



Cairo University

Faculty of Engineering

Computer Engineering Department

## CMPS458 Reinforcement Learning - Assignment 2

Team Name/Number: Team 1

First member name: Mariam Mahrous 1210301

Second member name: Menna Salah 1210032

Third member name: Farida Ahmed 1210276

*Supervisor:* Ayman AboElhassan

November 13, 2025

# Deliverables

Repo link: <https://github.com/Mennasalah140/Reinforcement-Learning>

Video record link:

<https://drive.google.com/drive/folders/1xS7qhv08kY4WoeNomgsLMmIuZCMSeY8V?usp=sharing>

# Discussion

## 0.1 Question Answers

### Q1 — Difference Between DQN and DDQN in Terms of Training Time and Performance

Deep Q-Network (DQN) and Double Deep Q-Network (DDQN) are both value-based methods that use neural networks to approximate the Q-function. The main distinction lies in how each method estimates the target Q-value during learning.

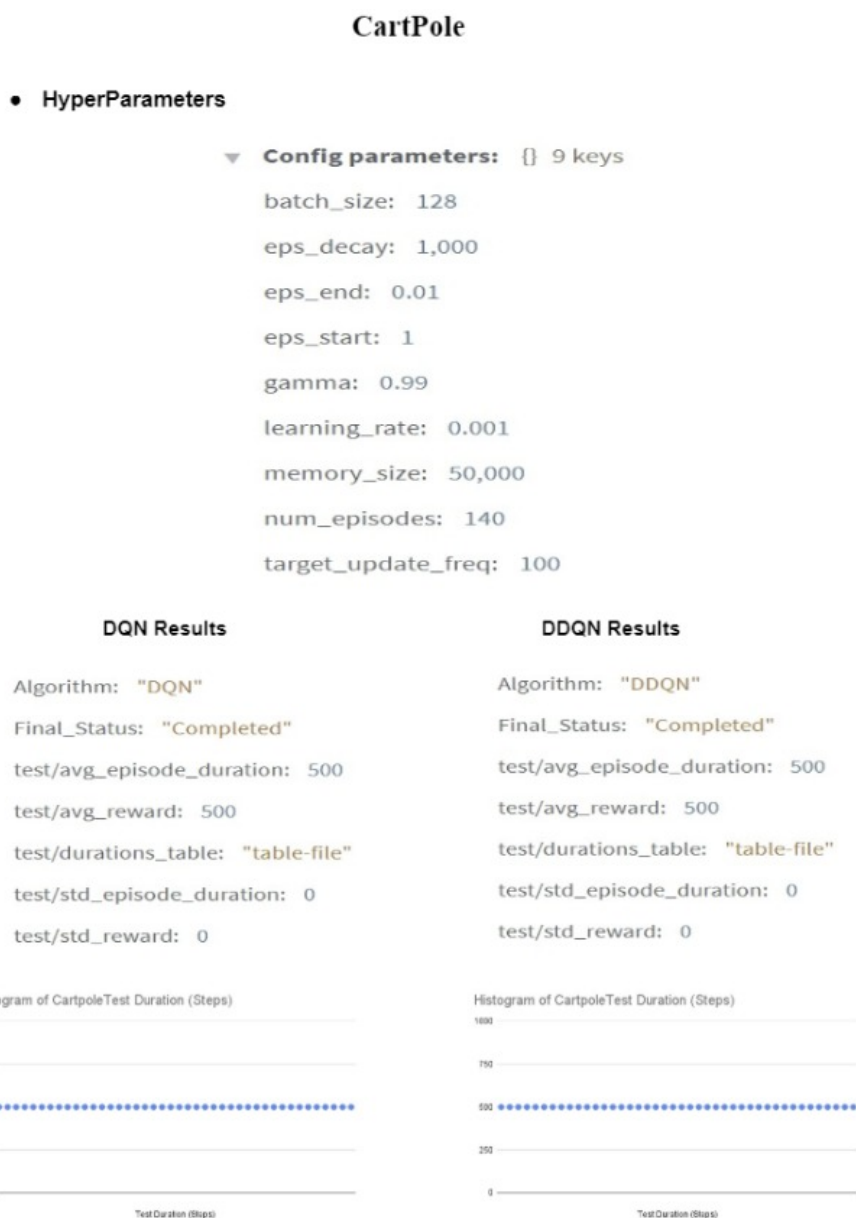
Table 1: Comparison between DQN and DDQN

Aspect	DQN	DDQN
Training Time	Slightly faster due to fewer updates	Slightly slower, but more consistent performance
Performance	Can reach high rewards quickly but unstable	Achieves higher long-term rewards and generalization

**Q2 — How stable are the trained agents? Show with test episode duration figures.**

### CartPole-v1

Figure 1: CartPole-v1



Acrobot-v1

Figure 2: Acrobot-v1

Acrobot

HyperParameters

▼ Config parameters: {} 9 keys

batch\_size: 128

eps\_decay: 1,000

eps\_end: 0.01

eps\_start: 1

gamma: 0.999

learning\_rate: 0.0002

memory\_size: 50,000

num\_episodes: 700

target\_update\_freq: 500

DQN Results

▼ Summary metrics: {} 9 keys

Algorithm: "DQN"

Final\_Status: "Completed"

test/avg\_episode\_duration: 97.89

test/durations\_table: "table-file"

test/std\_episode\_duration: 54.73972871690177

DDQN Results

▼ Summary metrics: {} 9 keys

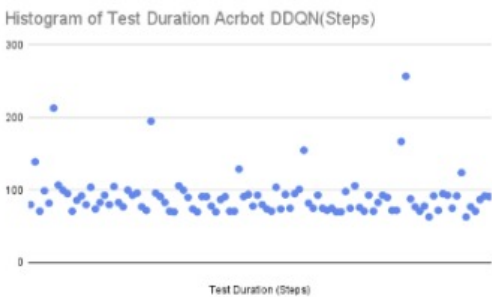
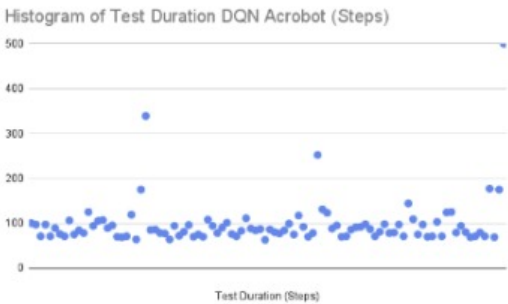
Algorithm: "DDQN"

Final\_Status: "Completed"

test/avg\_episode\_duration: 90.93

test/durations\_table: "table-file"

test/std\_episode\_duration: 29.07997764785936



MountainCar-v0

Figure 3: MountainCar-v0

MountainCar

• HyperParameters

▼ Config parameters: {} 9 keys

batch\_size: 128

eps\_decay: 20,000

eps\_end: 0.05

eps\_start: 1

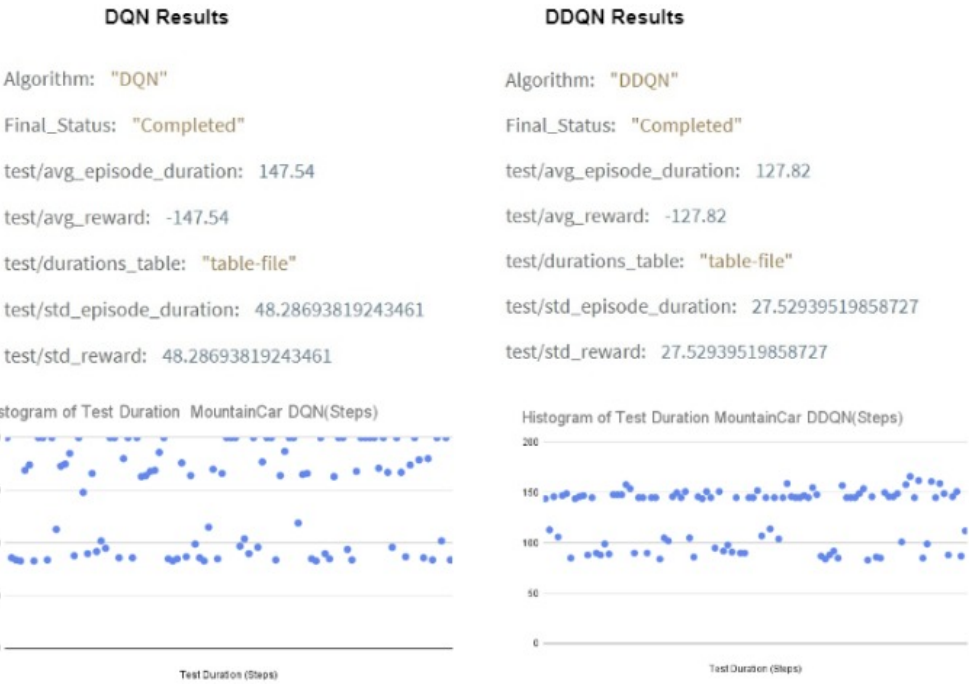
gamma: 0.9999

learning\_rate: 0.001

memory\_size: 2,000,000

num\_episodes: 1,500

target\_update\_freq: 500



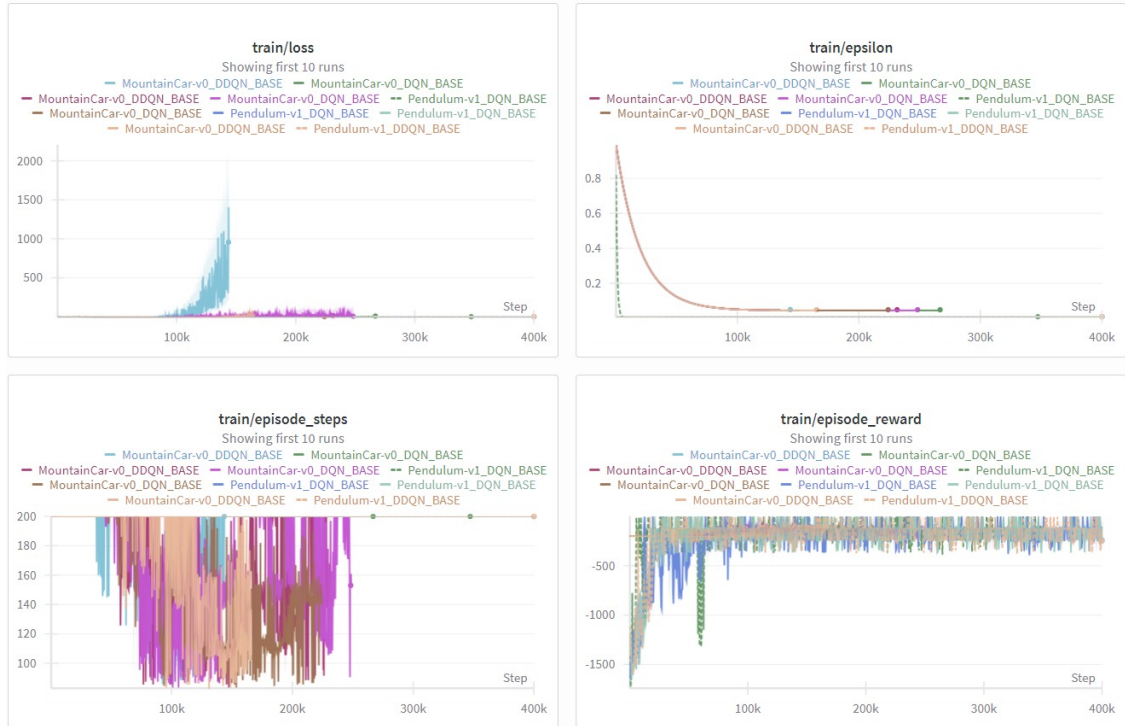
Pendulum-v1

Figure 4: Pendulum-v1



## All Training Runs

Figure 5: All Training Runs





### Q3 — Effect of Each Hyperparameter on RL Training and Performance

Table 2: Effect of Hyperparameters on DQN/DDQN Training

Hyperparameter	Effect on Training	Typical Observation
Discount Factor ( $\gamma$ )	Balances importance of future vs. immediate rewards	High $\gamma$ (0.99) leads to stable, long-term policies; low $\gamma$ causes short-sighted actions
Epsilon Decay Rate	Controls the speed of exploration decay	Fast decay $\rightarrow$ exploitation; slow decay $\rightarrow$ better exploration and stability
Learning Rate ( $\alpha$ )	Determines how quickly weights are updated	High $\alpha$ may cause divergence; low $\alpha$ slows convergence
Replay Memory Size	Stores past experiences for sampling	Large memory improves diversity and stability; small memory repeats transitions
Batch Size	Number of samples per gradient update	Small batch $\rightarrow$ noisy updates; large batch $\rightarrow$ smoother but slower learning

### Q4 — How Well DQN/DDQN Solve Each Classical Environment

For CartPole-v1, Acrobot-v1, and Mountain Car-v0:

- DQN/DDQN are **extremely well-suited** for these.
- **Reason:** These environments all have **discrete action spaces** (e.g., "push left" or "push right"). DQN is fundamentally designed to handle this. It works by having the neural network output one Q-value for each possible discrete action, and the agent simply picks the action with the highest Q-value.

For Pendulum-v1:

- DQN/DDQN are **fundamentally unsuited** for this environment in their standard form.
- **Reason:** Pendulum has a **continuous action space**. The action isn't "left or right," but a specific torque value (e.g., -1.56, 0.23, 1.88, etc.). There are infinitely many possible actions.
- A standard DQN network cannot output an infinite number of Q-values. Trying to use DQN here would force you to *discretize* the action space (e.g., treat it as just "hard left," "nothing," "hard right"), which is an approximation that often performs very poorly.