

# CSC\_52081\_EP, Reinforcement Learning Project

Rafael Fernandes Pignoli Benzi, Mateus Henrique Galvão, Guilherme Nunes Trofino

**Abstract**—Reinforcement Learning (RL) trains autonomous agents to interact with complex environments. This study examines the CarRacing-v3 environment from Gymnasium, featuring high-dimensional observations and both discrete and continuous actions. We compare RL algorithms: DQN, SARSA, CEM, PPO, and SAC. Our experiments evaluate the effects of visual variability, action-space design, and hyperparameter tuning on performance. A baseline is established with a default agent. This work analyzes algorithmic trade-offs, offering insights into RL strategies for continuous control in visually complex settings.

## I. INTRODUCTION

Reinforcement Learning (RL) has become a powerful paradigm for developing autonomous agents that learn optimal behaviors through interactions with their environments. In this study, we employ the CarRacing-v3 environment provided by Gymnasium [1], which presents a challenging control task in a racing scenario. The environment is characterized by a high-dimensional observation space and two distinct modes for the action space. Specifically, the observation space consists of a top-down  $96 \times 96$  RGB image capturing both the car and the racetrack, thus requiring the use of deep convolutional neural networks (CNNs) for effective feature extraction.



Fig. 1. A top-down  $96 \times 96$  RGB image of the car and racetrack.

The primary objective of this project is to investigate and compare different RL policies across both discrete and continuous action modalities. For discrete action control, we implement methods such as Deep Q-Network (DQN) and SARSA. In contrast, for continuous action control, we explore approaches like the Cross-Entropy Method (CEM), and also incorporating policy gradient techniques (Proximal Policy Optimization (PPO) and Soft Actor-Critic (SAC)). This comparative analysis is driven by our interest in understanding the strengths and limitations of each method in handling complex decision spaces.

The dual nature of the action space in CarRacing-v3 presents a significant challenge. When dealing with high-dimensional visual inputs, the necessity of effective feature extraction becomes paramount. To address this, our approach includes the development of a convolutional neural network architecture tailored to process the  $96 \times 96$  RGB images,

reducing their dimensionality while preserving essential spatial features required for decision making. Additionally, transitioning between discrete and continuous representations of actions requires careful algorithmic design and parameter tuning to ensure stable learning and convergence.

While previous studies have applied various RL techniques in simulated environments, many have tended to focus on either discrete or continuous action spaces separately. In our work, we adopt a comparative approach by evaluating different agents within the same CarRacing-v3 environment. This allows us to assess the performance of each method under similar conditions, examining aspects such as learning stability, computational complexity, and overall policy effectiveness.

The code for this project is available at GitHub.

## II. BACKGROUND

### A. Environment

As previously stated, the CarRacing-v3 environment's observation space comprises a top-down  $96 \times 96$  RGB image, depicting both the vehicle and the track. The high dimensionality of this input necessitates the implementation of deep convolutional neural networks (CNNs) to facilitate effective feature extraction.

Regarding the action space, CarRacing-v3 supports both continuous and discrete control modalities. In the continuous mode, the agent outputs three real-valued commands: steering, where values range from  $-1$  (full left) to  $+1$  (full right); gas; and braking. Conversely, in the discrete mode, the action space is reduced to five actions: do nothing, steer left, steer right, gas, and brake. This duality in action representation allows for a comprehensive evaluation of various RL algorithms under different control settings.

The reward structure of the environment underscores the challenge by combining two components: a penalty of  $-0.1$  per frame and a reward of  $+\frac{1000}{N}$  for each new track tile visited, where  $N$  represents the total number of track tiles. For example, completing the race after visiting all  $N$  tiles in 732 frames, results in a reward of  $1000 - 0.1 \times 732 = 926.8$  points, as shown in [1]. This scheme incentivize the agent to balance exploration (visiting tiles) with efficiency (minimizing frame usage), aligning its learning objectives with the task's overarching goal.

### B. Discrete Action Space

DEEP Q-Network (DQN) and SARSA are powerful reinforcement learning algorithms to solve discrete action control problems, suitable to our first approach to the Car Racing environment. Both approaches are based on the Q-learning algorithm, which is a model-free reinforcement learning algorithm that aims to learn the optimal action-value function

$Q(s, a)$ , where  $s$  is the state and  $a$  is the action. However, DQN uses a deep neural network to approximate the Q-function, while SARSA uses a table to store the Q-values [4]. To account for the limitations of the SARSA approach in limited high-dimensional state spaces (as the Car Racing environment with an observation space of  $96 \times 96 \times 3$ ), we will explore a modern approach called Deep SARSA. [5]

Deep SARSA combines the on-policy nature of traditional SARSA with neural network architectures to handle large state spaces effectively. Unlike DQN's off-policy approach, Deep SARSA maintains SARSA's fundamental characteristics while scaling to complex environments.

On-policy learning methods, such as SARSA, updates the Q-values using actions selected by the current policy. (usually, the  $\epsilon$ -greedy). The  $\epsilon$ -greedy policy selects a random action (explore) with probability  $\epsilon$  and the best action (exploit) with probability  $1-\epsilon$ .

On the other hand, off-policy learning methods, such as DQN, updates the Q-values using actions selected by a different policy. (often by greedy selection).

On-policy methods tend to be more stable, but they can be less sample-efficient and must actively explore during training. In contrast, off-policy methods can learn from experiences, which can lead to better exploration and exploitation of the environment, but may require careful tuning to ensure stability.

### C. Continuous Action Space

For the second approach to the Car Racing environment, we will explore algorithms that can handle continuous action spaces. Starting from an evolutionary approach, we will test the performance of the Cross-Entropy Method (CEM) in the Car Racing environment.

The Cross-Entropy Method is a simple optimization algorithm that iteratively samples policies from a Gaussian distribution and updates the distribution parameters to maximize the expected return. It generates multiple candidate solutions (policies), ranks them based on performance, and updates the policy distribution using the best-performing candidates. It is efficient in high-dimensional spaces and can discover complex policies through population diversity. It follows the steps:

- 1) Sample  $N$  policies from a Gaussian distribution.
- 2) Evaluate the policies in the environment.
- 3) Select the top  $M$  policies.
- 4) Update the distribution parameters to fit the selected policies.
- 5) Repeat until convergence.

Secondly, we will explore and compare two policy-based methods designed for continuous control: Proximal Policy Optimization (PPO) [6] and Soft Actor-Critic (SAC) [7] in the Car Racing environment.

PPO is primarily an on-policy gradient method because it collects trajectories using the current policy, and updates it using only the data from the most recent rollout (episode or batch). Once data is used for training, it is discarded (unlike fully off-policy algorithms like DQN that store and reuse old experiences). It uses an actor-critic architecture to learn a parameterized policy and, even though PPO directly updates

the policy, it still needs a value function (critic network) for advantage estimation (GAE - Generalized Advantage Estimation). The advantage function helps reduce variance in the policy updates.

The classic PPO algorithm follows the steps:

- 1) Collect trajectories using the current policy in a buffer.
- 2) Compute the advantage function using the critic network.
- 3) Update the policy using the advantage function and the policy gradient.
- 4) Repeat until convergence.

SAC is an off-policy (meaning it reuses past experiences for learning) actor-critic method that uses the maximum entropy framework to encourage exploration and improve sample efficiency. It learns a stochastic policy that maximizes the expected return while maximizing the entropy of the policy (avoiding premature convergence to suboptimal policies). It follows the steps:

- 1) Collect trajectories using the current policy.
- 2) Compute the advantage function using the critic network.
- 3) Update the policy using the advantage function and the policy gradient.
- 4) Update the critic network using the temporal difference error.
- 5) Repeat until convergence.

### D. Objective

In this study, we will compare the performance of Deep SARSA and DQN in the Car Racing environment to understand the trade-offs between on-policy and off-policy learning methods in the context of discrete action space. Conversely, we also aim to compare the performance of CEM, PPO, and SAC to understand the trade-offs between evolutionary algorithm and policy-based methods in continuous action spaces.

## III. METHODOLOGY / APPROACH

### A. Environment and Agent Implementation

This study employs the CarRacing-v3 environment from Gymnasium [1], a challenging benchmark characterized by high-dimensional visual observations and multiple action modalities. In order to transform the observation space (a  $3 \times 96 \times 96$  picture frame) in the action space supported by the environment, we implemented a Convolutional Neural Network (CNN) inspired by [8] with the following characteristics:

- **Image:** the frames are transformed from an RGB state representation to grayscale, and the resolution is down-scaled to  $84 \times 84$  pixels. These steps facilitate the network computations, making them less complex. Additionally, to allow the network to accurately track the car's dynamics and improve the agent's ability to adjust actions based on the current state, each frame is stacked with the previous three frames. As a result, every state is represented by a sequence of four consecutive grayscale frames at  $84 \times 84$  resolution (see Figure 2).
- **Input Layer:** 4-channel  $84 \times 84$  images.

- **First Convolutional Layer:** 16 filters, 8x8 kernel, stride 4, ReLU activation. After applying the convolution, the output have dimensions 20x20.
- **Second Convolutional Layer:** 32 filters, 4x4 kernel, stride 2, ReLU activation. After applying the convolution, the output have dimensions 9x9.
- **Flatten Layer:** The output from the second convolutional layer is flattened into a 1D vector. The size of this vector is  $32 * 9 * 9 = 2592$ .
- **First Fully Connected (FC) Layer:** 256 neurons, ReLU activation.
- **Second Fully Connected (FC) Layer:** The last fully connected layer defines the action space of the environment. For discrete action-space trials (DQN and Deep SARSA), this layer consists of five neurons with a linear activation function. In contrast, for continuous action-space tests (PPO, SAC, and CEM), the output layer comprises three neurons with different activation functions: *tanh* for the steering action (ranging from [-1, 1]) and *sigmoid* for the gas and brake actions (both ranging from [0, 1]).

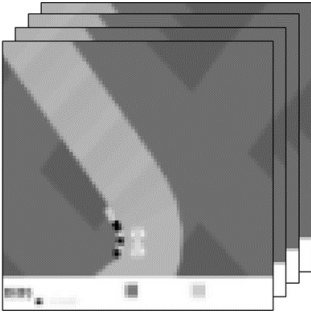


Fig. 2. Input observation space for the CNN, composed by four grayscale consecutive frames from the car racing environment. Reference: [8]

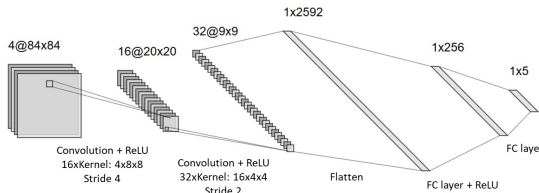


Fig. 3. Convolutional Neural Network (CNN) architecture for processing the observation from the car racing environment. Reference: [8]

### B. Replay Buffer

A key innovation introduced in the original DQN paper is the concept of experience replay. This technique involves storing experiences in a replay memory buffer, allowing the agent to break the temporal dependencies between consecutive experiences. During training, random minibatches are sampled from this buffer, which improves the stability of the learning process.

### C. Algorithms Design

1) *DEEP Q-Network (DQN)*: The model is trained by optimizing a loss function based on the temporal difference error between predicted and target Q-values.

For the DQN, the target network consists in maintaining two separate networks: the main (or online) network, which is used for learning and selecting actions, and the target network, which is updated less frequently. The target network is a copy of the online network, and its parameters are periodically updated by copying the parameters of the online network to it. This approach helps stabilize the learning process by providing a fixed target for the updates, preventing oscillations and divergence in the Q-value estimates.

Training is conducted on minibatches of state-action-reward-next state sequences sampled from the replay buffer. The update rule of the algorithm is given by:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (1)$$

where  $\max_{a'} Q(s', a')$  is the maximum Q-value at the next state  $s'$  (greedy selection).

2) *DEEP SARSA*: For the SARSA algorithm, learning is performed using a single Q-network instead of maintaining a separate target network. The algorithm follows an on-policy approach, where the agent updates its Q-values based on the actions it actually takes, rather than using a target derived from the maximum possible future reward. This ensures that updates remain consistent with the agent's current policy.

Training is conducted on minibatches of state-action-reward-next state-next action sequences sampled from the replay buffer. The update rule of the algorithm is given by:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', \pi(s')) - Q(s, a)] \quad (2)$$

where  $Q(s', \pi(s'))$  corresponds to the Q-value of the next state-action pair, following the agent's current policy ( $\epsilon$ -greedy). This approach ensures that the learning process accounts for the agent's actual behavior, rather than assuming a greedy action selection at every step.  $\epsilon$  is then decayed for better balance between prioritizing exploration at the beginning, and more exploitation after.

3) *Proximal Policy Optimization (PPO)*: The training for the PPO is described as follows:

- First, we collect trajectories from the environment:

$$\mathcal{D} = \{(s_t, a_t, r_t, s_{t+1})\}_{t=0}^T$$

And the experience tuple is stores in the replay buffer, until termination.

- Once the buffer has at least `batch_size` samples, PPO updates are performed over multiple epochs, and the current replay buffer is cleared.
- The update occurs sampling a batch of transitions from the buffer, and the target values are computed using the Bellman equation:

$$V_{target} = r + \gamma V(s) \quad (3)$$

And the value loss is then calculated using Mean Squared Error (MSE).

- The advantage function is computed using Generalized Advantage Estimation (GAE):

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

And the advantages are normalized for stable training

- Policy update computes the probability ratio between the new and old policies (using PPO clipped surrogate objective after):

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$$

To encourage exploration, an entropy bonus is added. The policy loss is backpropagated, and the optimizer updates the policy network.

- Repeat for multiple Epochs.

Unlike other algorithms, we used two Neural Networks (NNs): one for the policy and another for the value function. The Policy Network follows the CNN architecture described in Section III-A, with its output representing the mean and standard deviation of the action distribution. Within the PPO class, the update function constructs a normal distribution based on these parameters, from which an action is sampled. The action is then clipped to ensure it stays within the boundaries of the continuous action space.

The Value Network, in contrast, uses a simplified version of the CNN, as it does not require extensive feature extraction. Its output is a single scalar representing the estimated state value.

4) *Cross-Entropy Method (CEM)*: The Cross-Entropy Method (CEM) is a population-based optimization algorithm used to iteratively refine a distribution over the parameter space, converging towards an optimal policy. It operates by maintaining a probability distribution over the policy parameters and updating it based on the performance of sampled candidates.

At each iteration, a set of candidate solutions is sampled from a multivariate normal distribution:

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \forall i \in \{1, \dots, m\} \quad (4)$$

where  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  represent the mean vector and covariance matrix, respectively, defining the current search distribution. Each candidate solution  $\mathbf{x}_i$  is evaluated using the objective function, which corresponds to the negative cumulative reward in reinforcement learning. The best-performing  $m_{\text{elite}}$  candidates (elite set) are selected to update the parameters of the search distribution:

$$\boldsymbol{\mu} \leftarrow \frac{1}{m_{\text{elite}}} \sum_{j=1}^{m_{\text{elite}}} \mathbf{x}_j \quad (5)$$

$$\boldsymbol{\Sigma} \leftarrow \frac{1}{m_{\text{elite}}} \sum_{j=1}^{m_{\text{elite}}} (\mathbf{x}_j - \boldsymbol{\mu})(\mathbf{x}_j - \boldsymbol{\mu})^T \quad (6)$$

This process iterates until convergence, refining the distribution towards regions of higher reward. Unlike gradient-based methods, CEM does not require differentiability of the

objective function, making it particularly suitable for policy search in high-dimensional and non-differentiable reinforcement learning problems.

As mentioned for PPO, two NNs are also used in this algorithm. The first one processes the high-dimensional visual observations produced by the CarRacing environment. It takes as input the preprocessed frames (in this case, four stacked grayscale images of size  $84 \times 84$ ) and extracts spatial features by applying a series of convolutional filters and nonlinear activations. The second one is a policy that the agent can use to interact with the environment, handling the necessary preprocessing of the input observation. Moreover, this policy provides utility methods to extract all network parameters into a single flattened vector and to update the network parameters from such a vector.

#### 5) SAC:

### D. Comparison and Evaluation Methods

1) *Action-Space Design*: We compare the efficiency of discrete versus continuous action spaces by analyzing convergence rates, stability, and final policy performance across different action representations.

2) *Hyperparameter Sensitivity Analysis*: We investigate the stability and convergence dynamics under different step-size configurations by adapting the learning rate.

3) *Performance Metrics*: Our evaluation framework incorporates several key indicators to comprehensively assess the performance of the implemented methodologies. Firstly, cumulative rewards are used to quantify the overall policy efficiency and long-term reward accumulation. Secondly, convergence speed is measured by the number of training episodes required to reach a stable performance threshold. Thirdly, sample efficiency is evaluated by assessing the learning progress per unit of environmental interaction, which helps determine the algorithmic effectiveness. Lastly, robustness evaluation is conducted through perturbation tests under varying environmental conditions to gauge the model's adaptability and resilience.

### E. Visual Analysis and Interpretation

To further refine our understanding of agent behavior and policy effectiveness, we employ several visualization techniques:

1) *Learning Curves*: We analyze episodic reward trends over the course of training to identify key inflection points and phases in the learning process. This involves plotting the cumulative rewards obtained in each episode, which helps in diagnosing the learning stability, convergence behavior, and potential over-fitting or under-fitting issues. By examining these curves, we can infer the efficiency of the learning algorithm and the impact of different hyperparameter settings.

2) *Final Policy Comparisons*: We conduct a comparative analysis of the final policies learned under different training configurations. This involves contrasting the strategies and behaviors exhibited by the agent after the training process is complete. In order to compare algorithm performance under

the same conditions of an environment (which is randomly generated), we fixed a seed variable upon its generation to maintain the main characteristics.

By systematically comparing these policies, we can infer the impact of various algorithmic choices, such as different exploration strategies, reward structures, or network architectures, on the overall performance and robustness of the learned policies. This comparative analysis provides a deeper understanding of the strengths and limitations of each approach, guiding future improvements and refinements.

#### F. Reproducibility

All implementations utilize OpenAI Gymnasium, PyTorch/TensorFlow for training, and Stable-Baselines3 for baseline comparisons. The full codebase is available at GitHub to facilitate reproducibility.

### IV. RESULTS AND DISCUSSION

#### A. DQN and SARSA Comparison

The evolution curve of the models can be seen in Figures 4 and 5, respectively. The blue curve shows the episodic rewards, while the orange line represents the moving average, providing a clearer view of the overall learning trend.

The DQN was trained over 2300 episodes. The reward curve shows an upward trend, indicating that the DQN is learning and improving its policy over time. The reward stabilizes around  $820 \pm 171.98$  by episode 2364, suggesting that the model has likely converged to a near-optimal policy, and the steady improvement suggests that the learning rate is likely well-tuned. Given that the maximum reward reached in the environment is 926.8, the results are quite satisfactory, but the convergence takes more time to saddle.

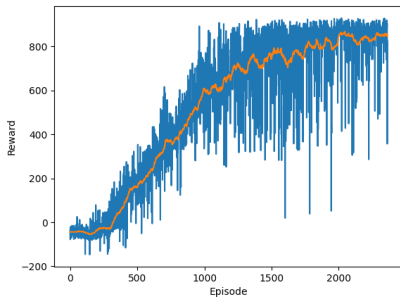


Fig. 4. Evolution curve showing the rewards over episodes - DQN for 2300 episodes.

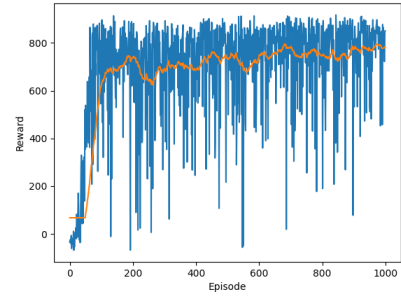


Fig. 5. Evolution curve showing the rewards over episodes - Deep SARSA for 1000 episodes.

For the SARSA model, 1000 episodes were simulated. The agent starts with low rewards in the early episodes, as expected, while it explores the environment. Rapid improvement occurs within the first 100-200 episodes, showing that the agent is learning a useful policy. After this phase, rewards continue to increase but with fluctuations, which suggest that the learning rate may be moderately high, allowing for quick learning, but eventually stabilizing around  $784 \pm 126.7$ .

Overall, this proved to be a good training and the final performance is strong, indicating that the policy has reached a near-optimal state, while also fast and sample-efficient, followed by stable convergence.

In order to compare their performance under the same conditions (same environment), we performed five simulations for over 25 seeds. The graph for comparison can be seen in Figure ??.

In summary, Deep SARSA prioritizes fast adaptation but requires more time to stabilize, while DQN takes longer to learn but ultimately achieves higher performance and stability.

#### B. CEM

Figure 6 shows the evolution of the training rewards over iterations for the CEM algorithm. As expected from the minimization strategy applied in our implementation, the reward value starts at a high level and decreases steadily, eventually stabilizing between 5 and 10. This behavior indicates that the algorithm is correctly minimizing the cost, since our objective function returns the negative reward.

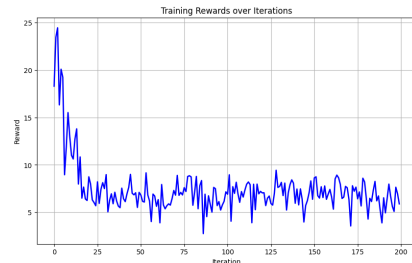


Fig. 6. Training plot showing the evolution of the reward over iterations.

However, the desired outcome was to achieve a much lower (more negative) reward value as training progressed, which would have translated into better performance by the agent when navigating the track. Unfortunately, the final model



did not produce the expected results in the environment, as demonstrated by the video reproduction of the car on the track.

This may be the result of some problems with the model. One possibility is that the reward structure, when negated for minimization, may be inadvertently favoring a strategy that minimizes penalization rather than encouraging forward progress. Additionally, the parameters of the CNN or the action scaling (especially for the steering, gas, and brake outputs) might not be appropriately tuned, causing the policy to output nearly constant or ineffective actions. Overall, these factors combined may result in the agent exploiting a trivial strategy that minimizes the cost, rather than discovering a policy that drives effectively along the track. But, despite this problem, the training plot confirms that the minimization mechanism is operating as designed.

### C. PPO and SAC

For both algorithms, over approximately 1000 episodes, we did not observe a training that converged to an optimal policy. Figure 7 shows the poor performance of the model, that can be explained by higher complexity of the continuous action-space possibilities.

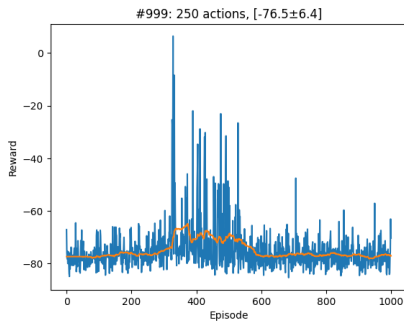


Fig. 7. Evolution curve of the rewards per episode for the PPO algorithm.

We then evaluated both performances of PPO and SAC under the same conditions (same environment), which can be seen in Figure ??.

Many of the reasons that could explain the poor performance of the algorithms are that continuous action spaces introduce infinitely many possible actions, making exploration harder compared to discrete spaces. If the policy does not explore effectively, it may fail to find optimal actions. PPO and SAC are highly sensitive to hyperparameters like learning rate, entropy regularization (SAC), and clipping ratio (PPO). Inappropriate values can lead to either premature convergence to suboptimal policies or unstable training. Additionally, the neural networks may struggle to approximate the optimal policy in complex continuous spaces, especially if architecture, activation functions, or regularization are not well chosen, and there are many possibilities.

Overall, we don't have enough evidence to claim that the algorithm completely failed, because analysing the videos after 1000 episodes, it does seem to make reasonable choices, but is not fast enough. Maybe with more training it would start to perform better, but we did not have the time to train the algorithm (usually it requires more than 10 hours).

## V. CONCLUSIONS

### A. Discrete Action-Space

Deep SARSA and DQN are both powerful reinforcement learning algorithms for discrete action-space in the Car Racing environment, but they exhibit different learning dynamics and convergence properties.

Deep SARSA learns faster in the early stages of training due to its on-policy nature, which allows it to adjust its policy dynamically based on the latest experiences. This results in quicker adaptation and initial reward improvements. However, because it updates its policy with actions influenced by exploration, it tends to be more unstable and takes longer to fully converge.

On the other hand, DQN follows an off-policy approach, learning from past experiences stored in a replay buffer. This makes training more stable and helps DQN discover a more optimal policy in the long run. However, DQN generally takes more time to reach this optimal policy due to its need for extensive experience replay and value function approximation.

### B. Continuous Action-Space

## REFERENCES

- [1] Gymnasium. *Car Racing environment*, [https://gymnasium.farama.org/environments/box2d/car\\_racing/](https://gymnasium.farama.org/environments/box2d/car_racing/), accessed 04 March 2025.
- [2] Sutton and Barto. *Reinforcement Learning*, MIT Press, 2020.
- [3] Read. Lecture IV - Reinforcement Learning I. In *INF581 Advanced Machine Learning and Autonomous Agents*, 2024.
- [4] AI Mind. Popular Reinforcement Learning algorithms and their implementation, 2023 <https://pub.aimind.so/popular-reinforcement-learning-algorithms-and-their-implementation-7adf0e092464>
- [5] van Seijen et al. Deep SARSA: A Novel Approach to Deep Reinforcement Learning, 2014.
- [6] Schulman et al. Proximal Policy Optimization Algorithms, 2017.
- [7] Haarnoja et al. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, 2018.
- [8] GitHub - wiitt/DQN-Car-Racing: Implementation of DQN and DDQN algorithms for Playing Car Racing Game <https://github.com/wiitt/DQN-Car-Racing> GitHub, 2024.

## APPENDIX