# DDPG and TD3: Roboschool Environments

Amisi Fikirini
afikirini@aimsammi.org

Fay Majid Elhassan
elhassan@aimsammi.org

Vongai Mitchell Makuwaza
vmmakuwaza@aimsammi.org

Abigail Naa Amankwaa Abeo
aabeo@aimsammi.org

**AIMS-AMMI Senegal**

## Abstract

This report delves into the capabilities of Reinforcement Learning (RL) algorithms, DDPG and TD3, in training robotic agents for complex control tasks. Originally designed for Roboschool, the study transitioned to PyBullet due to Roboschool's depreciation in 2019. However, the complete adaptation of the physics environment remains a work in progress, limiting our testing in the race stadium environment. The research highlights the implementation and evaluation of these algorithms, aiming to offer insights into their practical applications and further the research in RL for robotics. The methodology, related works, and results emphasize the algorithms' performance, concluding with potential future research avenues in robotic RL applications.

## 1 Introduction

Reinforcement learning (RL) has demonstrated significant promise in training agents to perform tasks across various environments. The ability to train robotic agents using RL techniques holds immense potential for real-world applications, from autonomous vehicles to industrial automation. However, the complex and dynamic nature of robotic control tasks presents unique challenges that demand sophisticated RL algorithms.

Advanced algorithms like DDPG (Deep Deterministic Policy Gradient) and TD3 (Twin Delayed Deep Deterministic Policy Gradient) have emerged as powerful tools for addressing these challenges. Specifically designed to handle environments with continuous action spaces, these algorithms have shown remarkable adaptability and performance. In this context, Roboschool, a set of open-source environments, offers a valuable testing ground for evaluating the capabilities of DDPG and TD3 agents.

The motivation behind this project is twofold: first, to explore the potential of DDPG and TD3 algorithms in training robotic agents for complex control tasks, and second, to assess their adaptability and generalization capabilities in realistic environments provided by Roboschool. By implementing, training, and evaluating DDPG and TD3 agents in selected Roboschool environments, we aim to contribute to the growing body of research in RL for robotics and provide insights into the practical applications of these algorithms.

The rest of this paper is organized as follows: In Section 2, we review related work in the field of reinforcement learning, providing context for our research. Section 3 presents the methodology and architecture of our implementation, detailing our approach to training and evaluation. Section 4 outlines the results and discussions, shedding light on the performance of DDPG and TD3 agents in Roboschool environments. We conclude our findings in Section 5 and discuss potential directions for future research in Section 6.

## 2 Literature Review

The Deep Q Network (DQN) algorithm by [1] Mnih et al. (2015) marked a significant milestone in RL. DQN combined advances in deep learning for sensory processing with reinforcement learning, enabling agents to learn directly from raw sensory input. This work introduced the concept of using convolutional neural networks (CNNs) to process visual input and demonstrated impressive results on a range of Atari 2600 games. The success of DQN highlighted the potential of deep learning for RL and inspired subsequent research.

Robotic control tasks often involve continuous action spaces, which pose challenges for traditional RL algorithms designed for discrete action spaces. [2] Lillicrap et al. (2016) introduced the Deep Deterministic Policy Gradient (DDPG) algorithm, which extended the actor-critic architecture to effectively handle continuous action spaces. DDPG, through the use of neural networks, demonstrated remarkable capabilities in complex control tasks with high-dimensional state spaces.

While DDPG was a breakthrough in handling continuous action spaces, it had limitations related to overestimation bias and instability. [3] Fujimoto et al. (2018) addressed these challenges by introducing the Twin Delayed Deep Deterministic Policy Gradient (TD3) algorithm. TD3 incorporates twin critics and a noise-regularized target policy smoothing technique, resulting in improved training stability and sample efficiency.

The application of RL-trained agents in diverse environments is a topic of growing interest. Research on transfer learning has explored how agents can adapt their learned policies to new scenarios. Furthermore, visualization of RL-trained agents in different environments provides insights into their adaptability and generalization capabilities.

In our project, we aim to bridge these concepts by transferring agents trained in Roboschool environments to a racing environment, thereby assessing their ability to generalize and perform in novel scenarios.

### 2.1 Research Gap and Project Objectives

While significant progress has been made in RL for robotic control, there is still a need for comprehensive studies that apply DDPG and TD3 to Roboschool environments. This project aims to address this gap by implementing and evaluating these algorithms in challenging robotic control tasks.

In the following sections, we will delve into the details of the Roboschool environments, algorithm implementations, training procedures, results, and conclusions of our investigation.

## 3 Methodology

### 3.1 Roboschool Environments

The Roboschool environment is an essential component of this project. It offers realistic physics simulations for a variety of robotic control tasks. While Roboschool provides a framework for training agents, the choice of algorithms and training strategies significantly impacts the agents' performance and adaptability. The integration of DDPG and TD3 algorithms into the Roboschool framework represents a novel approach to training robotic agents for complex control tasks. (It is important to note that the Roboschool environment has been deprecated since 2019, hence we used the Pybullet environment [4]).

In light of this we selected four distinct Roboschool environments: humanoid, hopper, walker2D and half cheetah, to train our DDPG and TD3 agents as follows:

1. **Humanoid Environment:**
   - Simulates a 3D bipedal humanoid robot.
   - Challenges agents with complex balance and walking behaviors.
   - Demands high-dimensional action space control and poses significant computational demands.

2. **Hopper Environment:**

- Features a single-legged robot with a hopping task.
- Offers an intermediate level of complexity, testing balance and forward motion control.

3. **Walker2D Environment:**

- Involves a two-dimensional bipedal robot tasked with walking forward.
- Requires control of two legs for balance and locomotion.
- Tests algorithms' ability to handle multi-joint and multi-leg control.

4. **Half Cheetah Environment:**

- Simulates a half-cheetah robot designed for running.
- Presents a high level of complexity, demanding fast and dynamic locomotion control while maintaining balance.
- Computationally demanding due to precise motor control requirements.

These environments offer varying levels of complexity and challenges and the agents were trained with varying set of hyperparameters.

## 3.2 JAX and Haiku

Leveraged JAX and Haiku in our implementation. JAX is a versatile machine learning library that combines Autograd for automatic differentiation with XLA for GPU and TPU acceleration and has the ability to automatically differentiate Python and NumPy functions, support for reverse and forward-mode differentiation, and just-in-time compilation (jit) for efficient execution on hardware accelerators. JAX's composable function transformations, such as grad for differentiation and jit for compilation, enable users to express complex algorithms in Python. JAX offers tremendous potential for machine learning research and development.

Haiku is a high-level neural network library built on JAX, simplifying the creation and management of deep neural networks. It offers modularity, parameter handling, and seamless integration with JAX's powerful capabilities, streamlining deep learning model development. Haiku enhances code structure and readability, making it a valuable tool for efficient neural network research and implementation.

## 3.3 Scheduled Noise

The scheduled noise plays an important role in enhancing exploration and learning, allowing agents to adapt to dynamic environments. Noise scheduling in our implementation gradually reduces the magnitude of exploration noise added to an agent's action over the course of the training. The key components are the initial standard deviation of the noise, minimum allowed standard deviation and decay rate which governs the rate at which noise reduction occurs. By decreasing the noise level over time, the agent gradually moves from an exploratory phase to a more deterministic policy, improving its learning process. Proper tuning of hyperparameters such as initial standard deviation, minimum deviation and decay rate is essential and can significantly impact the agent's performance.

## 3.4 Policy Gradients

Policy gradients are a class of reinforcement learning algorithms used for training policies to maximize cumulative rewards in a given environment. Unlike value-based methods (e.g., Q-learning), which estimate the value of states or state-action pairs, policy gradients focus on directly optimizing the policy itself. The objective in policy gradient methods is to find the parameters of the policy that maximize the expected cumulative reward, often represented as the expected return $J(\theta)$.

The policy gradient algorithms that we shall look more into are the Deep Deterministic Policy Gradient (DDPG) and the Twin Delayed Deep Deterministic Policy Gradient (TD3). DDPG and TD3 are considered off-policy reinforcement learning algorithms because they make use of experience replay, target networks, and separate policy and value function updates. These characteristics enable them to learn efficiently from past experiences and contribute to the stable and effective training of agents, particularly in environments with continuous action spaces.

## 3.5 DDPG Implementation

Our implementation adheres to the core principles of DDPG while utilizing the JAX library for efficient computation and consisted of actor and critic neural networks with suitable architectures for the state and action spaces. In DDPG, the actor network is responsible for selecting continuous-valued actions based on the current state, while the critic network(s) estimate the action-value function to assess the quality of the chosen actions. The use of target networks and the training of the actor and critic networks are essential components of DDPG, which helps in training agents for reinforcement learning tasks with continuous action spaces.

One of the fundamental aspects of our DDPG implementation is the use of soft updates. Soft updates, as introduced by [2]Lillicrap et al. (2016), play a pivotal role in stabilizing the training process. We employ the soft update function to blend the parameters of the target networks with those of the online networks. This technique reduces sudden changes in target values during training, mitigating training instability. Hyperparameters, such as learning rates and batch sizes, were tuned for each environment. We employed Ornstein-Uhlenbeck noise for exploration. In the DDPG algorithm documentation provided by OpenAI [5], detailed information about the Deep Deterministic Policy Gradient algorithm is found as follows.

---

**Algorithm 1** DDPG Algorithm with Soft Updates

---

1: **Input:** State dimension, action dimension, learning rates, buffer size, etc.
2: Initialize actor network $f_\theta$ and critic network $Q_\phi$ with random weights.
3: Initialize target networks $f_{\theta'}$ and $Q_{\phi'}$ with weights identical to the actor and critic.
4: Initialize replay buffer $\mathcal{B}$.
5: Initialize noise schedule $\mathcal{N}$.
6: **for** each episode **do**
7:     Observe initial state $s_0$
8:     **for** each timestep **do**
9:         Select action $a_t$ using exploration policy derived from $f_\theta(s_t)$
10:         Execute $a_t$, observe next state $s_{t+1}$, reward $r_t$ and done signal $d_t$
11:         Store $(s_t, a_t, r_t, s_{t+1}, d_t)$ in $\mathcal{B}$
12:         **if** len($\mathcal{B}$) > batch size **then**
13:             Sample a batch from $\mathcal{B}$
14:             Compute target: $y = r + \gamma(1-d)Q_{\phi'}(s_{t+1}, f_{\theta'}(s_{t+1}))$
15:             Update critic: $\phi \leftarrow \phi - \alpha_\phi \nabla_\phi (Q_\phi(s_t, a_t) - y)^2$
16:             Update actor: $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta f_\theta(s_t) Q_\phi(s_t, f_\theta(s_t))$
17:             Soft update target networks: $\theta' \leftarrow \tau\theta + (1-\tau)\theta'$
18:             $Q_{\phi'} \leftarrow \tau Q_\phi + (1-\tau)Q_{\phi'}$
19:         **end if**
20:         Update noise schedule $\mathcal{N}$
21:     **end for**
22: **end for**

---

## 3.6 TD3 Implementation

In the TD3 algorithm documentation by OpenAI [6], the authors provide detailed information about the Twin Delayed Deep Deterministic Policy Gradient algorithm. TD3 is an improvement over DDPG designed to enhance stability and reduce overestimation bias. It employs two critics to estimate the action-value function and uses the minimum of their values to reduce overestimation.

In the TD3 algorithm, the actor network is responsible for selecting actions, while the twin critic networks estimate the action-value function. The presence of two critic networks and the use of target networks are key features of TD3 that contribute to its improved stability and learning performance in continuous action spaces. These neural networks are trained iteratively to improve the agent's policy and Q-value estimates, ultimately leading to better decision-making in reinforcement learning tasks.

**Algorithm 2** TD3 Algorithm with Soft Updates

---

1: **Input:** State dimension, action dimension, learning rates, buffer size, etc.
2: Initialize actor network $f_\theta$ and critic networks $Q_{\phi_1}, Q_{\phi_2}$ with random weights.
3: Initialize target networks $f_{\theta'}$ and $Q_{\phi'_1}, Q_{\phi'_2}$ with weights identical to the actor and critics.
4: Initialize replay buffer $\mathcal{B}$.
5: Initialize noise schedule $\mathcal{N}$.
6: **for** each episode **do**
7:     Observe initial state $s_0$
8:     **for** each timestep **do**
9:         Select action $a_t$ using exploration policy derived from $f_\theta(s_t)$
10:         Execute $a_t$, observe next state $s_{t+1}$, reward $r_t$ and done signal $d_t$
11:         Store $(s_t, a_t, r_t, s_{t+1}, d_t)$ in $\mathcal{B}$
12:         **if** len($\mathcal{B}$) > batch size **then**
13:             Sample a batch from $\mathcal{B}$
14:             Compute target: $y = r + \gamma(1 - d)\min_{i=1,2} Q_{\phi'_i}(s_{t+1}, f_{\theta'}(s_{t+1}))$
15:             Update critics: $\phi_i \leftarrow \phi_i - \alpha_\phi \nabla_{\phi_i}(Q_{\phi_i}(s_t, a_t) - y)^2$
16:             **if** timestep % policy delay == 0 **then**
17:                 Update actor: $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta f_\theta(s_t) Q_{\phi_1}(s_t, f_\theta(s_t))$
18:                 Soft update target networks: $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$
19:                 $Q_{\phi'_1} \leftarrow \tau Q_{\phi_1} + (1 - \tau)Q_{\phi'_1}$
20:                 $Q_{\phi'_2} \leftarrow \tau Q_{\phi_2} + (1 - \tau)Q_{\phi'_2}$
21:             **end if**
22:         **end if**
23:         Update noise schedule $\mathcal{N}$
24:     **end for**
25: **end for**

---

## 3.7 Training Process

The training process includes the following steps:

- Initialization of networks, target networks, replay buffers, and hyperparameters.

- Exploration strategies, such as Ornstein-Uhlenbeck noise, to encourage exploration.

- Training loop to collect experiences and improve the policy.

- Periodic target network updates to stabilize training.

- Convergence monitoring and stopping criteria.

### 3.7.1 Initial Focus: Inverted Pendulum

We began by implementing DDPG and TD3 algorithms in the Inverted Pendulum environment. This provided a straightforward and computationally efficient starting point to understand algorithm behavior, fine-tune parameters, and verify implementations. It served as a performance benchmark before advancing to more complex tasks in the Roboschool framework.
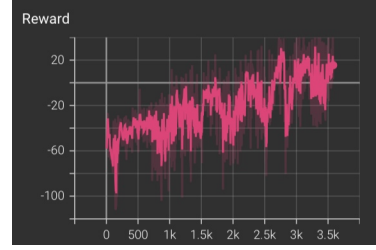
## 4 Results and Discussions

We tried various experiments and played with different kinds of hyperparameters. Due to computational power and memory, we could not run as many episodes as compared to what we found in papers. We noticed that most papers run their codes for 1 million episodes. The maximum number of episodes we could run was 100000, even this was only achieved for the hopper environment, as among the environments, the hopper trained faster.

### 4.1 Results

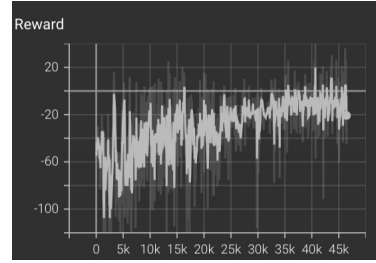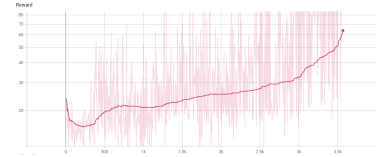| Hyperparameter | Value |
|---|---|
| Learning Rate (Actor) | 1e-5 |
| Learning Rate (Critic) | 1e-5 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.001 |
| Batch Size | 1024 |
| Replay Buffer Size | 1000000 |

(a) Hyperparameters Used in Humanoid Training



(b) Training Performance on the Humanoid

Figure 1: Hyperparameters and Training Performance for humanoid using DDPG

| Hyperparameter | Value |
|---|---|
| Learning Rate (Actor) | 3e-4 |
| Learning Rate (Critic) | 3e-4 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

(a) Hyperparameters Used in Humanoid Training



(b) Training Performance on the Humanoid

Figure 2: Hyperparameters and Training Performance for humanoid using TD3

| Hyperparameter | Value |
|---|---|
| Learning Rate (Actor) | 3e-4 |
| Learning Rate (Critic) | 3e-4 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

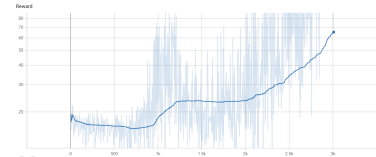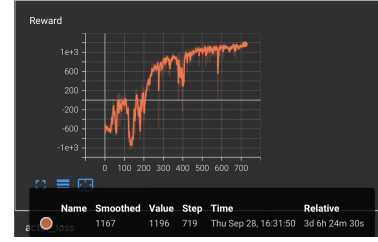(a) Hyperparameters Used in Hopper Training



(b) Training Performance on the Hopper Environment

Figure 3: Hyperparameters and Training Performance for Hopper using DDPG

| Hyperparameter | Value |
|---|---|
| Learning Rate (Actor) | 3e-4 |
| Learning Rate (Critic) | 3e-4 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

(a) Hyperparameters Used in Hopper Training



(b) Training Performance on the Hopper Environment

Figure 4: Hyperparameters and Training Performance for Hopper using TD3

| Hyperparameter | Value |
| --- | --- |
| Actor Network | 256 |
| Critic network | 512 |
| Learning Rate (Actor) | 1e-3 |
| Learning Rate (Critic) | 1e-3 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 1000000 |

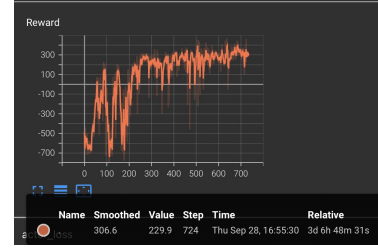(a) Hyperparameters Used in HalfCheetah Training



(b) Training Performance on the HalfCheetah Environment

Figure 5: Hyperparameters and Training Performance for HalfCheetah using DDPG

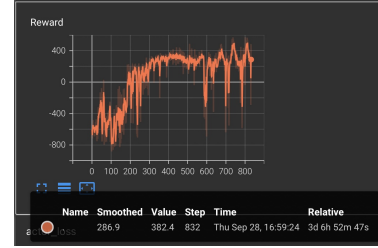| Hyperparameter | Value |
| --- | --- |
| Learning Rate (Actor) | 3e-4 |
| Learning Rate (Critic) | 3e-4 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 1000000 |

(a) Hyperparameters Used in HalfCheetah Training



(b) Training Performance on the Halfcheetah Environment

Figure 6: Hyperparameters and Training Performance for HalfCheetah using TD3

| Hyperparameter | Value |
| --- | --- |
| Learning Rate (Actor) | 1e-3 |
| Learning Rate (Critic) | 1e-3 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

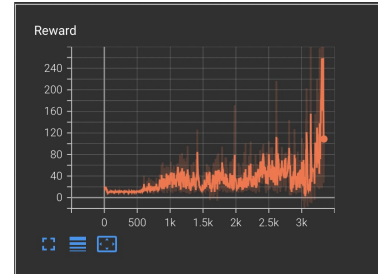(a) Hyperparameters Used in HalfCheetah Training



(b) Training Performance on the Halfcheetah Environment

Figure 7: Hyperparameters and Training Performance for HalfCheetah using TD3 Buffer size change

| Hyperparameter | Value |
| --- | --- |
| Actor Network | 256 |
| Critic network | 512 |
| Learning Rate (Actor) | 1e-4 |
| Learning Rate (Critic) | 1e-4 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.001 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

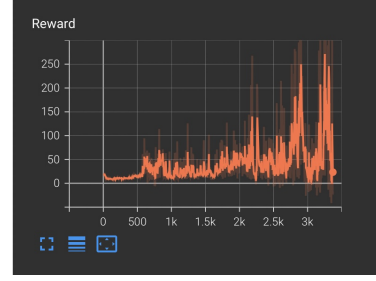(a) Hyperparameters Used in Walker Training



(b) Training Performance on the Walker Environment

Figure 8: Hyperparameters and Training Performance for Walker using DDPG

| Hyperparameter | Value |
|---|---|
| Learning Rate (Actor) | 1e-3 |
| Learning Rate (Critic) | 1e-3 |
| Discount Factor ($\gamma$) | 0.99 |
| Soft Update Rate ($\tau$) | 0.005 |
| Batch Size | 256 |
| Replay Buffer Size | 200000 |

(a) Hyperparameters Used in Walker Training



(b) Training Performance on the Walker Environment

Figure 9: Hyperparameters and Training Performance for Walker using TD3

## 4.2 Discussions

1. **Humanoid Environment:**
   - The training outcomes of our experiments revealed distinct performance characteristics for DDPG and TD3.
   - DDPG demonstrated a more favorable learning trajectory, achieving positive rewards within the range of 0 to 20. However, the experiment was halted before exploring whether the agent could attain higher rewards.
   - TD3, despite extensive hyperparameter tuning, exhibited challenges in achieving convergence and maintained only minimal positive rewards throughout training.

2. **Hopper Environment:**
   - DDPG was able to exhibit better performance compared to TD3.
   - One hyperparameter we found very interesting is the buffer size. We saw how a larger replay buffer (in our case, 200000 as compared to 100000) helped to improve sample efficiency. This means that the agent can reuse a larger and more diverse set of past experiences during training, reducing the need for collecting new data from the environment. With a small buffer, the agent may forget valuable experiences quickly, leading to inefficient learning.

3. **Walker2D Environment:**
   - Involves a two-dimensional bipedal robot tasked with walking forward.
   - DDPG demonstrated a more favorable learning trajectory, achieving positive rewards within the range of 20 to 250. However,with notice in improvement but due to limit resources the training was taking long with no remarkable improvement the experiment was halted before exploring whether the agent could attain higher rewards.
   - TD3 seemed to be challenging with fluctuations o the positive and slight remarkable improvement

4. **Half Cheetah Environment:**
   - The DDPG algorithm, under the specified parameters, demonstrates distinctive learning characteristics. Observing the convergence speed, DDPG starts to stabilize around 800 episodes, suggesting an efficient learning process in the early stages of training. However, fluctuations in reward values are observed, indicating variations in the algorithm's performance and its ability to maintain a stable policy.
   - Adjustments to the actor and critic networks, including a reduced actor and an expanded critic, were pivotal, affecting DDPG's learning stability and overall efficiency. The smaller actor facilitated more precise policy approximations, enhancing reliability in action selections. In contrast, the larger critic provided refined value estimations, offering richer feedback for policy improvements. These alterations significantly refined the learning process, allowing for optimal and stable policies amidst the environmental complexities, underscoring the influence of network architecture in achieving robust learning outcomes in DDPG implementations.

- The exploration of buffer size revealed its critical role in learning efficiency. Increasing it to 200,000 from 100,000 improved sample efficiency by enabling the utilization of a more diverse set of past experiences, reducing the reliance on new data and mitigating the loss of valuable experiences, thus optimizing learning outcomes.

These findings underscore the complexity of reinforcement learning tasks and emphasize the importance of algorithm selection and hyperparameter tuning in achieving desired training outcomes. Future work may involve further experimentation and exploration of alternative algorithms and hyperparameter settings to address the challenges encountered in this study.

## 5    Conclusion

This study embarked on a comprehensive exploration of the DDPG and TD3 algorithms across multiple environments, revealing nuanced insights into the algorithms' learning trajectories, stability, and overall performances. In various environments like Hopper and Half Cheetah, DDPG exhibited favorable learning trajectories and managed to achieve positive rewards, illustrating its adaptability and efficiency in complex tasks. However, for Humanoid and Walker2D, the training process had to be halted in some cases before fully exploring the potential for higher rewards due to resource constraints.

TD3, conversely, struggled with convergence and demonstrated minimal positive rewards, indicating the challenges associated with its implementation in the given tasks. These disparities underscore the inherent complexities in reinforcement learning tasks and the pivotal role of algorithm selection in navigating these complexities.

In particular, the study highlighted the significance of architectural adjustments and hyperparameter tuning in optimizing learning outcomes. The modification of the actor and critic networks in DDPG, along with the strategic enlargement of the replay buffer, proved instrumental in enhancing learning stability and sample efficiency. These adjustments facilitated more precise policy approximations and richer, more informative feedback for policy improvements, enabling the algorithm to navigate the environmental complexities more effectively and converge to more optimal and stable policies.

These findings emphasize the multifaceted nature of reinforcement learning and the intricate interplay between algorithm selection, architectural considerations, and hyperparameter tuning in determining the success of the learning process. They provide a foundation for future research endeavors, prompting further experimentation and exploration of alternative algorithms and hyperparameter configurations to address the identified challenges and continue the pursuit of more robust and efficient reinforcement learning models.

## 6    Future work

Future endeavors may delve deeper into resolving the observed challenges and exploring the untapped potentials of these algorithms. Further refinements in architectural configurations and hyperparameter settings, coupled with extensive experimentation, hold the promise of unveiling more advanced strategies to harness the capabilities of reinforcement learning in diverse and complex environments, paving the way for groundbreaking advancements in the field.

# References

[1] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

[2] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2016.

[3] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.

[4] OpenAI. Roboschool deprecated, pybullet used now, 2019.

[5] OpenAI. Ddpg: Deep deterministic policy gradient, 2023.

[6] OpenAI. Td3: Twin delayed deep deterministic policy gradient, 2023.

[7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 2012.

# 7   Appendix C - Code

Link to our github repo