

# **Integrating Deep Reinforcement Learning in Game AI**

## **Project Topic Analysis**

Submitted for the MSc in  
Computer Science for Games Development

May 2023

by

**Ronit Mishra**

Word Count: 2998

## Table of Contents

|    |   |    |
|----|---|----|
| 1  | Introduction.....                                     | 3  |
| 2  | Aim and Objectives .....                              | 4  |
| 3  | Background Context and Overview .....                 | 5  |
| 4  | Expanded Specification and Analysis .....             | 7  |
| 5  | Project Management .....                              | 12 |
| a. | Task List .....                                       | 12 |
| b. | Gantt Chart .....                                     | 14 |
| c. | Risk Analysis .....                                   | 15 |
|    | Appendix A: Deep Reinforcement Learning Methods ..... | 17 |
|    | References .....                                      | 19 |

# 1 Introduction

This project aims to explore how Deep Reinforcement Learning (DRL) can be integrated into Game AI, with a focus on designing an intelligent and adaptive Game AI. The objective is to develop game agents that possess cognitive abilities, reasoning capabilities, and problem-solving skills, allowing them to make informed decisions within the game environment.

Intelligence in game agents refers to their capacity to analyze the game state, evaluate potential actions, anticipate consequences, and select optimal strategies. These agents employ various techniques such as pattern recognition, decision trees, neural networks, and reinforcement learning to simulate human-like intelligence and maximize their chances of success.

Adaptability in game agents refers to their ability to adjust and modify their behavior based on changing circumstances in the game. An adaptable game agent can dynamically respond to new challenges, unexpected events, and evolving game states by altering its strategies, decision-making processes, or action selection. This adaptability enhances the agent's performance and improves its chances of achieving better outcomes.

Deep Reinforcement Learning (DRL) has emerged as a powerful approach in training game agents. By combining deep learning and reinforcement learning techniques, DRL enables agents to learn and excel at various games. Algorithms such as Deep Q-Networks (DQN) and AlphaGo have demonstrated significant success in this field.

The project will involve investigating different Reinforcement Learning algorithms, Neural Network architectures, and Game AI techniques to determine the most effective approach for integrating reinforcement learning into Game AI. The goal is to design and implement a Game AI system that employs Deep Reinforcement Learning, enabling the creation of intelligent and responsive game agents that can learn from their interactions with the game environment.

To achieve this, the Game AI system will need to perceive the game state, take actions, and receive feedback on its performance. Over time, the system will improve its decision-making abilities through continuous learning and optimization. The project will also involve evaluating the performance of the Game AI system with multiple agents and testing its effectiveness in various game scenarios.

Ultimately, this project aims to contribute to the advancement of intelligent and adaptive game agent technologies by exploring Deep Reinforcement Learning.

## 2 Aim and Objectives

*Design and implement Deep Reinforcement Learning based Game AI System*

### **Objective 1 – Comprehensive literature exploration**

Build upon the existing collection of papers collected while laying the groundwork as part of module 700111 and further explore the literature on DRL algorithms and their applications in Game AI systems, ensuring a comprehensive understanding within the first three weeks while allocating sufficient time and resources for the task.

### **Objective 2 – Create a detailed architectural design**

Design the Game AI system utilizing established design principles and best practices within three weeks. The design should include detailed architecture that incorporates DRL algorithms, agent-environment interactions, and other learning components.

### **Objective 3 – Develop and integrate all required functionalities based on the design plan**

Develop and integrate all components, including the neural network, reward mechanisms, and training processes. Utilize suitable programming languages and frameworks to implement the system.

Complete the implementation phase within six weeks.

### **Objective 4 – Testing the performance of the system**

Utilize appropriate testing methodologies to create and execute test cases for assessing system's behavior and validate its adaptability. The results should be recorded using appropriate metrics so that useful conclusion could be easily derived from the test results.

This phase should start immediately after the design objective has been successfully achieved.

Once the test cases are ready, remainder tasks from the testing phase should be conducted in parallel to the implementation phase.

### **Objective 5 – Reporting and Documentation**

Prepare comprehensive documentation that effectively communicates the design and implementation details of the project and publish it on GitHub to facilitate easy understanding, collaboration, and future reference.

### 3 Background Context and Overview

#### Evolution of Deep Reinforcement Learning Integration in Game AI

In 2013, Volodymyr Mnih, a research scientist from Google DeepMind published his seminal work on Deep Reinforcement Learning, "Playing Atari with Deep Reinforcement Learning," established the foundation for integrating deep neural networks with reinforcement learning algorithms, particularly Q-learning. This pioneering study demonstrated the impressive accomplishment of a learning agent achieving human-level performance across various Atari 2600 games (Mnih et al., 2013).

In 2015, van Hasselt's study addressed the prevalent issue of action value overestimation in reinforcement learning and introduced Double Q-learning as a significant breakthrough to mitigate this problem. This work provided invaluable insights into improving the stability and performance of learning agents in the field of Reinforcement Learning (van Hasselt et al., 2015).

Expanding on prior research, Mnih et al.'s paper in 2015, "Human-level control through Deep Reinforcement Learning," demonstrated the exceptional capabilities of deep Q-networks (DQNs) to learn directly from raw pixel inputs and achieve human-level performance across a diverse range of Atari 2600 games. This study revealed the profound potential of Deep Reinforcement Learning in the domain of Game AI (Mnih et al., 2015).

In 2016, Mnih et al. introduced the concept of asynchronous reinforcement learning algorithms in their paper, "Asynchronous Methods for Deep Reinforcement Learning". This research showcased the Asynchronous Advantage Actor-Critic (A3C) approach, which enabled concurrent and asynchronous learning of multiple agents, revolutionizing the scalability and efficiency of training Deep Reinforcement Learning agents. Additionally, the study shed light on critical parallelization techniques, further advancing the field (Mnih et al., 2016).

OpenAI's significant contribution in 2017 came in the form of the influential paper "Proximal Policy Optimization Algorithms". Authored by John Schulman, a research scientist and co-founder of OpenAI, this paper introduced the Proximal Policy Optimization (PPO) algorithm, which revolutionized reinforcement learning in continuous control domains. The state-of-the-art technique laid a robust foundation for developing adaptive Game AI agents capable of learning and refining their policies over time (Schulman et al., 2017).

Another notable publication in 2017 from OpenAI presented the innovative work "Deep Reinforcement Learning from Human Preferences", introducing the Deep Reinforcement Learning from Human Feedback (DRLHF) approach. By integrating human preferences with reinforcement learning, this methodology enabled agents to adapt their behavior to align with human-like tendencies, representing a significant step toward cultivating game agents capable of emulating human-like decision-making (Christiano et al., 2017).

In 2019, OpenAI established itself as a world leader in the field of Deep Reinforcement Learning with "OpenAI Five", which showcased the convergence of reinforcement learning with other cutting-edge techniques. This paper outlined the development of OpenAI Five, a system comprising a team of AI agents proficient in playing the complex

game Dota 2 at a highly competitive level. OpenAI Five's remarkable performance against professional human players underscored the successful integration of reinforcement learning in the realm of complex multiplayer games (OpenAI et al., 2019).

The past decade has seen significant advancements in Deep Reinforcement Learning in Game AI, with the introduction of DQN, A3C, PPO and DRLHF algorithm. These innovations have enabled the development of adaptive Game AI agents capable of learning, refining policies, and emulating human-like decision-making.

### **Practical Challenges in Implementing Deep Reinforcement Learning for Game AI**

There are still several challenges and problems in the implementation of DRL

- Sample Efficiency:

Training DRL agents can require a large number of interactions with the environment, making it time-consuming and computationally expensive. Improving sample efficiency is an ongoing research focus to reduce the number of interactions required for effective training.

- Generalization and Transfer Learning:

Agents trained in one game or environment often struggle to transfer their learned skills to new, unseen scenarios. Generalizing DRL agents to different games or tasks remains a significant challenge.

- Exploration and Reward Design:

Effective exploration strategies are crucial for discovering optimal policies. Designing appropriate reward functions that guide the agent's behavior and provide useful feedback is another open problem in DRL.

- Stability and Reproducibility:

Training DRL models can be unstable, making it difficult to reproduce results and achieve consistent performance across different runs. Researchers are actively working on developing more stable and robust algorithms.

Overall, while DRL has made remarkable advancements and found applications in various industries, there are still several areas that require further research and development to improve the implementation and address existing challenges (Shao et al., 2019).

## 4 Expanded Specification and Analysis

### System specification

The following high-level module specification outlines the key components of the system. The actual implementation may require further breakdown of modules and additional components based on specific requirements and design considerations.

1. World Interface:

- Responsible for interfacing with the game environment and providing necessary inputs to the AI system.
- Handles game state observations, rewards, and actions.

2. Learning Module:

- Core module that implements the Deep Reinforcement Learning algorithm.
- Utilizes neural networks to approximate the value or policy functions.
- Includes components such as experience replay, action selection, and learning updates.

3. Training Module:

- Manages the training process of the game agents.
- Implements the necessary algorithms to optimize the agent's behavior over time.
- Handles exploration-exploitation trade-offs and policy updates.

4. Game Agents:

- Represents an individual game agent that interacts with the game environment.
- Receives observations from the environment and selects actions based on the learned policy.
- Interacts with the Deep Reinforcement Learning Module to update its knowledge and improve performance.

5. User Interface:

- Provides a user-friendly interface for interacting with the Game AI system.
- Allows users to configure training parameters, monitor training progress, and visualize game agent behavior.

## Algorithm design considerations for the learning module

### Deep Q-Networks (DQN)

DQN excels in learning directly from raw sensory inputs, making it ideal for game environments where pixel-level information is pivotal. DQN employs a type of deep neural network known as a convolutional neural network (CNN). CNNs are particularly effective for processing visual input data, such as raw pixel inputs from images or frames in video games.

#### Strengths

- DQN is a model-free algorithm that can learn directly from raw sensory inputs, such as pixels in a game environment.
- It has been successful in solving complex tasks in reinforcement learning, including Atari games.
- DQN incorporates experience replay, which helps to stabilize learning and improve sample efficiency.

#### Weaknesses

- DQN suffers from high sample complexity, as it requires a large number of interactions with the environment to learn effective policies.
- It can be unstable during training and may experience difficulties in converging to optimal solutions.
- DQN is known to overestimate action values, which can lead to suboptimal policies.

### Asynchronous Advantage Actor-Critic (A3C)

A3C takes a distinct approach by leveraging multiple asynchronous agents that learn and explore concurrently. This parallel training setup augments exploration capabilities and reduces sample correlation. A3C combines the merits of policy gradient and value-based methods, endowing it with versatility for gaming tasks.

A3C also utilizes a type of deep neural network, specifically a combination of two components: an actor network and a critic network. The actor network is typically a feedforward neural network that approximates the policy, while the critic network estimates the value function. Both networks can be implemented using fully connected layers or any other appropriate architecture.

#### Strengths

- A3C's asynchronous nature allows for parallel training of multiple game-playing agents, leading to improved exploration and reduced correlation between samples.
- It combines the advantages of both policy gradient methods and value-based methods, making it a versatile approach for gaming tasks.
- A3C has been successful in solving complex game scenarios.



## Weaknesses

- The asynchronous nature of A3C introduces complexities in terms of communication and synchronization among the agents, which can be challenging to manage.
- A3C may suffer from high variance in gradient estimates, resulting in slower convergence or instability during training.
- Scaling A3C to games with high-dimensional observations and large action spaces can be difficult.

## Proximal Policy Optimization (PPO)

PPO addresses stability and sample efficiency concerns through a trust region approach. By ensuring stable and monotonically improving training, it emerges as an appealing option for gaming applications.

PPO can utilize various types of deep neural networks as function approximators. Typically, PPO employs feedforward neural networks with fully connected layers to represent the policy and value function. However, other architectures, such as recurrent neural networks (RNNs), can also be used depending on the specific application.

## Strengths

- PPO ensures stable and monotonic improvement during training by using a trust region approach to update policy parameters.
- It is effective in a wide range of gaming tasks, including continuous control and robotic manipulation.

## Weaknesses

- PPO can be sensitive to hyperparameter choices, and finding the right hyperparameters may require extensive tuning.
- It may suffer from slow convergence in some scenarios, especially when dealing with complex environments or large action spaces.
- PPO typically requires a large amount of interaction data to achieve good performance.

## Deep Reinforcement Learning from Human Feedback (DRLHF)

DRLHF leverages human demonstrations or feedback to expedite the learning process. By harnessing human expertise, it aims to enhance sample efficiency and furnish agents with valuable guidance. This approach holds promise in gaming domains where human knowledge is readily available.

DRLHF integrates Deep Reinforcement Learning with human preferences or feedback. The neural network architecture used in DRLHF can vary depending on the specific implementation. It may involve feedforward neural networks or other architectures, depending on how the human feedback is incorporated into the learning process.

## Strengths

- DRLHF leverages human demonstrations or feedback to guide the learning process, enabling faster learning and improved sample efficiency in gaming tasks.
- It allows for the integration of human expertise into the game-playing agents, making it suitable for tasks where human knowledge is available.

## Weaknesses

- Obtaining high-quality human demonstrations for gaming tasks can be challenging and time-consuming.
- DRLHF assumes that human demonstrations are reliable and represent optimal behavior, which may not always hold true in the gaming industry.
- Combining human feedback with reinforcement learning adds complexity to the algorithm design and integration process.

## Model Training Considerations:

Consider the following aspects while designing the training process:

### Exploration Strategies:

Implement exploration strategies, such as  $\epsilon$ -greedy exploration, Boltzmann exploration, or noisy networks, to encourage the AI agent to explore the game environment effectively.

### Reward Design:

Develop appropriate reward functions that provide meaningful feedback to guide the agent's behavior. Consider shaping the rewards to encourage desired actions and discourage undesired ones. Explore techniques like reward shaping, intrinsic motivation, or curriculum learning to improve training efficiency.

### Sample Efficiency Techniques:

Investigate sample-efficient approaches, such as prioritized experience replay, Hindsight Experience Replay (HER), or model-based reinforcement learning, to reduce the number of interactions required for effective training.

### Transfer Learning and Generalization:

Incorporate transfer learning techniques to improve the agent's ability to generalize its learned skills to new games or scenarios. Consider methods like domain adaptation, pretraining on related tasks, or using auxiliary tasks to transfer knowledge.

## Evaluation and Testing Module:

Assess the performance of the Game AI system based on the following criteria:

### Responsiveness:

Evaluate how quickly the Game AI system reacts to changes in the game environment and the actions of the opposing game agents.

### Decision Quality:

Analyze the effectiveness of the decisions made by the Game AI system and their impact on the game outcome.

### Adaptability:





Test the ability of the Game AI system to adapt its strategy based on different game scenarios and opponent behaviors.

### Resource Usage:

Measure the computational resources consumed by the Game AI system during gameplay.

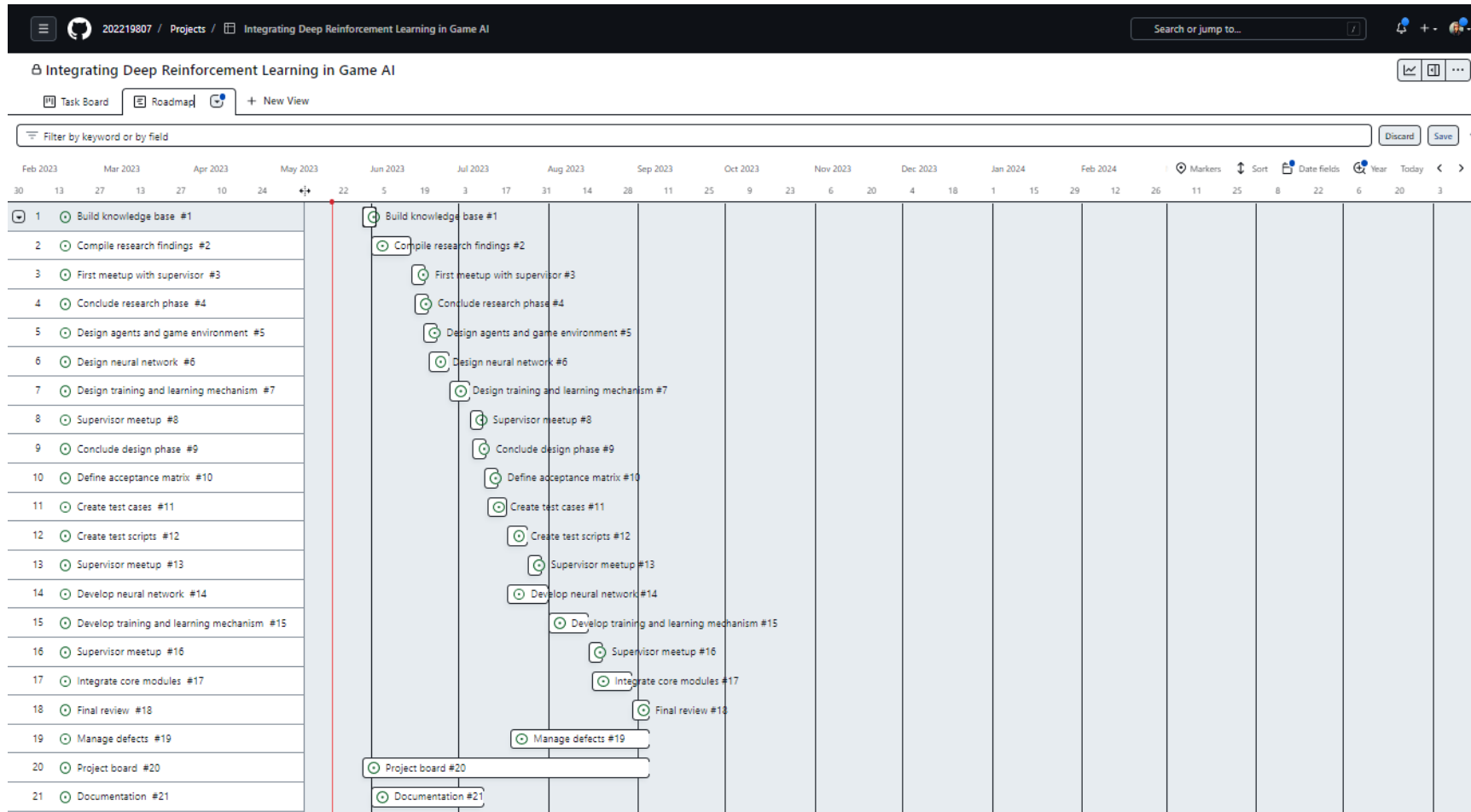
## 5 Project Management

### a. Task List

| #  | Task Name                               | Description   | Duration |
|--|---|---|----------|
| 1.   | Build knowledge base                    | Review the existing collection of papers and identify additional relevant literature sources and collect papers | 3 days   |
| 2.   | Compile research findings               | Take notes and summarize key findings from each paper   | 2 weeks  |
| 3.   | First meetup with supervisor            | Discuss findings and insights with supervisor   | 1 day    |
| 4.   | Conclude research phase                 | Organize collected papers and create a comprehensive bibliography   | 1 day    |
|  <b>Milestone 1</b>   |   |   |          |
| 5.   | Design agents and game environment      | Define agent-environment interactions and state representations   | 2 days   |
| 6.   | Design neural network                   | Design the neural network architecture  | 1 week   |
| 7.   | Design training and learning mechanism  | Determine training methodologies and reinforcement learning algorithms to be used                               | 1 week   |
| 8.   | Supervisor meetup                       | Discuss findings and insights with supervisor   | 1 day    |
| 9.   | Conclude design phase                   | Refine and finalize the architectural design based on feedback received   | 2 days   |
|  <b>Milestone 2</b> |   |   |          |
| 10.  | Acceptance matrix                       | Define acceptance criteria for relevant modules   | 1 day    |
| 11.  | Create test cases                       | Prepare efficient test case/suites  | 1 week   |
| 12.  | Create test scripts                     | Create scripts for loading various scenarios and agent settings   | 1 week   |
| 13.  | Supervisor meetup                       | Discuss progress  | 1 day    |
|  <b>Milestone 3</b> |   |   |          |
| 14.  | Develop neural network                  | Implement neural network architecture   | 2 weeks  |
| 15.  | Develop training and learning mechanism | Implement training algorithms, reward mechanisms and reinforcement learning processes                           | 2 weeks  |
| 16.  | Supervisor meetup                       | Discuss progress  | 1 day    |
|  <b>Milestone 4</b> |   |   |          |
| 17.  | Integrate core modules                  | Integrate neural network, reward mechanisms, and training processes   | 2 weeks  |

|     |                |   |              |
|-----|----------------|---|--------------|
| 18. | Manage defects | Log and track identified defects                      | Till closure |
| 19. | Documentation  | Update README.md and other design documents regularly | Till closure |
| 20. | Project board  | Update project board regularly on GitHub              | Till closure |
| 21. | Final Review   | Get final review and make relevant changes            | 1 week       |

## b. Gantt Chart



### c. Risk Analysis

| <b>Risk</b>   | <b>Severity<br/>(L/M/H)</b> | <b>Likelihood<br/>(L/M/H)</b> | <b>Significance<br/>(Sev. x Like.)</b> | <b>How to Avoid</b>   | <b>How to Recover</b>  |
|---|-----------------------------|-------------------------------|--|---|--|
| Data loss   | H                           | M                             | HM                                     | Keep Backups, Update GitHub repository regularly                                    | Reinstate from backups   |
| Loss of backups                                     | H                           | L                             | HL                                     | Multiple Backups, maintain separate branches  | Use alternate  |
| Delay in research                                   | M                           | M                             | MM                                     | Develop a clear and realistic timeline for conducting research                      | Adjust project schedule, seek assistance if needed   |
| Limited knowledge transfer from supervisor          | M                           | L                             | ML                                     | Clearly communicate expectations and requirements, ask for clarification            | Seek additional guidance or support from supervisor or experts if needed                         |
| Project scope creep                                 | H                           | M                             | HM                                     | Define and freeze project scope, maintain clear communication                       | Evaluate impact, prioritize changes, and update project plan accordingly                         |
| Insufficient computational resources                | M                           | M                             | MM                                     | Estimate resource requirements, utilize cloud computing, or optimize resource usage | Prioritize and allocate resources effectively, consider alternative solutions or scaling options |
| Technical issues with neural network implementation | H                           | M                             | HM                                     | Validate and test implementation before integration                                 | Debug and resolve technical issues, consult with experts if needed                               |

|  |   |   |    |  |  |
|--|---|---|----|--|--|
| Inadequate algorithm performance       | H | M | HM | Thoroughly evaluate and validate algorithms, conduct rigorous testing and optimization | Refine algorithms, explore alternative approaches, consult experts if needed |
| Defects affecting system functionality | M | M | MM | Conduct rigorous testing and quality assurance   | Prioritize and resolve defects, implement mitigation measures if necessary   |



## Appendix A: Deep Reinforcement Learning Methods

TABLE I  
A GENERAL REVIEW OF RECENT DRL METHODS FROM 2017 TO 2018.

| DRL Algorithms         | Main Techniques   | Networks | Category                    |
|------------------------|---|----------|-----------------------------|
| DQN [14]               | experience replay, target Q-network                                   | CNN      | value-based, off-policy     |
| Double DQN [15]        | double Q-learning   | CNN      | value-based, off-policy     |
| Dueling DQN [17]       | dueling neural network architecture                                   | CNN      | value-based, off-policy     |
| Prioritized DQN [16]   | prioritized experience replay   | CNN      | value-based, off-policy     |
| Bootstrapped DQN [51]  | combine deep exploration with DNNs                                    | CNN      | value-based, off-policy     |
| Gorila [20]            | massively distributed architecture                                    | CNN      | value-based, off-policy     |
| LS-DQN [22]            | combine least-squares updates in DRL                                  | CNN      | value-based, off-policy     |
| Averaged-DQN [23]      | averaging learned Q-values estimates                                  | CNN      | value-based, off-policy     |
| DQfD [24]              | learn from the demonstration data                                     | CNN      | value-based, off-policy     |
| DQN with Pop-Art [18]  | adaptive normalization with Pop-Art                                   | CNN      | value-based, off-policy     |
| Soft DQN [29]          | KL penalty and entropy bonus  | CNN      | value-based, off-policy     |
| DQV [25]               | training a Quality-value network                                      | CNN      | value-based, off-policy     |
| Rainbow [26]           | integrate six extensions to DQN                                       | CNN      | value-based, off-policy     |
| RUDDER [27]            | return decomposition  | CNN-LSTM | value-based, off-policy     |
| Ape-X DQfD [28]        | transformed Bellman operator, temporal consistency loss               | CNN      | value-based, off-policy     |
| C51 [30]               | distributional Bellman optimality                                     | CNN      | value-based, off-policy     |
| QR-DQN [31]            | distributional RL with Quantile regression                            | CNN      | value-based, off-policy     |
| IQN [32]               | an implicit representation of the return distribution                 | CNN      | value-based, off-policy     |
| A3C [33]               | asynchronous gradient descent   | CNN-LSTM | policy gradient, on-policy  |
| GA3C [34]              | hybrid CPU/GPU version  | CNN-LSTM | policy gradient, on-policy  |
| PPO [42]               | clipped surrogate objective, adaptive KL penalty coefficient          | CNN-LSTM | policy gradient, on-policy  |
| ACER [43]              | experience replay, truncated importance sampling                      | CNN-LSTM | policy gradient, off-policy |
| ACKTR [44]             | K-FAC with trust region   | CNN-LSTM | policy gradient, on-policy  |
| Soft Actor-Critic [52] | entropy regularization  | CNN      | policy gradient, off-policy |
| UNREAL [35]            | unsupervised auxiliary tasks  | CNN-LSTM | policy gradient, on-policy  |
| Reactor [39]           | Retrace( $\lambda$ ), $\beta$ -leave-one-out policy gradient estimate | CNN-LSTM | policy gradient, off-policy |
| PAAC [36]              | parallel framework for A3C  | CNN      | policy gradient, on-policy  |
| DDPG [45]              | DQN with deterministic policy gradient                                | CNN-LSTM | policy gradient, off-policy |
| TRPO [41]              | incorporate a KL divergence constraint                                | CNN-LSTM | policy gradient, on-policy  |
| D4PG [46]              | distributed distributional DDPG                                       | CNN      | policy gradient, on-policy  |
| PGQ [37]               | combine policy gradient and Q-learning                                | CNN      | policy gradient, off-policy |
| IMPALA [40]            | importance-weighted actor learner architecture                        | CNN-LSTM | policy gradient, on-policy  |
| FiGAR-A3C [53]         | fine grained action repetition  | CNN-LSTM | policy gradient, on-policy  |
| TreeQN/ATreeC [47]     | on-line planning, tree-structured model                               | CNN      | model-based, on-policy      |
| STRAW [48]             | macro-actions, planning strategies                                    | CNN      | model-based, on-policy      |
| World model [49]       | mixture density network, variational autoencoder                      | CNN-LSTM | model-based, on-policy      |
| MuZero [54]            | representation function, dynamics function, and prediction function   | CNN      | model-based, off-policy     |



## References

- Mnih, V. *et al.* (2013) Playing Atari with Deep Reinforcement Learning, arXiv.org. Available at: <https://arxiv.org/abs/1312.5602> (Accessed: 18 May 2023).
- van Hasselt, H., Guez, A. and Silver, D. (2015) Deep Reinforcement Learning with double Q-learning, arXiv.org. Available at: <https://arxiv.org/abs/1509.06461> (Accessed: 18 May 2023).
- Mnih, V. *et al.* (2015) Human-level control through Deep Reinforcement Learning, Nature News. Available at: <https://www.nature.com/articles/nature14236> (Accessed: 18 May 2023).
- Mnih, V. *et al.* (2016) Asynchronous methods for Deep Reinforcement Learning, arXiv.org. Available at: <https://arxiv.org/abs/1602.01783> (Accessed: 18 May 2023).
- Schulman, J. *et al.* (2017) Proximal policy optimization algorithms, arXiv.org. Available at: <https://arxiv.org/abs/1707.06347> (Accessed: 18 May 2023).
- Christiano, P. *et al.* (2017) Deep Reinforcement Learning from human preferences, arXiv.org. Available at: <https://arxiv.org/abs/1706.03741> (Accessed: 18 May 2023).
- OpenAI *et al.* (2019) Dota 2 with large scale Deep Reinforcement Learning, arXiv.org. Available at: <https://arxiv.org/abs/1912.06680> (Accessed: 18 May 2023).
- Shao, K., Zhu, Y. and Tang, Z. (2019) A survey of deep reinforcement learning in video games, arXiv.org. Available at: <https://arxiv.org/pdf/1912.10944.pdf> (Accessed: 18 May 2023).