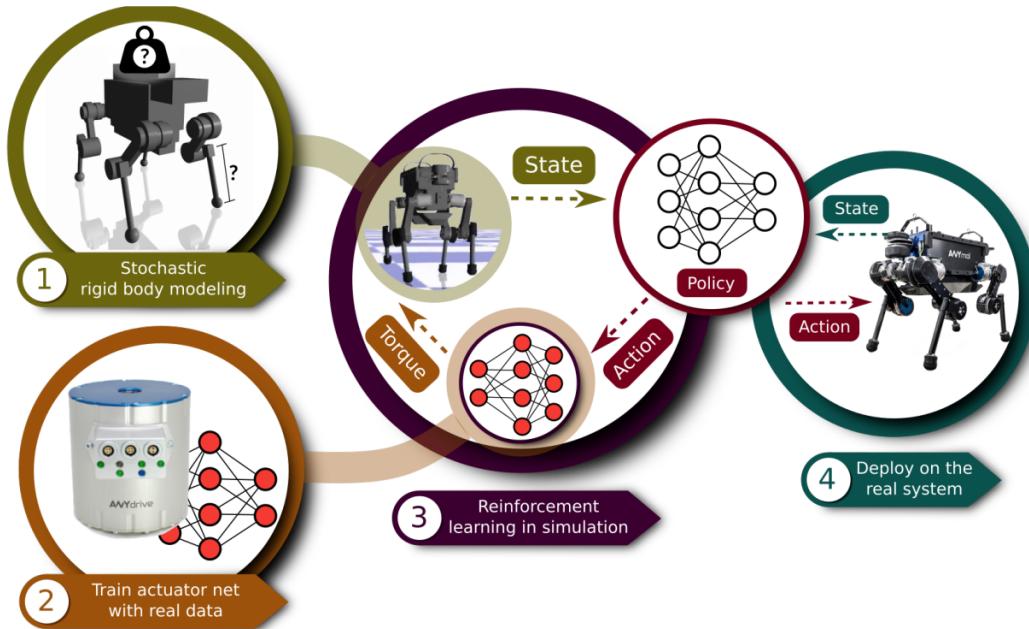


Figure 2.6: Sim-to-real framework for RL-based quadruped controller[20]

In conclusion, in the field of reinforcement learning-based control in robotics, researchers currently face three major challenges:

- 1. Sample Inefficiency in Model-Free Reinforcement Learning:**

Model-free reinforcement learning is known for its sample inefficiency, where the time required for random trial and error can be excessively long.

2. Immaturity of Model-based RL

Model-based reinforcement learning is a growing force but still in its infancy. Finding the right balance between model information and reinforcement learning remains an ongoing challenge.

3. Sim-to-Real Gap Limitations:

Deploying controllers learned through reinforcement learning on real systems is limited due to the significant sim-to-real gap. This limitation confines a substantial portion of research to simulations.

The research community has much work to do before reinforcement learning becomes a stable and widely accepted approach in the industry.

3 Methodology

The primary goal of this thesis is to achieve the swing-up movement in the pendubot or acrobot setup and to maintain stability around the highest points. Challenges have been revealed in initial training trials with reinforcement learning, including potential entrapment in local minima and difficulty in maintaining stability at the highest point for extended periods.

To address these challenges, two main strategies are adopted. For the stabilization issue, a combined controller is introduced. During the swing-up process, control is assumed by an RL-trained agent, utilizing the Soft Actor-Critic—a classic model-free reinforcement learning algorithm—based on its learned policies. However, as the system nears the maximum point, a seamless transition occurs. This shift allows for the takeover of control by a continuous-time LQR controller, ensuring the final stabilization required to maintain stability at the highest point.

3.1 Soft actor critic

Within the landscape of reinforcement learning, the Soft Actor Critic (SAC)[18] stands out as an algorithm specifically designed for environments with continuous state and action spaces. Such environments, exemplified by our double pendulum system where actuators can be adjusted to any value within the torque limit range, and state measurement can take any real number, influenced our decision to adopt SAC.

Like many other deep reinforcement learning algorithms, SAC optimizes a policy by maximizing the expected cumulative reward the agent obtains over time. This optimization is primarily achieved through an actor-critic structure[21].

The actor determines the best actions by interpreting the current environmental conditions and adhering to the existing policy. Typically, the actor is visualized as a shallow neural network that approximates the mapping between the input state and the output probability distribution over actions. Furthermore, SAC incorporates a stochastic policy within its actor, which fosters exploration and aids the agent in refining its policies.

3 Methodology

On the other hand, the critic evaluates the value of state-action pairs. It estimates the expected cumulative reward the agent can achieve by following a particular policy. More often than not, the critic is depicted as a neural network that processes state-action pairs as inputs to yield the estimated value.

A distinguishing feature of SAC, besides the actor-critic framework, is entropy regularization[1]. SAC utilizes a stochastic policy. This means that instead of always settling on a single best action for each state, the agent considers a probability distribution over potential actions. The incorporation of entropy in SAC aims to encourage exploration: high entropy signifies a more uniform distribution, implying the agent's uncertainty and tendency to explore diverse actions, while low entropy points to a concentrated distribution, suggesting the agent's confidence in a specific action. By definition, entropy quantifies randomness. Within SAC, it captures the unpredictability of the policy's action distribution. If x is a random variable with a probability density function P , the entropy H of x is defined as:

$$H(P) = \mathbb{E}_{x \sim P} [-\log P(x)] \quad (3.1)$$

By maximizing entropy, SAC encourages exploration and accelerates learning. It also prevents the policy from prematurely converging to a suboptimal solution. The trade-off between maximizing reward and maximizing entropy is controlled through a parameter, α . This parameter serves to balance the importance of exploration and exploitation within the optimization problem. Each interaction between the agent and the environment can be recorded as a tuple (s_t, a_t, R, s_{t+1}) . The optimal policy π^* can be defined as follows:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(R(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot | s_t)) \right) \right] \quad (3.2)$$

Where s_t and a_t represent the state and action at time t , and s_{t+1} represents the state at time $t + 1$. R denotes the immediate reward received by the agent after taking action a_t in state s_t , while γ is the discount factor that determines the agent's emphasis on long-term cumulative rewards over immediate ones.

During training, SAC learns a policy π_θ and two Q-functions Q_{ϕ_1}, Q_{ϕ_2} concurrently. The loss functions for the two Q-networks are ($i \in 1, 2$):

$$L(\phi_i, D) = \mathbb{E}_{(s,a,r,s',d) \sim D} \left[\left(Q_{\phi_i}(s, a) - y(r, s', d) \right)^2 \right], \quad (3.3)$$

where the temporal difference target y is given by:

$$\begin{aligned} y(r, s', d) &= r + \gamma(1-d) \times \left(\min_{j=1,2} Q_{\phi_{targ,j}}(s', \tilde{a}') - \alpha \log \pi_{\theta}(\tilde{a}' | s') \right), \\ \tilde{a}' &\sim \pi_{\theta}(\cdot | s') \end{aligned} \quad (3.4)$$

In each state, the policy π_{θ} should act to maximize the expected future return Q while also considering the expected future entropy H . In other words, it should maximize $V^{\pi}(s)$:

$$V^{\pi}(s) = \mathbb{E}_{a \sim \pi} [Q^{\pi}(s, a)] + \alpha H(\pi(\cdot | s)) \quad (3.5)$$

$$= \mathbb{E}_{a \sim \pi} [Q^{\pi}(s, a) - \alpha \log \pi(a | s)] \quad (3.6)$$

The interaction experiences are stored as tuples (s, a, r, s', d) in a replay buffer D . Here, d represents the signal for episode termination. When it is time to update the Q-value and policy, a batch of transitions $B = \{(s, a, r, s', d)\}$ is randomly sampled from D . The Q-functions are updated by one step of gradient descent using:

$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi_i}(s, a) - y(r, s', d))^2 \quad \text{for } i = 1, 2 \quad (3.7)$$

The policy is updated using one step of gradient ascent:

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} \left(\min_{i=1,2} Q_{\phi_i}(s, \tilde{a}_{\theta}(s)) - \alpha \log \pi_{\theta}(\tilde{a}_{\theta}(s) | s) \right) \quad (3.8)$$

where $\tilde{a}_{\theta}(s)$ is a sample of action derived from $\pi_{\theta}(\cdot | s)$ which is differentiable with respect to θ . The actor and critic neural networks' parameters are updated, resulting in policy adaptation.

3 Methodology

In conclusion, SAC's combination of stochastic policies, exploration through entropy regularization, value estimation, and gradient-based optimization make it a well-suited algorithm for addressing the challenges posed by continuous state and action spaces.

3.2 Linear quadratic regulator

The Linear Quadratic Regulator (LQR)[24] is an effective control method primarily designed for linear systems. Yet, when dealing with nonlinear dynamics, it remains applicable. The nonlinear system is linearized around a selected operating point, and based on this linearized version, the LQR controller can be sculpted.

Taking a step back, the general form of a nonlinear system defined by state vector x and input vector u can be expressed as:

$$\dot{x}(t) = f(x(t), u(t)) \quad (3.9)$$

In certain applications, such as pendubot or acrobot stabilization, it is important to select the appropriate operating point. Due to the tasks to be completed, the operating point is chosen around the upright position, specifically $x_{op} = [\pi, 0, 0, 0]^T$. Around this point, the system can be linearized, leading to:

$$\dot{\bar{x}}(t) = A\bar{x}(t) + Bu(t) \quad (3.10)$$

Here, the deviation from the desired state is given by $\bar{x} = x - x_{op}$. Its first derivative, $\dot{\bar{x}} = \dot{x} - \dot{x}_{op} = \dot{x}$, for x_{op} being a constant. The linearized state matrices A and input matrix B around the operation point are derived as:

$$A = \left. \frac{\partial f}{\partial x} \right|_{x=x_{op}}, \quad B = \left. \frac{\partial f}{\partial u} \right|_{x=x_{op}} \quad (3.11)$$

To derive an optimal control strategy, a quadratic cost function J is needed:

$$J = \int_0^T (x^T Q x + u^T R u) dt \quad (3.12)$$

The matrices Q and R must be chosen to be symmetric and positive definite, which means $Q = Q^T$ with all its eigenvalues positive, and similarly, $R = R^T$ with all positive eigenvalues.

The Hamilton-Jacobi-Bellman equation (HJB) characterizes the optimal value function $J^*(x)$, representing the minimum cost-to-go from state x to termination at T . The HJB equation is as follows:

$$\frac{dJ^*(x)}{dt} = \min_u \left(x^T Q x + u^T R u + \frac{dJ^*(x)}{dx} (Ax + Bu) \right) \quad (3.13)$$

By calculating $\frac{dJ^*(x)}{dt}$ from Equation 3.12, the HJB equation can be reduced to the continuous-time algebraic Riccati equation:

$$A^T S + S A - S B R^{-1} B^T S + Q = 0 \quad (3.14)$$

Finally, the LQR control law obtained with the above method for this linearized system is:

$$u(t) = -K\bar{x}(t) \quad (3.15)$$

where $K = R^{-1}B^T S$.

For a LQR controller like this, torques will always be applied to steer the system state toward the origin of the linear system.

3.3 Combining SAC and LQR with region of attraction

some other guy[sean gillen] has done a similar work, cite that.

In our approach, a combined control method is utilized for both the swing-up and stabilization tasks. During the swing-up phase, the SAC controller is employed. Once the state of the double pendulum approaches the vicinity of the desired goal state, the transition from the SAC controller to the LQR controller is made for final-stage stabilization. An essential aspect of

3 Methodology

any combined control strategy is the determination of the conditions under which the system is ready for a transition between control methods. In our SAC+LQR control strategy, the criterion for making this determination is provided by the region of attraction method[29].

The region of attraction (ROA) for a nonlinear system, represented by R_a , signifies the set \mathcal{B} of initial states surrounding a fixed point x_0 . If a state lies within this region, the system will gravitate towards x_0 as $t \rightarrow \infty$. In the context of an LQR controller, the ROA demarcates the area within the state space where the controlled system exhibits asymptotic stability. For complex systems, directly computing R_a can be challenging; it is often estimated instead. Therefore, analyzing the stability of the controlled system within this region becomes crucial. The state space representation of a controlled linear system can be expressed as:

$$\dot{x}(t) = (A - BK)x(t) \quad (3.16)$$

where $A_c = A - BK$.

The most straightforward way to analyze the stability of a system is using the Lyapunov method. We choose a Lyapunov function in a quadratic form:

$$V(\bar{x}) = \bar{x}^T S_{LQR} \bar{x} \quad (3.17)$$

where S_{LQR} is a positive definite matrix. This function serves as an "energy-like" metric, and our goal is to demonstrate its decrease over time for a stable system.

Next, we compute the time derivative of the Lyapunov function. For the infinite horizon LQR, $\frac{\partial S_{LQR}}{\partial t} = 0$ and $\frac{\partial x_0}{\partial t} = 0$, hence \dot{V} is:

$$\dot{V}(\bar{x}) = 2\bar{x}^T S_{LQR} \dot{\bar{x}} \quad (3.18)$$

For the system to be asymptotically stable at the equilibrium x_{op} , we require the system to satisfy the Lyapunov conditions:

$$\begin{cases} V(\bar{x}) > 0 \\ \dot{V}(\bar{x}) < 0 \quad \text{for all } x \neq x_{op} \end{cases} \quad (3.19)$$

To construct an optimizable region of attraction, we define an upper bound ρ for the Lyapunov function V . We seek the greatest ρ such that the above Lyapunov conditions are satisfied:

$$\begin{cases} \mathcal{B} = \{x \mid 0 < V(x) < \rho\} \\ \dot{V}(x) = 2\bar{x}^T S_{LQR} \dot{\bar{x}} < 0 \end{cases} \quad (3.20)$$

We computed the RoA similar to[29] but with a sums of squares method[39]. Consequently, the resulting shape of the ROA in a 4D state space resembles an ellipsoid.

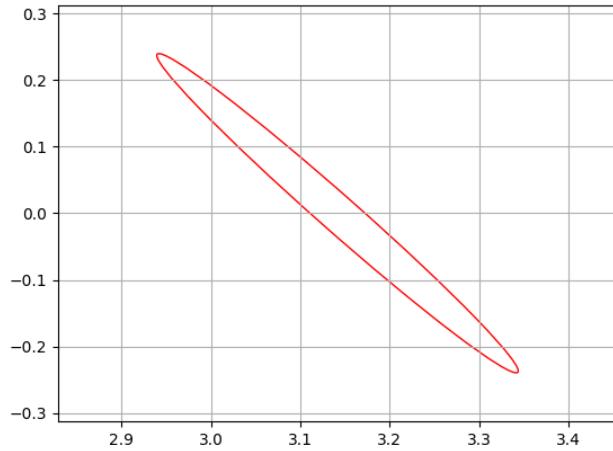


Figure 3.1: Region of Attraction

Once the RoA is computed, it can be checked whether a state x belongs to the estimated RoA of the LQR controller by calculating the cost-to-go of the LQR controller with the matrix S_{LQR} and comparing it with the scalar ρ .

3.4 Reward shaping

In theory, the reward function is designed to guide the agent's behavior towards achieving stability around the system's goal point. However, initial training attempts have revealed

3 Methodology

challenges. For instance, the agent might get stuck in a local minimum for an extended period, fail to perform a swing-up, or be unable to maintain stability around the upright position for a prolonged duration. The stabilization problem is addressed by the LQR controller and its region of attraction, which guarantees asymptotic stability. To tackle the swing-up issue, we designed a customized three-stage reward function to steer the agent away from problematic local minima and into the region of attraction of the LQR controller. The full equation for this reward function is:

$$\begin{aligned}
 r(x, u) = & - (x - x_g)^T Q_{train} (x - x_g) - u^T R_{train} u \\
 & + \begin{cases} r_{line} & \text{if } h(p_1, p_2) \geq h_{line} , \\ 0 & \text{else} \end{cases} \\
 & + \begin{cases} r_{LQR} & \text{if } (x - x_g)^T S_{LQR} (x - x_g) \geq \rho , \\ 0 & \text{else} \end{cases} \\
 & - \begin{cases} r_{vel} & \text{if } |v_1| \geq v_{thresh} , \\ 0 & \text{else} \end{cases} \\
 & - \begin{cases} r_{vel} & \text{if } |v_2| \geq v_{thresh} , \\ 0 & \text{else} \end{cases}
 \end{aligned} \tag{3.21}$$

In the initial stage, a quadratic reward function is employed to encourage smooth swinging of the entire system within a relatively small number of training sessions. The matrix $Q_{train} = diag(Q_1, Q_2, Q_3, Q_4)$ is a diagonal matrix, while R_{train} is a scalar. This is due to the nature of underactuated control in the double pendulum system, where only a single control input is available.

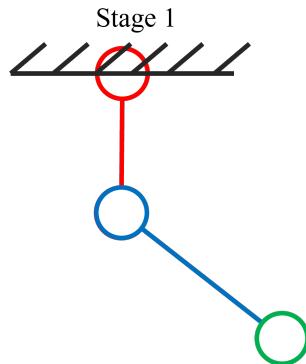


Figure 3.2: Swing up stage 1

As the end effector reaches a threshold line $h_{line} = 0.8(l_1 + l_2)$, we introduce a second level of reward r_{line} . The end effector height is given by

$$h(p_1, p_2) = -l_1 \cos(p_1) - l_2 \cos(p_1 + p_2). \quad (3.22)$$

with the link lengths l_1 and l_2 . This reward provides the agent with a fixed value but is carefully designed to prevent the system from spinning rapidly in either clockwise or counterclockwise directions. To discourage the agent from exploiting rewards by spinning at excessive speeds, a significant penalty $-r_{vel}$ is implemented for any speed exceeding $v_{thresh} = 8 \text{ rad/s}$ in absolute value. This penalty effectively compels the agent to approach the maximum point while adhering to the predefined speed interval. The speed penalty was only needed for the acrobot.

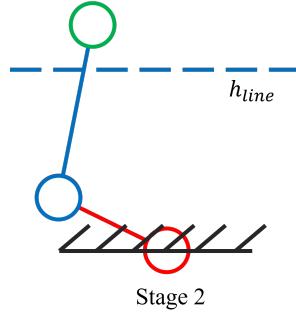
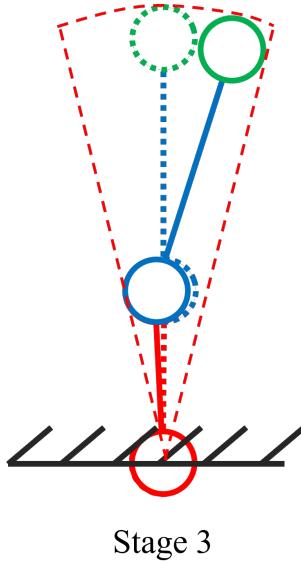


Figure 3.3: Swing up stage 2

The third level of reward r_{LQR} aims to provide a substantial reward to the agent when it remains within the Region of Attraction (RoA) of the LQR controller. By this we want to achieve that the policy learns to enter the LQR controller RoA so that there can be a smooth transition between both controllers. For details on the LQR controller and its region of attraction, we refer to these lecture notes. The parameters, we used in the cost matrices of the LQR controller are listed in Table 5.1. We computed the RoA similar to but with a sums of squares method. Once the RoA is computed, it can be checked whether a state x belongs to the estimated RoA of the LQR controller by calculating the cost-to-go of the LQR controller with the matrix S_{LQR} and comparing it with the scalar ρ .



Stage 3

Figure 3.4: Swing up stage 3

3.5 Introduction to leaderboard metrics

To compare the performances of various controllers developed for the double pendulum test-bench, we have established three distinct leaderboards: the simulation leaderboard, the robustness leaderboard, and the real system leaderboard. Each of these leaderboards is further subdivided into two categories based on the pendubot and acrobot setups.

3.5.1 Performance Leaderboard in Simulation and Real system

In our evaluation of controller performance, we employ a set of metrics that go beyond just measuring task success. They delve deeper into the nuances of controller operation. For experiments conducted both in simulation environments and on real hardware, the performance evaluation metrics remain consistent. These measures span from assessing the fundamental success of a swing-up maneuver to understanding the intricate details of energy and torque usage. Below is a detailed breakdown of each metric:

- **Swingup Success** c_{success} : Determines if the end-effector successfully remains above the predefined threshold by the simulation's conclusion.
- **Swingup Time** c_{time} : Measures the duration taken for the pendubot or acrobot to achieve and maintain its position above the threshold line. The metric only considers the swingup successful if the end-effector remains above the threshold until the simulation's end.
- **Energy** c_{energy} : Quantifies the total mechanical energy expended during the task.
- **Max Torque** $c_{\tau,\text{max}}$: Captures the highest torque applied at any point during the task.
- **Integrated Torque** $c_{\tau,\text{integ}}$: Represents the cumulative torque applied throughout the task's duration.
- **Torque Cost** $c_{\tau,\text{cost}}$: A quadratic metric that weighs the torques used, defined as $c_{\tau,\text{cost}} = \sum u^T R u$, where $R = 1$.
- **Torque Smoothness** $c_{\tau,\text{smooth}}$: Reflects the variability or fluctuations in the torque signals by measuring their standard deviation.
- **Velocity Cost** $c_{\text{vel},\text{cost}}$: A metric assessing the joint velocities achieved, computed as $c_{\text{vel}} = \dot{q}^T Q \dot{q}$, with Q being the identity matrix.

We utilize the subsequent criteria to determine the cumulative RealAI Score based on the specified formula:

$$S = c_{\text{success}} \left(w_{\text{time}} \frac{c_{\text{time}}}{n_{\text{time}}} + w_{\text{energy}} \frac{c_{\text{energy}}}{n_{\text{energy}}} + w_{\tau,\text{max}} \frac{c_{\tau,\text{max}}}{n_{\tau,\text{max}}} + w_{\tau,\text{integ}} \frac{c_{\tau,\text{integ}}}{n_{\tau,\text{integ}}} + w_{\tau,\text{cost}} \frac{c_{\tau,\text{cost}}}{n_{\tau,\text{cost}}} + w_{\tau,\text{smooth}} \frac{c_{\tau,\text{smooth}}}{n_{\tau,\text{smooth}}} + w_{\text{vel},\text{cost}} \frac{c_{\text{vel},\text{cost}}}{n_{\text{vel},\text{cost}}} \right) \quad (3.23)$$

The weights and normalizations are:

3 Methodology

Criteria	Normalization n	Weight w
Swingup time	10.0	0.2
Energy	100.0	0.1
Max. Torque	6.0	0.1
Integrated Torque	60.0	0.1
Torque Cost	360	0.1
Torque Smoothness	12.0	0.2
Velocity Cost	1000.0	0.2

Table 3.1: Weights and normalizations for performance leaderboards

In the simulation experiments, the pendubot is modeled using a Runge-Kutta 4 integrator with a timestep of $dt = 0.002s$ over a span of $T = 10s$. We initiate the pendubot in a hanging down configuration, represented as $x_0 = (0, 0, 0, 0)$, and aim to reach the unstable fixed point of the upright configuration, denoted as $x_g = (\pi, 0, 0, 0)$. The double pendulum is deemed to have achieved its upright position once the end-effector surpasses the threshold line situated at $h = 0.45m$, with the origin being the mounting point.

When it comes to real hardware experiments, there's a torque limit of 0.5 Nm on the passive joint, which serves to offset the motor's friction. The actuators can operate with a control frequency as high as 500Hz, and each experiment lasts for 10 seconds. The pendubot starts from a hanging down position, with the objective being the unstable fixed point in the upright configuration. Successful attainment of the upright position is confirmed when the end-effector crosses the threshold line set at $h = 0.45m$, measured from the mounting point's origin.

3.5.2 Simulation Robustness Leaderboard

In addition to performance metrics, we also consider robustness metrics. As the ultimate goal is to transfer successful models from a simulation environment to real hardware, it's essential to assess the robustness of controllers developed within the simulation. This helps determine the types of perturbations that affect each controller.

- **Model Inaccuracies c_{model} :** Model parameters determined through system identification are never perfectly accurate. To evaluate potential inaccuracies, we modify each model parameter individually within the simulator, while keeping the controller's parameters unchanged.

- **Velocity Measurement Noise** $c_{vel,noise}$: Controller outputs heavily rely on the captured system state. Particularly with the Quasi Direct Drives (QDDs), online velocity measurements can be noisy. It's crucial for successful transferability that a controller can handle this inherent noise. We test controllers with both the presence and absence of a low-pass noise filter.
- **Torque Noise** $c_{\tau,noise}$: It's not only the measurements that can exhibit noise. Sometimes the torque output from the controller might deviate from the intended value.
- **Torque Response** $c_{\tau,response}$: The torque requested by the controller is dynamic and can vary throughout execution. Due to mechanical limitations, the motor might not always adjust immediately to abrupt torque changes, leading to overshooting or undershooting the desired torque value. We model this with the equation $\tau = \tau_{t-1} + k_{resp}(\tau_{des} - \tau_{t-1})$, where τ_{des} is the desired torque. In this model, a k_{resp} value of 1 indicates flawless torque response, while any deviation from 1 indicates imperfect motor responses.
- **Time Delay** c_{delay} : Operating in a real-world environment inevitably introduces time delays due to communication lag and system reaction times. It's essential to account for these when evaluating controller performance.

The above criteria are employed to compute the comprehensive RealAI Score using the given formula:

$$S = w_{model}c_{model} + w_{vel,noise}c_{vel,noise} + w_{\tau,noise}c_{\tau,noise} + w_{\tau,response}c_{\tau,response} + w_{delay}c_{delay} \quad (3.24)$$

The weights are:

$$w_{model} = w_{vel,noise} = w_{\tau,noise} = w_{\tau,response} = w_{delay} = 0.2 \quad (3.25)$$

4 Agent training and experiments in simulation environment

Chapters 4 and 5 focus on the experimental phase of our project. This chapter details the training procedure of our SAC agent for the swing-up task and then transitions to a simulation phase to validate results for both the swing-up and stabilization tasks. We have structured this chapter into three distinct subsections: training setup, training process, and simulation results.

The first subsection describes the foundation of a reinforcement learning interaction rooted in a stable baseline3-based RL algorithm. It also touches upon a customized environment inherited from the OpenAI Gym environment. The second subsection centers on hyperparameter tuning and highlights the challenges encountered during training. The third subsection displays the results acquired from both the pendubot and acrobot setups.

4.1 Training setup

Stable Baseline 3(SB3)[31] is an open-source implementation of deep reinforcement learning algorithms based on the Python language. This library includes seven commonly used model-free deep reinforcement learning algorithms including SAC. Prior implementations of deep RL algorithms often encountered a problem wherein small implementation details could greatly affect performance, typically exceeding the differences between algorithms[some paper]. The developers of SB3 have done a commendable job stabilizing the performance of these deep RL algorithms by benchmarking each one on common environments and comparing them to prior implementations. Owing to its user-friendly and reliable nature, we chose Stable Baseline 3 as our RL library, which greatly simplified our research process.

OpenAI Gymnasium[40] is a toolkit for developing and comparing reinforcement learning algorithms, and it is widely used in the research community. Among its many advantages are its open-source nature and standardized environments, which facilitate the rapid testing and benchmarking of new algorithms. Additionally, it offers easy visualization and monitoring. Our decision to use the Gym library is based on its extensibility—from standard environments to highly customized ones—as well as its capability to integrate with tools like PyTorch and TensorFlow for GPU-accelerated computations. Furthermore, Gym provides the ability to construct stacked training environments for parallel training.

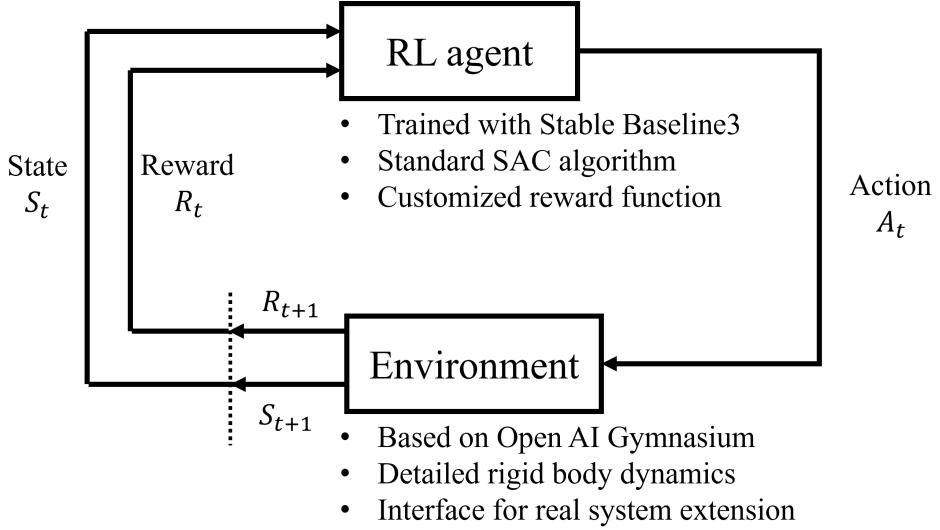


Figure 4.1: Interaction between reinforcement learning agent and customized environment

To construct the customized environment, we begin with the symbolic dynamic function that captures the nonlinear dynamics of the underactuated double pendulum. This function uses the current state observation from the simulation and the action determined by the control policy, producing a prediction of the next state via a Runge-Kutta integrator.

This predicted state is then input into the reward function, which yields a scalar output. This output is subsequently relayed to the SAC algorithm for policy evaluation and update. After these computations, the simulation provides an estimation of the current state, and the policy identifies the most probable action. Both of these values are then directed back to the dynamics function, initiating a new cycle.

It's widely recognized that reinforcement learning algorithms tend to converge more effectively when using normalized state and action spaces [37]. Therefore, we've designed a scaling mechanism to translate the state and action spaces from their normalized versions to real-world measurements. Users can choose to activate or deactivate this scaling. For example, while a normalized state resides within the interval $[-1, 1]$, we might wish to map it to real-world measurements $[pos1, pos2, vel1, vel2]$, such as $pos1$ in $[0, 2\pi]$, $pos2$ in $[-\pi, \pi]$, and velocities $vel1$ and $vel2$ in $[-20, 20]$ rad/s, while keeping the torque within $[-\tau_{\text{limit}}, \tau_{\text{limit}}]$. When this scaling mechanism is activated, it ensures that states and actions in the SAC algorithm always stay within the $[-1, 1]$ boundary, which often leads to faster convergence.

The logic of the pipeline of the interaction of customized environment is described in the picture below.

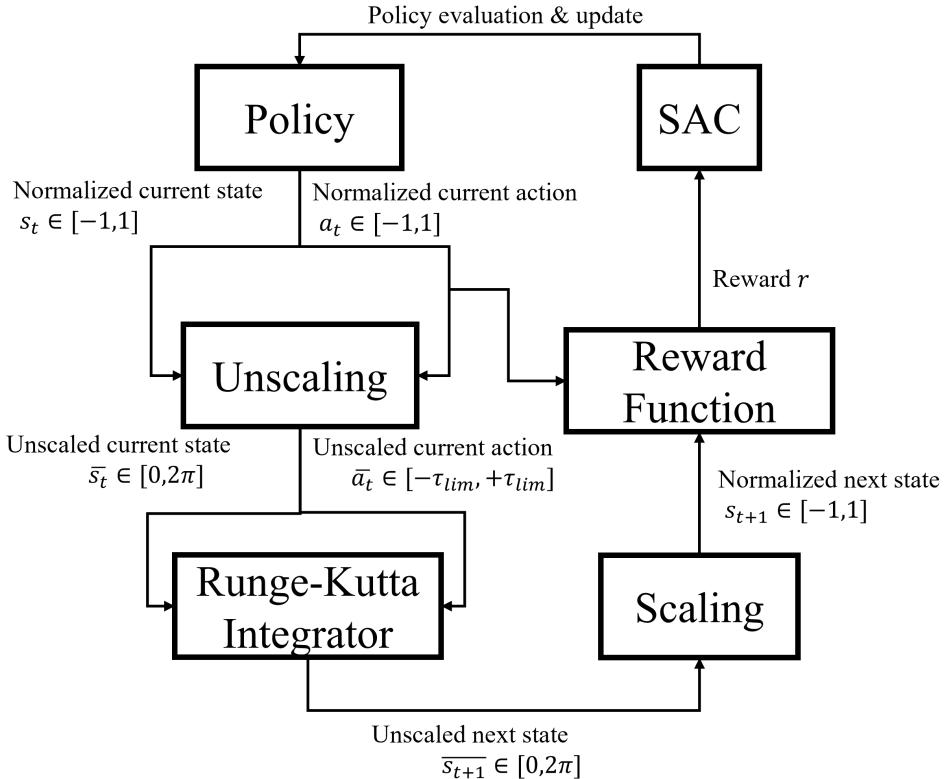


Figure 4.2: Detailed scaling pipeline in training process

4.2 Training phase

The training phase started immediately after we finished setting up the training pipeline. During the training of the SAC controller for both the acrobot and pendubot, our primary focus was on tuning several key hyperparameters. These encompassed the learning rate, control frequency, episode length, and learning timestep. We invested much effort in adjusting hyperparameters for the reward function and the LQR controller, given their pivotal roles in the learning process, as detailed in Table 4.1.

4 Agent training and experiments in simulation environment

We set the learning rate at 0.01 to promote effective learning and adaptation. Additionally, we configured the control frequency of the simulation to 100Hz, ensuring frequent updates and heightened responsiveness in the control process. We selected an episode length of 1000 for both the Acrobot and Pendubot, translating to 10-second-long episodes, to provide sufficient exploration and learning opportunities. To harness the full training potential, we executed a total of 2e7 learning time steps for the pendubot and 5e7 for the acrobot, allowing the agent to accumulate vast experience and enhance its performance.

Robot	Quadratic Reward	Constant Reward	LQR
Pendubot	$Q_1 = 8.0$		$Q_1 = 1.92$
	$Q_2 = 5.0$	$r_{line} = 500$	$Q_2 = 1.92$
	$Q_3 = 0.1$	$r_{vel} = 0.0$	$Q_3 = 0.3$
	$Q_4 = 0.1$	$r_{LQR} = 1e4$	$Q_4 = 0.3$
	$R = 1e-4$		$R = 0.82$
Acrobot	$Q_1 = 10.0$		$Q_1 = 0.97$
	$Q_2 = 10.0$	$r_{line} = 500$	$Q_2 = 0.93$
	$Q_3 = 0.2$	$r_{vel} = 1e4$	$Q_3 = 0.39$
	$Q_4 = 0.2$	$r_{LQR} = 1e4$	$Q_4 = 0.26$
	$R = 1e-4$		$R = 0.11$

Table 4.1: Hyper parameters used for the SAC training and the LQR controller.

The learning curve for a total of 2e7 timesteps for the pendubot is illustrated below. As depicted in the figure, from 0 to 1e7 timesteps, the learning curve experiences a gradual ascent. Between 1e7 and 1.2e7 timesteps, the learning curve surges rapidly, peaking in reward around one million. After 1.2e6 timesteps, the curve begins to stabilize with occasional fluctuations, and no substantial increase in reward is observed.

4.2 Training phase

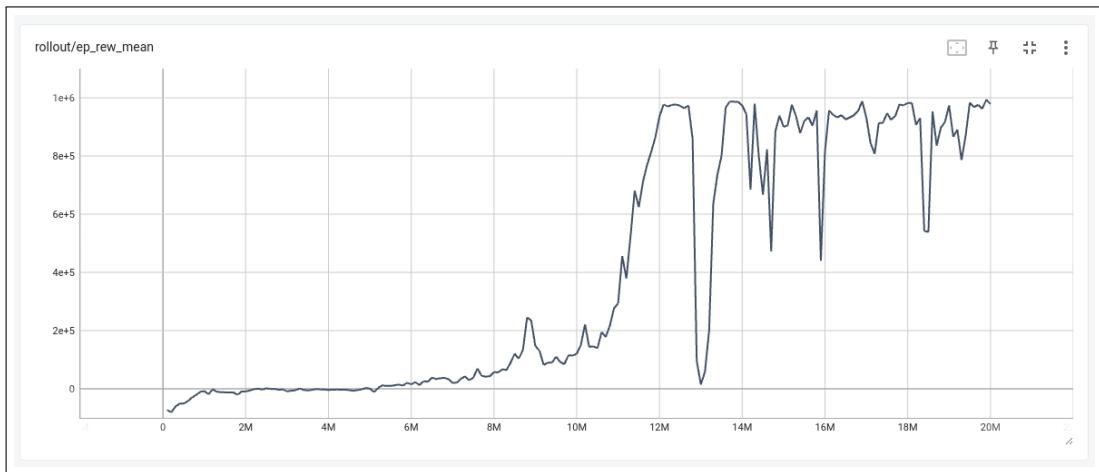


Figure 4.3: Pendubot learning curve

As shown in the learning curve over a total of $3e7$ timesteps for the acrobot, the curve experienced a relatively steady growth until $1.8e7$ timesteps. After a sudden drop in reward between $1.8e7$ and $2e7$ timesteps, the learning curve began to increase drastically, with two major fallbacks. By $3e7$ timesteps, the learning curve hadn't reached its peak. We extended the learning period to $5e7$ using a warm start from the model obtained at $3e7$ timesteps, and the learning curve began to stabilize around $4e7$ time steps.

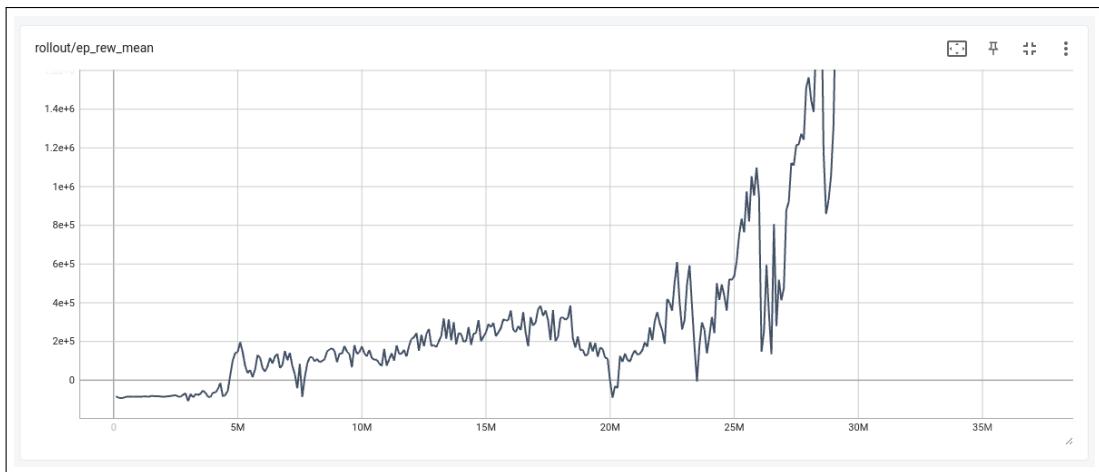


Figure 4.4: Acrobot learning curve

In conclusion, the training for the pendubot converges more quickly than that for the acrobot. Additionally, based on the training curves, the pendubot's training process is more stable than the acrobot's.

4.3 Simulation results

In this section, the simulation results for the acrobot and pendubot are presented separately. The model trained using the SAC algorithm is utilized during the swing-up stage and switches to the LQR controller when nearing the upright position. The primary success criteria for both swing-up and stabilization involve swinging up the double pendulum and maintaining its stability around the upright position for a prolonged period. Therefore, a total simulation time of 10 seconds is used. If the double pendulum fails to swing up within these 10 seconds, the result is deemed unsuccessful. Similarly, if the double pendulum loses stability within this time frame, the outcome is still considered a failure.

4.3.1 Pendubot simulation in ideal environment

A testing environment, identical to the customized learning environment, has been established to validate the training results and accurately replicate the agent's behavior learned during the reinforcement learning process. The pendulum is initialized at its lowest point with zero velocity and begins its swing-up using the control policy derived solely from the SAC.

As depicted in the figure, the swing-up time is under 1 second. After this, the state of the pendubot enters the Region of Attraction (ROA) of the LQR controller. The transition between the two controllers is both seamless and effective, with the system stabilizing towards the desired state under the LQR controller's influence. This validates the effectiveness of the ROA method for the LQR takeover.

A significant highlight from this successful simulation is the impressively short swing-up time, despite having strict torque limit in place. However, a concerning feature observed from this simulation is the noisiness of the input control signal. The torque alternates signs rapidly, and the gradient of the torque tends to have a high absolute value.

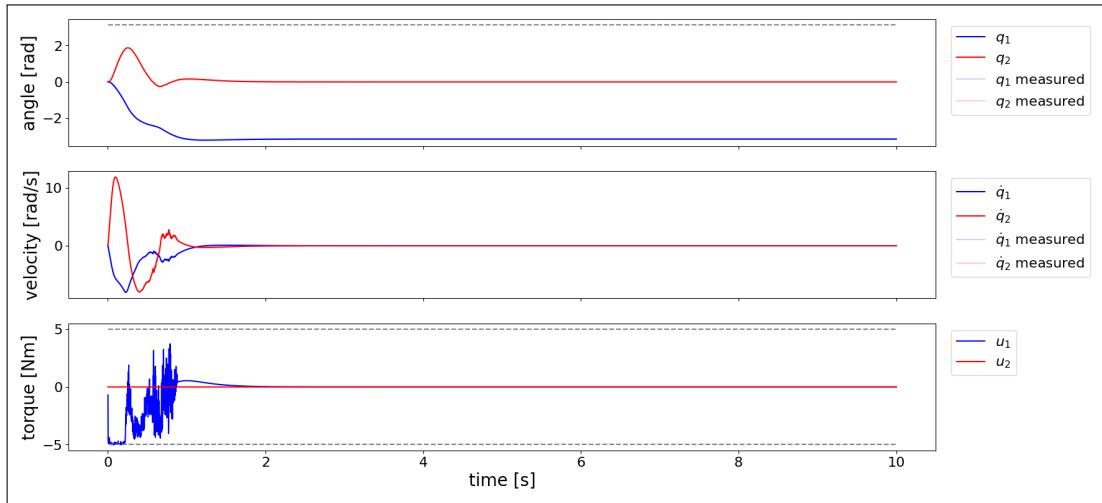


Figure 4.5: Pendubot simulation result

4.3.2 Acrobot simulation in ideal environment

Just like the pendubot simulation, the acrobot setup is tested in an ideal environment identical to its training environment. The acrobot begins its swing from the downward position with the goal of stabilizing around its highest point.

The accompanying image depicts a successful swing-up and stabilization within a 10-second window. The swing-up process takes about 2 seconds before the LQR effectively assumes control, maintaining an asymptotic stability around the target state with minimal fluctuations.

In comparison to the pendubot's simulation curves, the swing-up time for the acrobot is roughly twice as long. This aligns with the notion that controlling the acrobot is a more challenging task than the pendubot. A shared drawback observed in both simulations is the lack of torque smoothness. For the acrobot, the input torque exhibits several significant jumps from one torque limit to the other, which might cause difficulties when translating to real-world hardware control.

4 Agent training and experiments in simulation environment

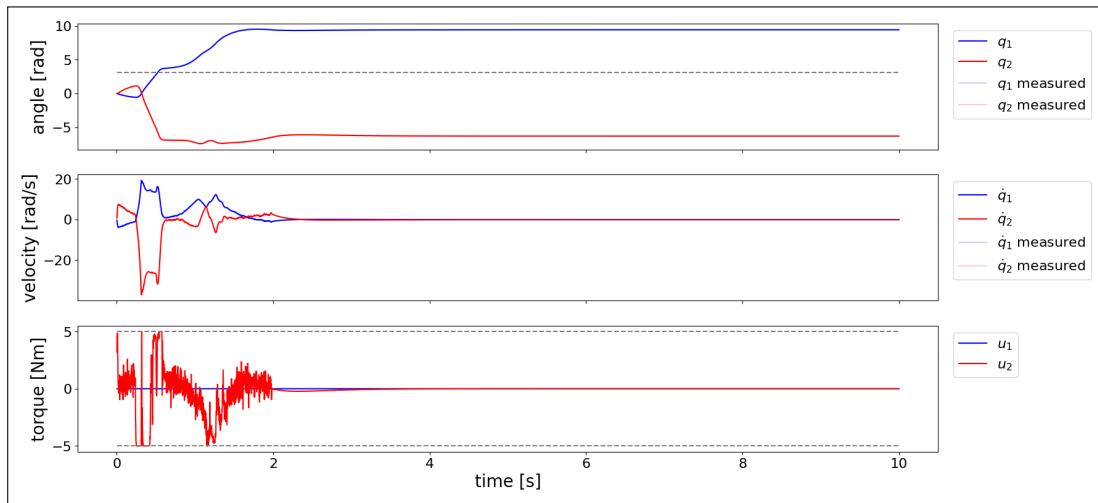


Figure 4.6: Acrobot simulation result

5 Experiments on real hardware

In this chapter, we discuss experiments conducted on the hardware system. The content is organized into four sections. The first section provides an overview of the hardware setup for the double pendulum system. The second section delves into the system identification of the hardware. The third section details our approach to addressing the sim-to-real gap challenge. The final section presents the successful outcomes of our hardware experiments.

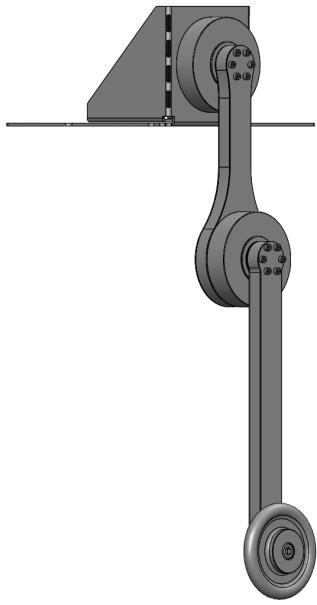
5.1 Hardware setup

Mechanically, the double pendulum system is a straightforward 2-R linkage. The first revolute joint attaches to the base, while the second one connects the two links. Quasi-direct drive motors are mounted on each joint to provide torque, and a counterweight is positioned at the end of the second link.

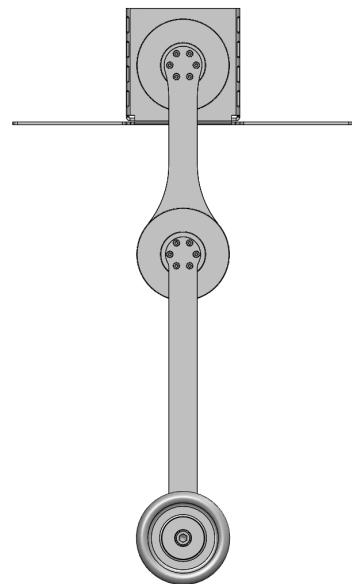
Our mechanical design for the double pendulum underwent two iterations. In the initial design, the base consisted of a bent aluminum plate, and the links featured a sandwich structure with aluminum on the outside and engineering plastic on the inside. These links were of homogeneous size, meaning that the cross-sectional areas at the link intersections were consistent. We abandoned this design due to rotational imbalances. Specifically, the rotation plane of the links wasn't consistently perpendicular to the motor axes, leading to vibrations and accelerated wear during regular use. In extreme cases, when the links rotated at high angular velocities, this imbalance was exacerbated. This not only resulted in the links bending but also led to catastrophic system failures. Such failures presented significant safety risks to personnel.

In the second iteration, we addressed the issues encountered in the previous design, leading to two major modifications. First, we replaced the aluminum-plastic combination with a carbon fiber-foam blend. While the incorporation of carbon fiber marginally increased the cost, it significantly boosted the yield strength. Second, we introduced triangular-shaped links with central cutouts. This design, while lightweight, also substantially enhanced the yield strength due to the changed intersections. Overall, this iteration resulted in a mechanical structure that is considerably more reliable than the first.

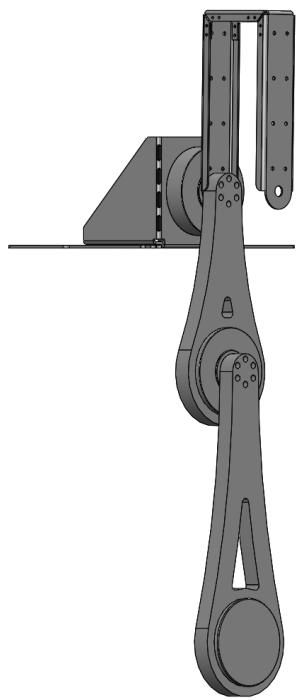
5 Experiments on real hardware



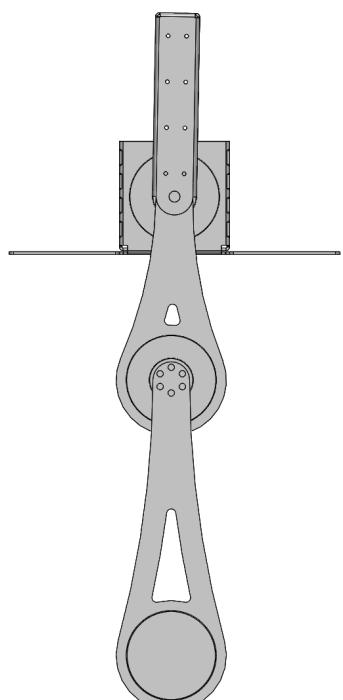
(a) Iteration 1 in isometric view



(b) Iteration 1 in front view



(c) Iteration 2 in isometric view



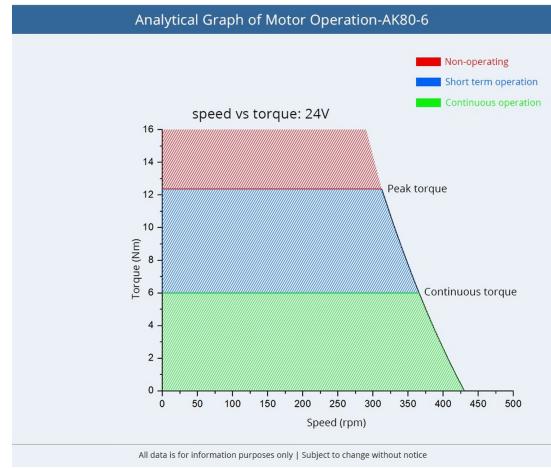
(d) Iteration 2 in front view

Figure 5.1: Double pendulum mechanical system iteration 1 and iteration 2

For the quasi-direct drive(QDD) motors, we selected the AK80-6 V100 motors from the company CubeMars. This motor's design facilitates easy mounting from both the front and rear ends. As depicted in Figure5.2(b), the motor's maximum torque during continuous operation is 6 Nm, which aligns well with our torque limit of 5 Nm. Additionally, the motor is designed for compatibility with both serial bus and CAN bus, simplifying the development process.



(a) AK80-6 V100 motor



(b) Speed-torque diagramm

Figure 5.2: Quasi direct drive motor AK80-6 V100

For communication, We have chosen CAN in this implementation. The Controller Area Network (CAN) bus is a robust, flexible, and efficient communication protocol that has been employed in various applications. There are many advantages to using the CAN bus for control. The CAN bus provides error checking and fault confinement capabilities. It operates in real-time, enabling high control frequencies with relatively simple wiring. In our configuration, there is only one master node (the PC) and two slave nodes (two motors). The control loop simply consists of a CAN-to-USB converter, one CAN high cable, and one CAN low cable, with a termination resistor of 120Ω at both ends. The detailed connection is depicted in the figure below:

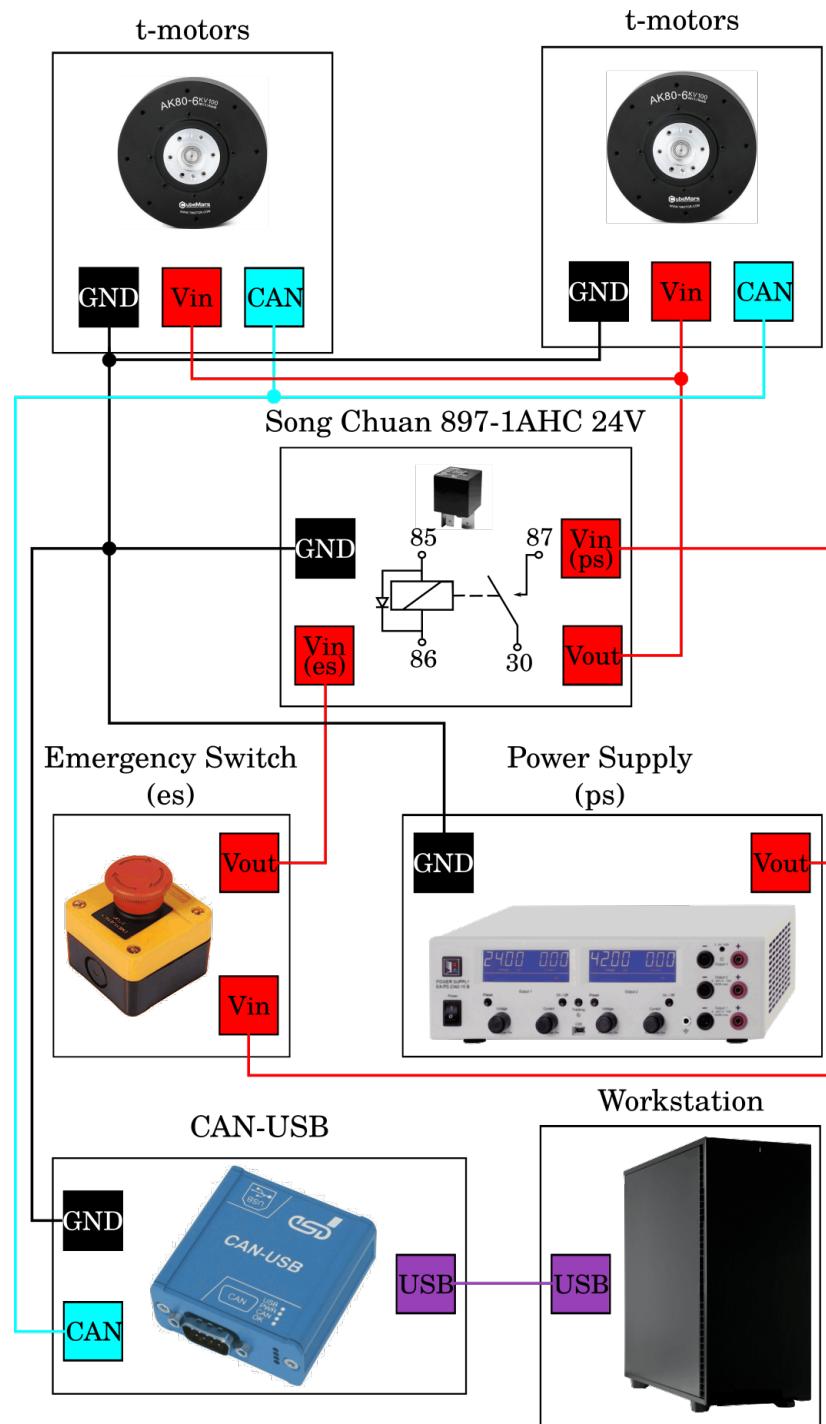


Figure 5.3: Wiring diagram

To ensure a higher control frequency of around 500 Hz, we chose the CAN-USB/2 product from ESD GmbH Hannover. This specialized CAN-to-USB interface utilizes USB 2.0, which supports a data rate of 480 Mbit/s. Its CAN capability is 1 Mbit/s, in accordance with ISO 11898-2. Additionally, it supports the SocketCAN interface included in the Linux Kernel 2.6, making it easier to use in a Linux development environment.



Figure 5.4: High speed CAN to USB interface[not the same version, correct it]

Due to accidents that occurred during testing with the mechanical systems from the first iteration, several safety protocols have been implemented to ensure the safety of both human lives and equipment. Four major measures have been taken.

An emergency stop is connected directly to the 24V power source. If the behavior of the double pendulum deviates from the expected range during testing, the power supply can be manually cut off immediately. The energy in the mechanical system will dissipate rapidly, and the system will return to its initial state automatically.

In scenarios where the system is moving at a very high speed when the emergency stop is engaged, the motors attached to the revolute joints act as generators, dissipating the mechanical energy from the system's motion. The generated current is fed back into the circuit, and in extreme cases, it could overload the power supply. To counteract this, a capacitor is connected to the power supply to absorb any electrical surge resulting from a sudden stop.

Further, speed and position limits have been set at the software level. As the links of the pendulum can experience vibrations and rotational imbalances at high speeds, potentially leading to structural disassembly, a speed limit of 20 rad/s has been established. Any speed exceeding this value will trigger a full system stop, equivalent to pressing the emergency stop button. The position limit is set to 2π for both pos1 and pos2. Excessive rotations could cause

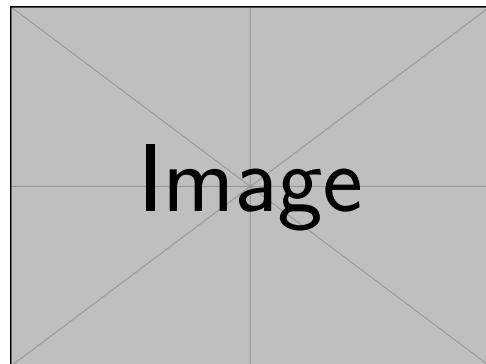


Figure 5.5: picture of the capacitor

the CAN and power cables to become entangled, leading to interference and potential cable damage.

Lastly, to guard against unpredictable system failures, a custom protective cage has been constructed using aluminum profiles and thick acrylic boards. This cage fully encloses the double pendulum hardware, significantly reducing the risk of accidents.

These measures have been instrumental in minimizing potential hazards during testing and operation.

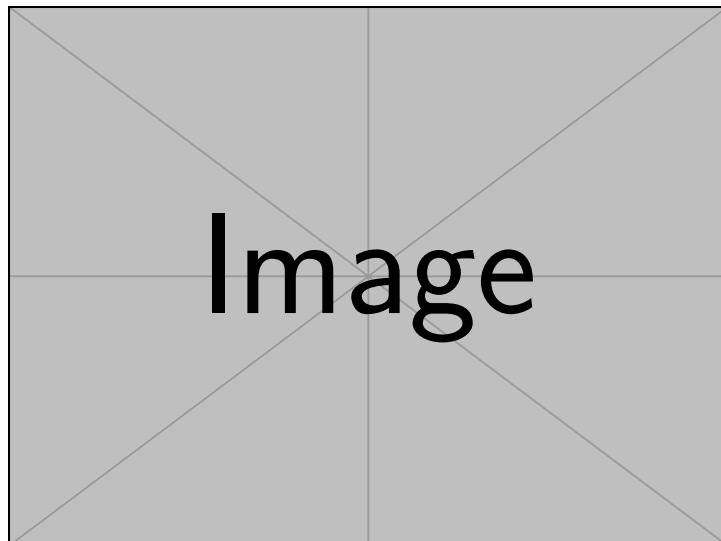


Figure 5.6: a overview of experiment setup

5.2 System identification

System identification is the process of deriving mathematical models of dynamic systems from observed input-output data. This method is fundamental in control theory, used to analyze, predict, and control the behavior of real-world systems. After setting up the hardware system and before transitioning the working model from a successful simulation to the real system, it's necessary to undergo a system identification process. This helps ascertain the real-world parameters governing the dynamics of the double pendulum.

Out of all model parameters, 15 are selected. While the naturally provided parameters g and g_r are held constant, the easily measurable parameters l_1 and l_2 are considered independent. The remaining system parameters, which are

$$m_1 r_1, m_2 r_2, m_1, m_2, I_1, I_2, I_r, b_1, b_2, c_{f1}, c_{f2}$$

need to be identified. The ultimate objective is to discern the parameters of the dynamic matrices present in the equations of motion.

$$M\ddot{q} + C(q, \dot{q})\dot{q} - G(q) + F(\dot{q}) - Du = 0 \quad (5.1)$$

By running excitation trajectories on the actual hardware, data tuples in the form $(q, \dot{q}, \ddot{q}, u)$ can be collected. To determine the most accurate system parameters, one can leverage the linearity of the dynamic matrices M , C , G , and F in relation to the independent model parameters. Consequently, a least squares optimization can be performed on the recorded data, relative to the dynamics equation.

The identified model parameters are shown in table below:

Parameter	Value
I_1	0.031887199591513114
I_2	0.05086984812807257
I_r	6.287203962819607e-05
b_1	0.001
b_2	0.001
c_{f1}	0.16
c_{f2}	0.12
g	9.81
g_r	6.0
l_1	0.2
l_2	0.3
m_1	0.5234602302310271
m_2	0.6255677234174437
r_1	0.2
r_2	0.25569305436052964

Table 5.1: Parameter values from system identification

5.3 Sim2Real problem

Transferring working models from simulation to real systems to produce similar performance has always been a challenge in controller design. This challenge is even more pronounced in model-free reinforcement learning for several reasons.

Firstly, model-free reinforcement learning relies solely on interaction with the environment to gain experience and select actions. While a simulation environment is merely a simplification of the real-world scenario, the agent in simulation might not capture all the factors, such as friction, sensor noise, or real-world dynamics, accurately. Therefore, a control policy optimized for a simplified model might not perform as expected in the more intricate real world.

Secondly, many simulations operate in discrete time and space, whereas the real world functions continuously. In our implementation, the control frequency presents a significant challenge. We use a control frequency of 100 Hz in simulation; however, it does not suffice in the real system. To enhance performance, we increased the control frequency to 400 Hz when experimenting on the real system, and this adjustment yielded positive results.

Thirdly, uncertainties in the environment can lead to substantial complications when attempting to apply learned strategies or actions. Generally speaking, the uncertainty of the environment is addressed by the robustness of the model. However, for highly chaotic systems like the pendubot or acrobot, even slight measurement errors can result in significant deviations from the planned behavior. Our model for controlling the pendubot and acrobot requires a higher level of robustness or a more effective sim2real transfer method.

5.3.1 Validation with noisy simulation

Our approach to addressing the sim2real challenge involves training multiple agents using the SAC algorithm under similar setups. Subsequently, we validate them in a noisy simulation. Only those agents that prove robust against perturbations in this noisy environment proceed to real system testing. Any agent failing these noisy simulation tests is deemed insufficiently robust and is discarded.

Throughout our experiments with real systems, we identified four critical factors contributing to successful swing-up and stabilization. Of these, friction emerged as the most significant. This is because our agent training was conducted in an ideal environment devoid of friction, thereby making friction the primary differentiator between the simulation and real-world conditions. To counter this, we adopted a friction compensation approach, beginning with modeling based on Coulomb's friction.

Coulomb's friction model is among the simplest existing friction models. It possesses two primary advantages: it's independent of both the contact area and the relative velocity. The frictional force (F_f) is directly proportional to the normal force (F_n) between the surfaces, represented by:

$$F_{fi} = c_{fi} \arctan(100\dot{q}_i), \quad i = 1, 2 \quad (5.2)$$

where c_f is the coefficient of kinetic friction. This frictional force opposes the relative motion between the surfaces. To neutralize this force, friction compensation applies torque in the same direction as the angular motion, energizing the system to mitigate friction's effects.

Though friction coefficients c_{f1} and c_{f2} were determined during the system identification phase, they seemed imprecise during real system tests. Consequently, we resorted to free fall tests to estimate the coefficients, fine-tuning them manually until the system behaved conservatively under friction's influence.

The second most critical factor is measurement noise. While the AK80-6 motors use built-in encoders to measure the position difference from the initial position with high accuracy, velocity measurement, which is derived from the first derivative of the position measurement, tends to have a relatively high error. In the simulation environment, we assume zero measurement error. However, in real-world applications, this error becomes significant.

To address this issue, we model the measurement error as a normal distribution where the mean corresponds to the true velocity value and the standard deviation is manually adjustable. Let's consider the measurement noise vector $\varepsilon = [\delta\text{pos}_1, \delta\text{pos}_2, \delta\text{vel}_1, \delta\text{vel}_2]^T$. We model this noise as following a multivariate normal distribution, given by:

$$\varepsilon \sim \mathcal{N}(\mu, \Sigma) \quad (5.3)$$

where μ is the mean vector (which represents the true values in the absence of noise), and Σ is the 4×4 covariance matrix that captures the spread or uncertainty of the noise in each dimension. Since the measurements of the four states are considered independent of each other, Σ is diagonal.

Latency is the third significant factor. Since any communication system requires time to transmit and receive data, and programs also take time to execute, latency is inevitable. This latency poses a substantial risk to control systems, especially those based on reinforcement learning. This is because reinforcement learning is grounded in the principles of Markov Decision Processes (MDPs) which adhere to the Markov property. This property dictates that the future state of a process is contingent only on the current state and action, not on the sequence of states that led to it. Latency undermines the Markov property by causing state mismatches and creating dependencies on historical data. If not addressed, this can have severe consequences.

The fourth factor to consider is torque responsiveness. We observed that when confronted with rapidly alternating control signals with significant differences between time steps, the motor struggles to produce torque that matches the control signal. In other words, the motor cannot respond timely to changes in torque. This lag means that the controller cannot operate at its full potential due to hardware constraints. Some potential solutions to this include reducing the control frequency and increasing torque smoothness.

In summary, we take into account friction, measurement error, latency, and torque responsiveness in the noisy simulation to select the most suitable agent.

As shown in the figure below, a trained pendubot agent successfully executed a swing-up and stabilization in the noisy simulation environment, indicating its robustness is sufficient for real-world application.

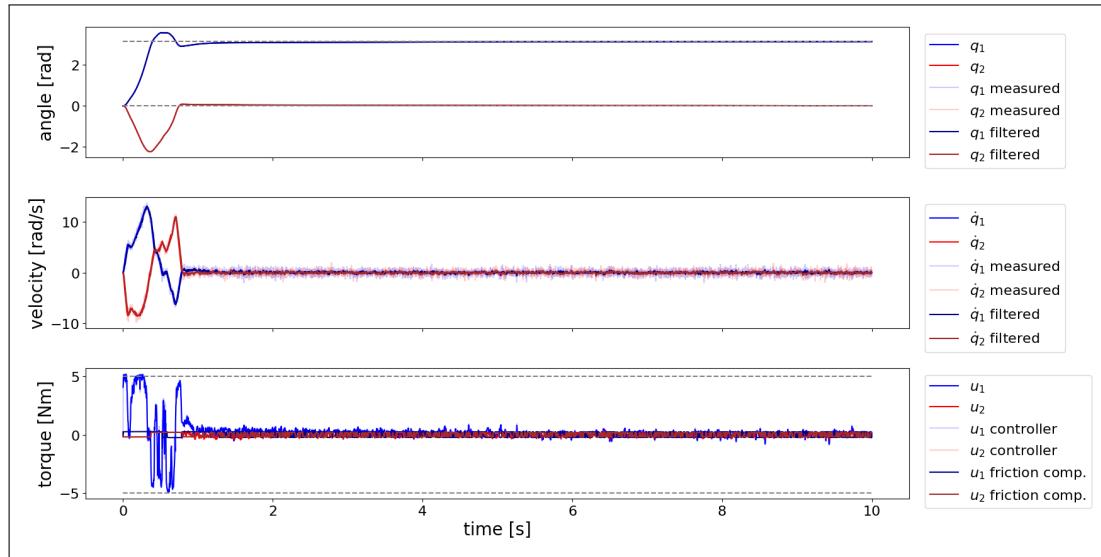


Figure 5.7: pendubot noisy simulation result

5.3.2 Training with domain randomization

Domain randomization[[some paper](#)] is a classic technique used to bridge the gap between simulation and reality. Its primary goal is to allow a model, when generalized in a simulated environment, to perform effectively in real-world scenarios without needing labeled real-world training data.

The concept of domain randomization originates from the field of computer vision. Instead of training a model in a single, fixed simulation, the model is exposed to a more varied and noisy environment by introducing random variations to the ideal simulation. Techniques commonly employed in vision-based systems include changing the colors and textures of objects, modifying the lighting conditions, adjusting object shapes and sizes, introducing random noise to sensor data, and perturbing physical properties like friction or mass.

In our effort to apply domain randomization to robotic manipulation tasks, we opted to introduce random measurement noise to both position and velocity and to introduce perturbations to model parameters. These perturbations are considered independent of one another and follow a normal distribution. The mean of this distribution is the true value, and there's a tunable standard deviation.

[results and more explanation]

5.4 Real hardware results

In this section, we present the results from the real hardware tests. Only agents trained for the pendubot setup successfully passed the noisy simulation check, so results are limited to this setup.

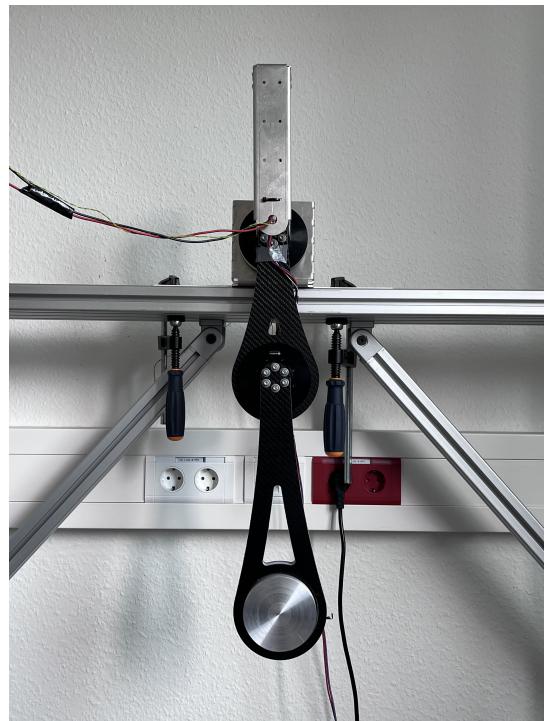


Figure 5.8: Double pendulum in real system

The complete real-world test procedure is as follows: power up, manually set the initial state, release the emergency button, run the initialization script (which enables the motor, sets the current position to zero, and tests the CAN connection), confirm the start of the test, record video, draw plots, and exit the test. This procedure is illustrated in the flow chart below:

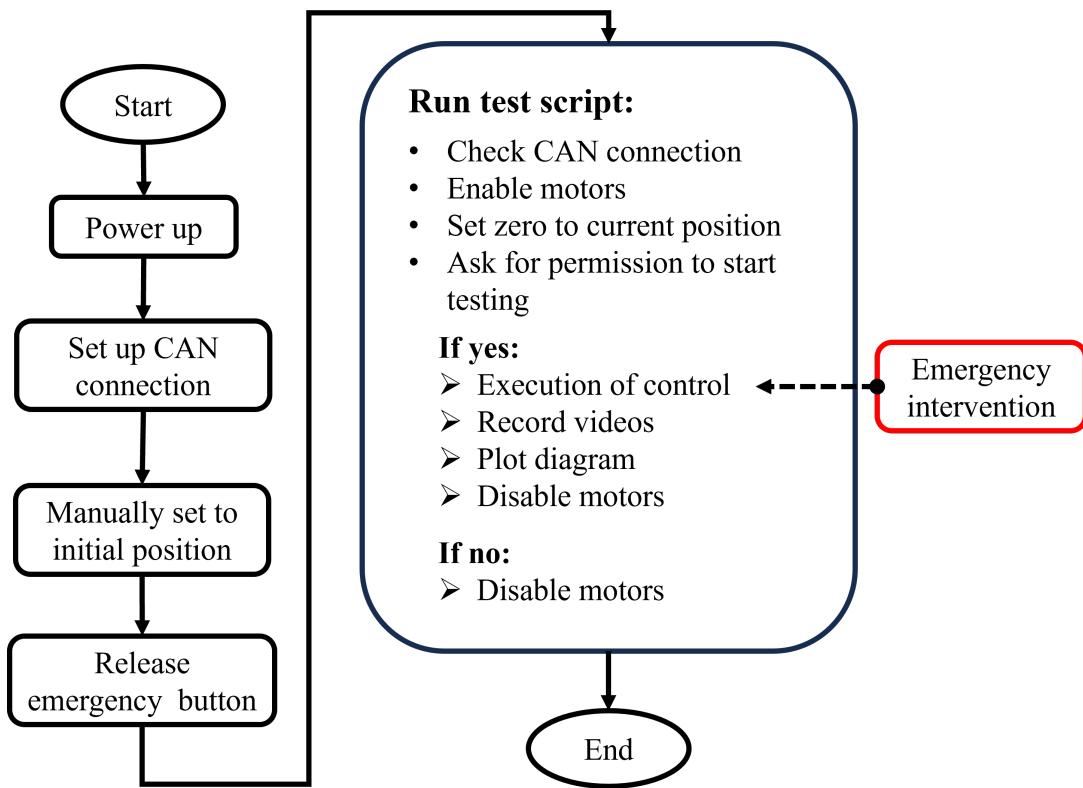


Figure 5.9: Real hardware system testing procedure

A remote testing system has been established to allow global access and sharing of our double pendulum testing hardware. This infrastructure was developed for the IJCAI 2023 competition, RealAIGym. Several groups have tested their algorithms on this remote system, including teams from the University of Padova and the Technical University of Darmstadt.

5.4.1 Pendubot results

Below is an example of one of our successful tests. [some explanation] We executed a sequence of 10 tests to determine the success rate and other key metrics for evaluating performance and robustness. These results will be discussed in the next chapter.

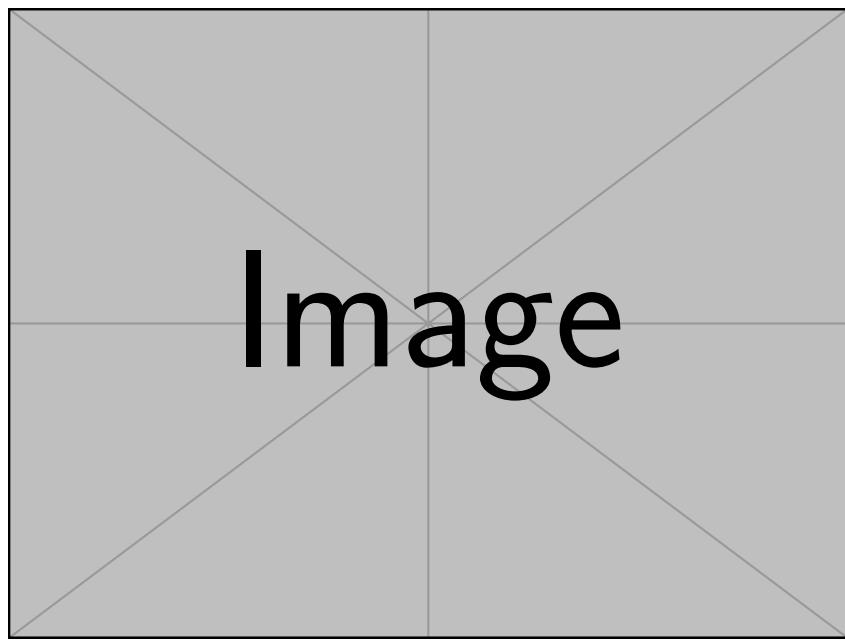


Figure 5.10: Pendubot results on real systems

5.4.2 acrobot results

acrobot:

6 Discussion and Future Work

In this chapter, we delve into the results of experiments conducted both in simulation and on the real system. The chapter is structured as follows: The first subsection provides a comprehensive introduction to the leaderboard metrics employed. Sections 6.2, 6.3, and 6.4 discuss the three facets of leaderboard analysis, namely simulation, robustness, and real hardware, respectively. The final section wraps up our current findings and hints at directions for future research.

6.1 Interpretation of simulation leaderboard

In the table below, the performance leaderboard results for both the pendubot and acrobot in simulation experiments are presented. Three major types of controllers are listed for comparison. The SAC+LQR controller is our design and is based on the model-free reinforcement learning method. MC-PILCO [4] [25], which stands for Monte Carlo Probabilistic Inference for Learning Control, is a model-based reinforcement learning method. It utilizes probabilistic models to predict the system's dynamics and employs Monte Carlo methods to optimize control policies based on these predictions. This method was implemented by a team from the University of Padova using a remote testing system. tvLQR is an extension of the standard Linear Quadratic Regulator (LQR) control design. It is tailored for systems with time-dependent state-space matrices or where the optimal control needs to be dynamic. Representing the optimal control method, it was implemented by a separate team from DFKI RIC.

Criteria	SAC+LQR		MC-PILCO		tvLQR	
	Pendubot	Acrobot	Pendubot	Acrobot	Pendubot	Acrobot
Swingup Success	success	success	success	success	success	success
Swingup time [s]	0.65	2.06	1.43	1.1	4.2	3.98
Energy [J]	9.4	29.24	12.67	9.81	9.06	10.92
Max. Torque [Nm]	5.0	5.0	2.4	2.82	2.82	5.0
Integrated Torque [Nm]	2.21	4.57	3.48	1.27	2.57	2.27
Torque Cost [N^2m^2]	8.58	12.32	7.77	2.27	2.0	2.47
Torque Smoothness [Nm]	0.172	0.954	0.07	0.057	0.031	0.077
Velocity Cost [m^2/s^2]	44.98	193.78	94.68	242.44	137.31	100.34
RealAI Score	0.801	0.722	0.891	0.869	0.827	0.8

Table 6.1: Performance scores of various controllers for pendubot and acrobot experiments.

6 Discussion and Future Work

All three controllers are successful with both the Pendubot and Acrobot setups.

In the Pendubot setup, the performance of the SAC+LQR controller is commendable, particularly with a swift swing-up time of 0.65s. The energy consumption of the SAC+LQR controller (9.4J) is significantly lower than that of the MC-PILCO controller (12.67J) and is nearly on par with the tvLQR controller (9.06J). Additionally, its overall RealAI score is competitive, closely trailing the scores of MC-PILCO and tvLQR. However, a notable drawback is its torque smoothness; it performs the worst among the three controllers, being 2.46 times that of MC-PILCO and 5.55 times that of tvLQR.

For the Acrobot setup, the SAC+LQR controller loses its edge in both swing-up time and energy consumption. Its deficit in torque smoothness becomes even more pronounced, leading to a considerably lower RealAI score compared to the other two controllers.

In general, the SAC+LQR controller demonstrates competitive performance in simpler tasks, such as the Pendubot, especially excelling in swing-up time. However, when faced with a more complex challenge like the Acrobot, its performance declines. The MC-PILCO consistently delivers the best overall performance across both setups and is notable for its remarkably low maximum torque input and consistent torque smoothness. Conversely, the tvLQR, a non-learning-based method, highlights its effectiveness in both scenarios. While its swing-up time is relatively extended, its energy consumption and torque smoothness are commendably low, leading to a moderate RealAI score.

6.2 Interpretation of robust leaderboard

In comparison, the SAC+LQR controller achieves a moderate overall robustness score among the three controllers. It exhibits a higher resistance to model inaccuracy (71.9% for pendubot and 76.7% for acrobot) compared to MC-PILCO (45.2% for pendubot and 40.5% for acrobot) and tvLQR (75.2% for pendubot and 59.0% for acrobot). While the other two controllers demonstrate a noticeable decline when tackling the more complex acrobot task, the performance of the SAC+LQR remains consistent. Additionally, SAC+LQR offers better resistance against velocity measurement noise compared to MC-PILCO, though the top score in this category is held by tvLQR. Apart from MC-PILCO, the other two controllers display consistent and strong robustness regarding torque noise and torque response.

When considering time delay, tvLQR outperforms both SAC+LQR and MC-PILCO. As previously predicted, time delay is the Achilles heel for RL-based methods, since high latency can disrupt the Markov decision process entirely.

Criteria	SAC+LQR		MC-PILCO		tvLQR	
	Pendubot	Acrobot	Pendubot	Acrobot	Pendubot	Acrobot
Model inaccuracy [%]	71.9	76.7	45.2	40.5	75.2	59.0
Velocity noise [%]	100.0	71.4	90.5	66.7	100.0	95.2
Torque noise [%]	100.0	100.0	100.0	81.0	100.0	100.0
Torque response [%]	100.0	100.0	100.0	90.5	100.0	100.0
Time delay [%]	76.2	61.9	90.5	19.0	100.0	76.2
Overall Score	0.896	0.820	0.852	0.595	0.950	0.861

Table 6.2: Robustness scores of various controllers for pendubot and acrobot experiments.

In general, tvLQR achieves the best robustness scores for both pendubot and acrobot setups, followed by SAC+LQR, with MC-PILCO ranking last. While SAC+LQR boasts consistency in robustness across both setups, time delay remains a significant issue, limiting the robustness of RL-based methods.

6.3 Interpretation of real system leaderboard

This section is about explaining the hardware results. [to be filled]

6 Discussion and Future Work

Criteria	SAC+LQR		MC-PILCO		tvLQR	
	Pendubot	Acrobot	Pendubot	Acrobot	Pendubot	Acrobot
Swingup Success	4/10	0/10	10/10	10/10	8/10	10/10
Swingup time [s]	0.67	-	1.37	1.55	4.12	4.03
Energy [J]	37.12	-	11.66	17.95	34.02	13.75
Max. Torque [Nm]	5.0	-	4.99	5.0	5.0	2.98
Integrated Torque [Nm]	24.87	-	3.72	5.93	19.06	5.61
Torque Cost [N^2m^2]	78.7	-	8.93	11.73	51.88	3.26
Torque Smoothness [Nm]	0.774	-	0.54	0.671	0.643	0.108
Velocity Cost [m^2/s^2]	114.04	-	84.61	118.38	242.34	109.77
Best RealAI Score	0.767	-	0.843	0.82	0.695	0.822
Average RealAI Score	0.298	-	0.839	0.817	0.547	0.821

Table 6.3: Real hardware performance scores of multiple controllers for pendubot and acrobot experiments.

6.4 Conclusion and future work

This section is to talk about things to be done.

Bibliography

- [1] J. Achiam. *Spinning up in deep reinforcement learning*. 2018.
- [2] C. Aguilar-Ibañez, M. S. Suárez-Castañón, and O. O. Gutiérres-Frias. “The direct Lyapunov method for the stabilisation of the Furuta pendulum”. In: *International Journal of Control* 83.11 (2010), pp. 2285–2293.
- [3] T. Albakhali, R. Mukherjee, and T. Das. “Swing-up control of the pendubot: an impulse-momentum approach”. In: *IEEE Transactions on Robotics* 25.4 (2009), pp. 975–982.
- [4] F. Amadio, A. Dalla Libera, R. Antonello, D. Nikovski, R. Carli, and D. Romeres. “Model-based policy search using monte carlo gradient estimation with real systems application”. In: *IEEE Transactions on Robotics* 38.6 (2022), pp. 3879–3898.
- [5] J. L. Ba, J. R. Kiros, and G. E. Hinton. “Layer normalization”. In: *arXiv preprint arXiv:1607.06450* (2016).
- [6] J. T. Betts. “Survey of numerical methods for trajectory optimization”. In: *Journal of guidance, control, and dynamics* 21.2 (1998), pp. 193–207.
- [7] L. Biagiotti and C. Melchiorri. *Trajectory planning for automatic machines and robots*. Springer Science & Business Media, 2008.
- [8] W. Bickley. “Piecewise cubic interpolation and two-point boundary problems”. In: *The computer journal* 11.2 (1968), pp. 206–208.
- [9] P. Biswal and P. K. Mohanty. “Development of quadruped walking robots: A review”. In: *Ain Shams Engineering Journal* 12.2 (2021), pp. 2017–2031.
- [10] A. Bogdanov. “Optimal control of a double inverted pendulum on a cart”. In: *Oregon Health and Science University, Tech. Rep. CSE-04-006, OGI School of Science and Engineering, Beaverton, OR* (2004).
- [11] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. “Openai gym”. In: *arXiv preprint arXiv:1606.01540* (2016).
- [12] B. S. Cazzolato, Z. Prime, et al. “On the dynamics of the furuta pendulum”. In: *Journal of Control Science and Engineering* 2011 (2011).
- [13] X. Cui and H. Shi. “A*-based pathfinding in modern computer games”. In: *International Journal of Computer Science and Network Security* 11.1 (2011), pp. 125–130.
- [14] K. Furuta, M. Yamakita, and S. Kobayashi. “Swing-up control of inverted pendulum using pseudo-state feedback”. In: *Proceedings of the Institution of Mechanical Engineers, Part I: Journal of Systems and Control Engineering* 206.4 (1992), pp. 263–269.

Bibliography

- [15] A. Gasparetto, P. Boscariol, A. Lanzutti, and R. Vidoni. “Path planning and trajectory planning algorithms: A general overview”. In: *Motion and Operation Planning of Robotic Systems: Background and Practical Approaches* (2015), pp. 3–27.
- [16] A. Gasparetto, P. Boscariol, A. Lanzutti, and R. Vidoni. “Trajectory planning in robotics”. In: *Mathematics in Computer Science* 6 (2012), pp. 269–279.
- [17] A. Gasparetto and V. Zanotto. “Optimal trajectory planning for industrial robots”. In: *Advances in Engineering Software* 41.4 (2010), pp. 548–556.
- [18] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor”. In: *International conference on machine learning*. PMLR. 2018, pp. 1861–1870.
- [19] T. Hiraoka, T. Imagawa, T. Hashimoto, T. Onishi, and Y. Tsuruoka. “Dropout q-functions for doubly efficient reinforcement learning”. In: *arXiv preprint arXiv:2110.02034* (2021).
- [20] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. “Learning agile and dynamic motor skills for legged robots”. In: *Science Robotics* 4.26 (2019), eaau5872.
- [21] V. Konda and J. Tsitsiklis. “Actor-critic algorithms”. In: *Advances in neural information processing systems* 12 (1999).
- [22] B. Kouvaritakis and M. Cannon. “Model predictive control”. In: *Switzerland: Springer International Publishing* 38 (2016).
- [23] A. Kumar, Z. Fu, D. Pathak, and J. Malik. “Rma: Rapid motor adaptation for legged robots”. In: *arXiv preprint arXiv:2107.04034* (2021).
- [24] N. Lehtomaki, N. Sandell, and M. Athans. “Robustness results in linear-quadratic Gaussian based multivariable control designs”. In: *IEEE Transactions on Automatic Control* 26.1 (1981), pp. 75–93.
- [25] D. Libera, A. Turcato, N. Giacomuzzo, G. Carli, R. Romeres, A. D. Libera, N. Turcato, G. Giacomuzzo, et al. “Athletic Intelligence Olympics challenge with Model-Based Reinforcement Learning”. In: 2023. URL: <https://api.semanticscholar.org/CorpusID:261487429>.
- [26] Y. Liu and H. Yu. “A survey of underactuated mechanical systems”. In: *IET Control Theory & Applications* 7.7 (2013), pp. 921–935.
- [27] T. Luukkonen. “Modelling and control of quadcopter”. In: *Independent research project in applied mathematics, Espoo* 22.22 (2011).
- [28] K. M. Lynch and F. C. Park. *Modern robotics*. Cambridge University Press, 2017.

-
- [29] L. J. Maywald, F. Wiebe, S. Kumar, M. Javadi, and F. Kirchner. “Co-optimization of Acrobot Design and Controller for Increased Certifiable Stability”. In: *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2022, pp. 2636–2641.
- [30] X. B. Peng, P. Abbeel, S. Levine, and M. Van de Panne. “Deepmimic: Example-guided deep reinforcement learning of physics-based character skills”. In: *ACM Transactions On Graphics (TOG)* 37.4 (2018), pp. 1–14.
- [31] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann. “Stable-Baselines3: Reliable Reinforcement Learning Implementations”. In: *Journal of Machine Learning Research* 22.268 (2021), pp. 1–8. URL: <http://jmlr.org/papers/v22/20-1364.html>.
- [32] S. Saeedvand, M. Jafari, H. S. Aghdasi, and J. Baltes. “A comprehensive survey on humanoid robot development”. In: *The Knowledge Engineering Review* 34 (2019), e20.
- [33] W. Schwarting, J. Alonso-Mora, and D. Rus. “Planning and decision-making for autonomous vehicles”. In: *Annual Review of Control, Robotics, and Autonomous Systems* 1 (2018), pp. 187–210.
- [34] T. Shinbrot, C. Grebogi, J. Wisdom, and J. A. Yorke. “Chaos in a double pendulum”. In: *American Journal of Physics* 60.6 (1992), pp. 491–499.
- [35] L. Smith, I. Kostrikov, and S. Levine. “A walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning”. In: *arXiv preprint arXiv:2208.07860* (2022).
- [36] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. “Dropout: a simple way to prevent neural networks from overfitting”. In: *The journal of machine learning research* 15.1 (2014), pp. 1929–1958.
- [37] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [38] R. Tedrake. “Underactuated robotics”. In: *Algorithms for Walking, Running, Swimming, Flying, and Manipulation* (2022).
- [39] R. Tedrake, I. R. Manchester, M. Tobenkin, and J. W. Roberts. “LQR-trees: Feedback motion planning via sums-of-squares verification”. In: *The International Journal of Robotics Research* 29.8 (2010), pp. 1038–1052.
- [40] M. Towers, J. K. Terry, A. Kwiatkowski, J. U. Balis, G. d. Cola, T. Deleu, M. Goulão, A. Kallinteris, et al. *Gymnasium*. Mar. 2023. DOI: 10.5281/zenodo.8127026. URL: <https://zenodo.org/record/8127025> (visited on 07/08/2023).
- [41] T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, et al. “Benchmarking model-based reinforcement learning”. In: *arXiv preprint arXiv:1907.02057* (2019).

Bibliography

- [42] X. Xin and M. Kaneda. “New analytical results of the energy based swinging up control of the Acrobot”. In: *2004 43rd IEEE Conference on Decision and Control (CDC)(IEEE Cat. No. 04CH37601)*. Vol. 1. IEEE. 2004, pp. 704–709.
- [43] M. Yamakita, M. Iwashiro, Y. Sugahara, and K. Furuta. “Robust swing up control of double pendulum”. In: *Proceedings of 1995 American Control Conference-ACC’95*. Vol. 1. IEEE. 1995, pp. 290–295.
- [44] Y. Zheng, S. Luo, and Z. Lv. “Control double inverted pendulum by reinforcement learning with double cmac network”. In: *18th International Conference on Pattern Recognition (ICPR’06)*. Vol. 4. IEEE. 2006, pp. 639–642.

Appendix

A An appendix

You can structure appendices, just like your thesis, with the \chapter, \section, and \subsection commands. Referencing also works as usual.

If your thesis does not contain an appendix, comment out the creation of the appendix at the appropriate place in the Thesis.tex file.