# INM707:
# Reinforcement Learning

Víctor Abia Alonso & Julian Jimenez Nimmo

student IDs: 230061642 & 230066319

May 12, 2024

Repository

# Contents

# 1 Environment and Problem

The environment to be solved by the Q-learning agent [1] is a variation of a maze layout. In a traditional maze environment, there are only walls and paths. In our environment, walls now represent fire, there exists a special state granting an additional reward and an exit state granting a final reward which terminates the game. We wish to model a situation where the goal of a robot is to save a person and reach the exit during a house fire in the shortest time possible.
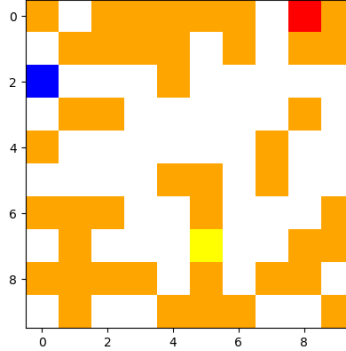


**Figure 1:** Fire maze simulation: the starting state is the blue cell, the person is located in the yellow cell,, and the exit is the red state, with fires represented by orange cells.

The maze is represented as a 10x10 grid. The starting state is cell (2,0), the person is located at cell (5,5) and the exit is located at cell (8,0). Other variations of the maze can be constructed, both through randomization or manually changing what is present at each cell. Our aim is to construct an agent which prioritizes retrieving the person and then reaching the exit rather than only heading towards the exit.

# 2 Transition and reward functions.

In reinforcement learning, an agent learns through interactions with an environment. The environment defines which actions are allowed for each state and provides a reward signal for each action, which the agent uses to learn. For our problem, we encapsulate rewards in a 3 dimensional matrix of shape (10, 10, 2) where the first two dimensions represent the position of the robot within the grid (xy coordinates) and the third dimension is a binary variable representing whether the person has been collected or not. Therefore, there are 200 possible states [1] the agent can find themselves in. For each state, in general, there are four available actions: [up, down, left, right], which make the agent transition from its current cell to the corresponding adjacent one. For edge cells and corners, out-of-bound actions make the agent stay in the grid. Also, if the person hasn't been collected yet, the actions that lead to cell (7, 5) change the last dimension of the matrix from 0 to 1. These descriptions are implemented through conditional statements in the **transition function** of the environment (`.transition_R()` method in the `maze_env.py` file). This environment is deterministic as every action unequivocally leads to a specific state at any state and time.

Additionally, the environment also provides a reward for every state and action which, as the transition function is deterministic, is in fact equivalent to receiving a reward for every next state. These rewards are encapsulated in the reward matrix R, which is a 4 dimensional array of size (10, 10, 2, 4) providing a reward for every action (4 possibilities) taken in each state (10x10x2). The unacceptable actions are conceived as NaN -Not a Number- values on the matrix. The Q matrix the agents uses to value a state-action pair has the same dimensions as this R matrix.

In order to explore different environment structures that would lead the agent to learn to collect the person and then reach the exit, we consider two different environments each with different reward and transition structures. Both structures summarize the objective we defined for the agent, and vary in the the implications of getting into a fire cell.

---

[1]In fact, there are just 199 because the robot cannot be in the location of the person unless it has been collected.

| Action | Limited movement | Terminal movement |
|---|---|---|
| Getting to the exit. | 1* | 1* |
| Collecting the person | 10 | 10 |
| Taking a step to a non-fire cell | -0.005 | -0.005 |
| Taking a step into a fire cell | None | -1* |

**Table 1:** Description of rewards in the two types of environments.

The asterisk (*) on the table means that the action terminates the episode. The choice of rewards for collecting the person and reaching the exit prioritizes collecting the person and provides extra reward for the exit. Safe steps are still characterized by some small negative reward in order to encourage the robot to finish the episodes quicker, which implicitly incentives them to get to the exit despite its relatively small reward. For actions that lead to the fire, the *limited* environment doesn't allow them at all -NaN value in the R matrix- while in the *terminal* environment it yields a significant negative reward and terminates the episode. These two ways of conceptualizing problematic states in environments are common in games, for example, either by giving negative reward (subtracting life points) or by directly prohibiting the dangerous actions in the first place. Therefore, both environments encourage the desired task of collecting the person and then exiting.

# 3 Q-learning parameters.

Q learning formula: $Q(s,a) \leftarrow Q(s,a) + \alpha \left[ r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$

The Bellman equation has two parameters alpha ($\alpha$) and gamma ($\gamma$), which are not defined by the environment, and the choice of these parameters can thus influence the performance and effectiveness of the agent. While finding the optimal parameters, it is important to show the evolution of rewards over time, as it showcases how significant each parameter was for convergence. The number of steps per episode is not an appropriate metric as terminating early, either by getting to the target without the coin or by dying in the fire (in the case of terminal environments) does not reflect the optimal performance.

The $\gamma$ parameter is the discount factor and reflects how much the agent values next step rewards compared to immediate rewards; it takes values between 0 and 1. It applies a geometric accumulation to future rewards; for instance, $\gamma = 0.1$ diminishes the value of rewards two steps away to just 0.01. Low $\gamma$ values render the agent myopic, especially in environments with rapid state changes where the value of future rewards decreases exponentially. Conversely, setting $\gamma$ to 1 is typically avoided because it theoretically eliminates the urgency for the agent to acquire rewards promptly. In order to find the optimal $\gamma$ for our environment we try different values and analyse convergence and rewards acrosss episodes.



**Figure 2:** Comparison of different gamma values. 200 episodes with softmax policy with $T_{min} = 0.001$, $T_0 = 50$ and $\lambda = 0.6$ in a limited reward environment. $\alpha = 1$.

Low values of $\gamma$, like 0.4, does not allow the agent to retrieve the person, and thus not obtain the optimal

reward. It considers this strategy too time consuming, preferring the immediate reward of the exit. Higher values of $\gamma$ seem to all converge to the optimal path, with 0.85 converging the fastest in around 60 episodes.

The $\alpha$ parameter is the learning rate of the update of the Q matrix at each step taking values between 0 and 1. An $\alpha$ value close to 1 means that the agent relies mostly on recent information (i.e. the last updates of that state-action pair) which may lead to unstable learning by underutilizing past information, although training can be quicker. An $\alpha$ close to 0 means Q values are more conservative, leading to more stable, but slower training as it incorporates information more gradually. In order to find the optimal value of $\alpha$, different values should be tried to analyse convergence.
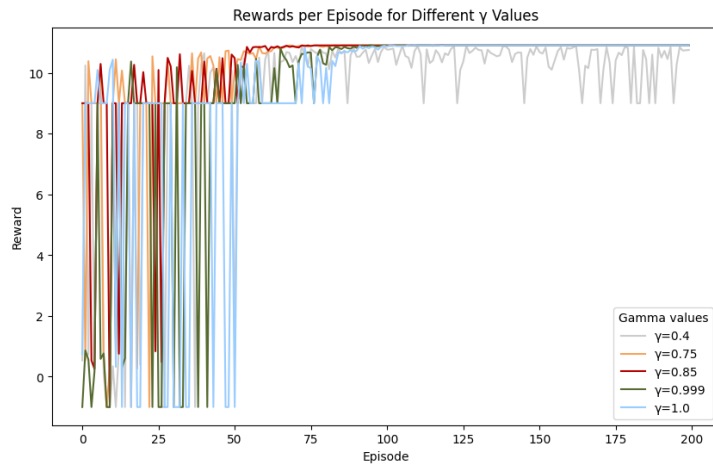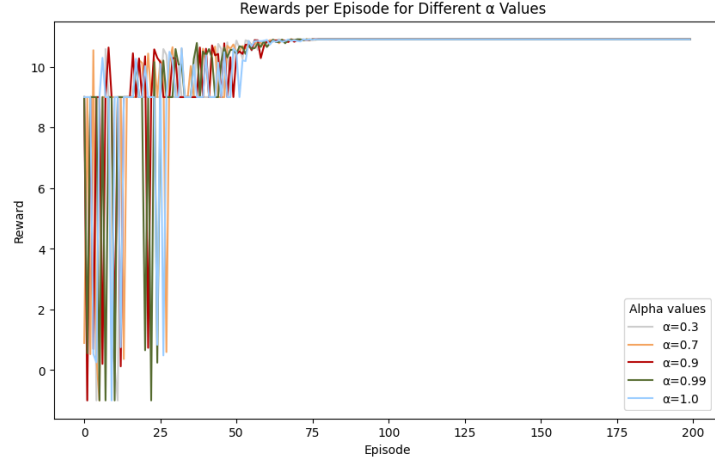


**Figure 3:** Comparison of different alpha values. 200 episodes with softmax policy with $T_{min} = 0.001$, $T_0 = 50$ and $\lambda = 0.6$ in a limited reward environment. $\gamma = 0.85$.

In this case, it seems like every value of $\alpha$ greater than 0.3 converges quite similarly, meaning that the dynamics of the environment in this specific setup seem to be indifferent to the update rate as learning will happen naturally over the epochs. However, this is not a guarantee that training is indifferent to $\alpha$ in every setup and its performance should be analysed case by case. We will choose $\alpha = 0.9$ from now on.

Note that almost from the beginning the agent is already getting high rewards, implying that the agent got the person. This is because it is the limited movement environment in which an agent never dies with the fire and it has 200 steps per episode to learn where the person and the exit are, which allows fast training (see more on the reward section). These results encouraged us to try the terminal reward structure, in which stepping into a fire kills you, and to tune the parameters of policies and compare them in this harder variant of the environment.

## 4 Different policies.

In reinforcement learning, agents choose which actions to take using a policy, which is a function mapping states to actions. Agents already have an internal map of the environment, in this case the Q matrix, which gives to every action in each state a value which aligns with the expected reward. As such, policies are a way of incorporating this knowledge from the Q-learning function into the actions the agents take. A key aspect is that the Q matrix is suboptimal at the beginning and requires exploration of the environment to be refined. Following the greedy policy of taking the action that maximizes rewards at the beginning may not give the best reward overall. This motivates the pursuit of an optimal exploration-exploitation ratio in which policies allow for sufficient exploration of the environment before focusing on choosing actions that maximize the expected reward. The two policies considered, epsilon-greedy and softmax, will encapsulate this trade-off with an exponential decay parameter $\lambda$ for which different values are considered.

The epsilon greedy policy chooses a random action with probability epsilon ($\epsilon$), and follows the greedy policy (choosing the action with the highest expected reward according to the Q function) with probability 1 - $\epsilon$. Thus, $\epsilon$ encapsulates the exploration and should be higher at the beginning and lower at the end. We choose $\epsilon = 1$ at the beginning. For the minimum $\epsilon$ (at the end) we first consider $\epsilon = 0.05$. However, this lead to degrading performance as even when the optimal path was identified, it "cannot" follow it as

sometimes it chooses actions randomly. In this case, the optimal path is 20 steps long, meaning that the agent has 64% chance ($= 1 - 0.95^{20}$) of taking a suboptimal step at some point in the path, which in the terminal environment frequently means death. Given the simple nature of our environment, and after monitoring training, we chose the minimum epsilon to be 0. The exponential decay rate describes the value of $\epsilon$ in each epoch, which decays from $\epsilon_{max} = 1$ in the first epoch to $\epsilon_{min} = 0$ after infinite epochs following the formula: $\epsilon_n = \epsilon_{min} - (\epsilon_{max} - \epsilon_{min}) \times e^{-\lambda n}$ where $n$ indicates the epoch.
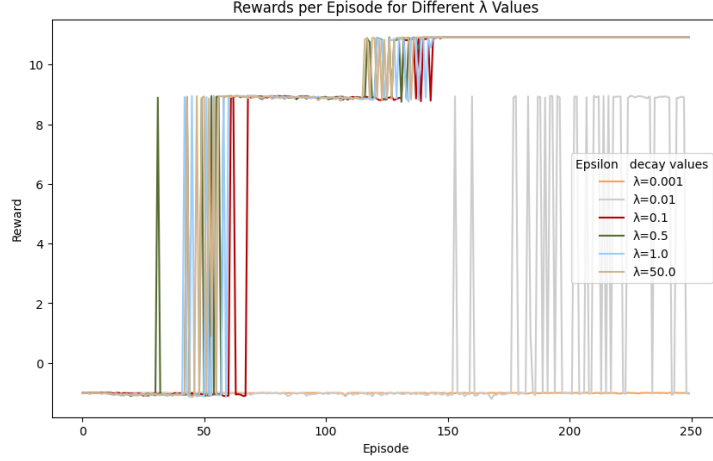


**Figure 4:** Comparison of different epsilon decay values on the epsilon greedy policy performance. 250 episodes. Terminal rewards. $\gamma = 0.85$. $\alpha = 1$. $\epsilon_{min} = 0$. $\epsilon_{max} = 1$.

The parameter $\lambda$ control the rate of decay, and different values of $\lambda$ are evaluated. The exponential decay rate $\lambda$ represents this exploration-exploitation tradeoff with a higher $\lambda$ meaning the agent explores less than with a lower $\lambda$. Figure 4 shows that values of $\lambda$ smaller than 0.01 cannot get the exit after the person in 250 episodes, and values smaller than 0.001 do not even reach the exit, making very small values of $\lambda$ quite inefficient. Bigger values of $\lambda$ all successfully converge to the optimal strategy roughly in less than 150 episodes, with the very large value of $\lambda = 50$ showing slightly better performance. This is interesting as if $\lambda = 50$, after one episode, the value of epsilon is $\epsilon_1 = 0 + (1 - 0) \times e^{-50 \times 1} = e^{-50} \approx 1.9 \times 10^{-19} \approx 0$. This means that after an initial random step, the agent is always following a pure greedy policy with no exploration but still performing optimally. The reason for this lies in the simple nature of the environment which, without intrinsic exploration in the policy, allows the agent to learn solely by trial and error. With this result we discovered that, for some parameters of $\alpha$ and $\gamma$, the environment is designed in a way that always guides the agent naturally towards the optimal path.

The softmax policy is an intrinsically stochastic policy which gives a probability for each action weighing their expected reward, giving non-zero chance for all actions to occur. The formula is

$$\pi(a) = \frac{e^{\frac{Q(a)}{T}}}{\sum_b e^{\frac{Q(b)}{T}}}$$

where $T \in [0, \infty)$ is the temperature which controls the action selection sampling, $a$ is the action considered and $b$ are all the possible actions. A high temperature means all actions are nearly equally likely to be chosen, thus prioritizing exploration. A low temperature prioritizes sharply the actions with highest expected reward and the agent is closer to greedy, thus prioritizing exploitation.

The decay of temperature $T$ is calculated in the same manner as the decay of epsilon. As before, it follows an exponential decay from $T_{max} = 5$ in the first epoch to $T_{min} = 0$ after infinite epochs. We tried different values of the decay parameter $\lambda$ for temperature. The results align with the analysis done for the epsilon decay parameter above, where $\lambda = 0.01$ does not reach the exit, $\lambda = 0.1$ reaches the person but not the exit, and higher values all converge in less than 200 episodes with the highest value $\lambda = 50$ showing the best performance.

In this environment, we can use high decay rates without problem as exploration is not necessary. We set decay parameters for both policy types to 1. The parameters should be finetuned for any specific environment to get optimal performance.
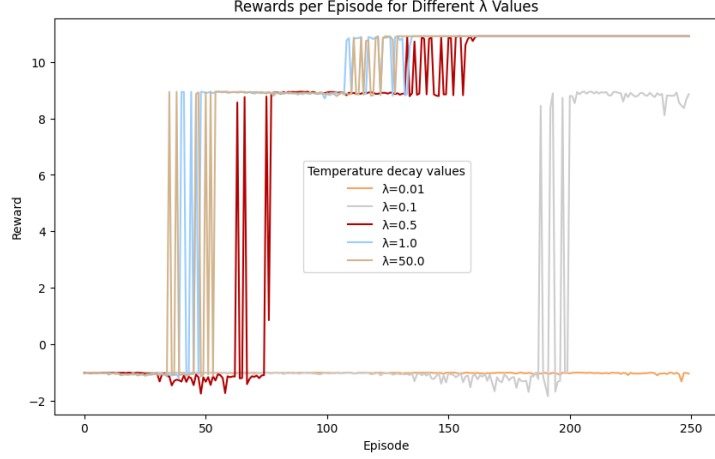
4

**Figure 5:** Comparison of different temperature decay values on the epsilon softmax policy performance. 250 episodes. Terminal rewards. $\gamma = 0.85$. $\alpha = 1$. $T_{min} = 0.001$. $T_{max} = 5$.

# 5   Different reward structures.



**(a)** Limited environment with softmax policy.



**(b)** Limited environment with greedy policy



**(c)** Terminal environment with softmax policy.



**(d)** Terminal environment with greedy policy.

**Figure 6:** Comparative evolution rewards over 250 episodes of the two policies in the two environments. $\gamma = 0.85$. $\alpha = 1$. $\epsilon_{min} = 0$. $\epsilon_{max} = 1$. $T_{min} = 0.001$. $T_{max} = 5$. $\lambda_T = \lambda_\epsilon = 1$.

After evaluating different parameters of the Q-learning agent and two different policies, we now explore its behaviour in the two different reward-transition structures of the environment. In the terminal environment -where the agent dies and is penalized every time it goes to a fire state-, as well as in the limited environment -where the agent cannot attempt to go towards the fire- the agent successfully learns to get the person and then go to the exit. If we look at the rewards per episode (Figure 12), the performance of the agent in a limited environment converges earlier to the optimal path (around 50 episodes) compared to the terminal environment where it takes three times longer (around 160 episodes).

This is a consequence of each specific design. At the beginning, the episodes on the terminal environment are very short with the agent dying often. In contrast, the limited environment episodes only can end when the exit is reached. They have a maximum timestep of 200, which allows the agent to get significant exploration from the beginning, thus learning "faster" in terms of rewards per episode. In terms of timesteps to solve the environment they are much closer with limited environment still outperforming, likely because the agent does not have reset the environment whenever a fire is encountered.

# 6   Conclusion of Q-learning

The results from the experiments show the design of these two environments successfully lead the agent to learn the intended objective, retrieve the target and each the exit. Both the small negative reward for each step by default and the discount factor gamma, pushed the agent to actually perform this task faster in the least amount of timesteps possible. Both policies where successful in making the agent learn. The finetuning of the decay rate -which showed that a purely greedy policy was enough to find the optimal path under certain settings- highlights the efficiency of the Bellman Equation for this environment. By updating the values associated to each state-action pair, the agent creates a meaningful representation of the rewards in the environment, that enables them to follow and find the optimal solution. There is a detailed map of the Q values the agent learned for each state in Appendix A.

# 7   Vanilla DQN and improvements.

## 7.1   Environment and problem.

A challenging environment which cannot be solved feasibly using tabular linear methods such as Q-Learning was chosen for this task. Atari environments with large state spaces are ideal, with extensive research applied to solve this environment. Tabular methods would crash from the state-space memory requirements, while deep learning approaches require a representative subsample to solve it.

Breakout is an Atari game [2] in which the agent with 5 lives must destroy blocks using a bouncing ball and moving a paddle. There are 6 different rows to destroy, with each pair of rows (1-2, 3-4, 5-6) formed by blocks with points 1, 4 and 7 respectively. There are 4 actions available: NOOP (no movement), FIRE, LEFT, RIGHT. A state in the environment is represented as an RGB frame with dimensions 210x160x3. A terminal state occurs when a player has lost all their lives. A combination of lives, ball position, paddle position and brick configuration leads to state space of over $10^{12}$ states [3].

## 7.2   Vanilla DQN

DQN [4] is an extension of the Q-learning algorithm , using Neural Networks (NN) to approximate the state-value function. The network learns by minimizing the loss between the estimated Q-values predicted by the NN and the ground truth Q-values calculated using the Bellman equation and the NN. This loss is propagated through the network and its weights are updated iteratively.

Like Q-Learning algorithm, DQN is an off-policy algorithm, allowing the algorithm to use samples belonging to different policies. Similarly to Q-Learning, exploration-exploitation is important to learn the environment and how to solve the objective. Improvements over this naive implementation of DQN were added [5], improving performance drastically. These improvements include:

1. Experience replay buffer which trains the agent using representative samples of recent past experiences, to decorrelate the agent's experiences and stabilize training.

2. Target network which get temporal difference target values in order for updates to move towards stationary target values, stabilizing training and reducing divergence of the policy.

The pseudocode for the DQN implementation can be seen in Algorithm 1, where $\gamma$ is the discount factor and $L$ the loss function (Mean Square Error):

## 7.3   Double DQN

An issue from DQN is the overestimation of the target Q-values. This occurs as bias may be introduced in the value calculations if the target network produces an error in its estimations of $\max_{a'} \hat{q}_2(s', a', \theta_2)$.

**Algorithm 1** DQN algorithm
_____
 1: Initialise Experience replay memory $M$ to capacity $N$
 2: Initialise policy network $\hat{q}_1$ with parameters $\theta_1 \in \mathbb{R}^d$ arbitrarily
 3: Initialise target action-value network $\hat{q}_2$ with parameters $\theta_2 = \theta_1$
 4: **for** each episode **do**
 5:     Initialise $S$
 6:     Choose action $A$ in state $S$ using policy derived from $\hat{q}_1(S, \cdot, \theta_1)$
 7:     Take action $A$, observe reward $R$ and next-state $S'$
 8:     Store transition $(S, A, R, S')$ in $M$
 9:     **for** transition $(S_j, A_j, R_j, S'_j)$ in minibatch sampled from $D$ **do**
10:         $y = \begin{cases} R_j(s,a) & \text{if } s' \text{ is terminal} \\ R_j(s,a) + \gamma \hat{q}_2(s', \max_{a'} \hat{q}_2(s', a', \theta_2)) & \text{otherwise} \end{cases}$
11:         $\hat{y} \leftarrow \hat{q}_1(S_j, A_j, \theta_1)$
12:         Perform gradient descent step $\nabla_{\theta_1} L(y, \hat{y})$
13:     **end for**
14:     Every $C$ time-steps, update $\theta_2 = \theta_1$
15: **end for**
_____

As we are maximizing over Q-values, an overestimated choice will always be preferred.

Double-DQN (DDQN) [6] mitigates this issue by introducing the policy network for action selection, removing positive bias and decoupling action selection from value estimation. This should lead to more stable convergence and improved policy learning.

Target value $y = \begin{cases} R(s,a) & \text{if } s' \text{ is terminal} \\ R(s,a) + \gamma \hat{q}_2(s', \arg\max_{a'} \hat{q}_1(s', a')) & \text{otherwise} \end{cases}$

## 7.4 Prioritized Experience Replay

Random sampling from a list of experiences is not the most optimal way of learning, as all experiences have equal probability of being sampled, regardless of their contribution to learning. Some experiences are more important than others, and should be sampled more frequently. Prioritized experience replay (PER) [7] is an improvement over regular experience buffers, sampling more often valuable experiences. The value of an experience is measured using the temporal difference error, in our case, between target and estimated Q-values.

This non-uniform sampling can add bias if not regulated, as samples do not belong to the same distribution as the expectation. Weights added to the backpropagated loss are calculated from the probabilities $P(i)$ to correct this bias, with an additional hyperparameter $\beta$ controlling the compensation against bias. These weights are normalized for stability purposes, as well as ensuring weights can only decrease the update downwards. The probabilities $P(i)$ are calculated from the stored priorities $p_i$ and hyperparameter $\alpha$ controls the prioritization used. Probabilities are added a small $\epsilon = 0.01$ constant to avoid probabilities with value 0. $N$ corresponds to the sample batch size.

$$Weights(i) = \left( \frac{1}{N} \cdot \frac{1}{P(i)} \right)^{\beta} \quad \text{and} \quad Probabilities(i) = \frac{p_i^{\alpha}}{\sum_k p_k^{\alpha}}$$

As additional computation is required for sampling and weight calculations, an efficient data structure is necessary. To reduce sampling and update complexity, priorities are stores in a sum-tree data structure, offering a complexity of $O(logN)$, unlike lists with a complexity of $O(N)$.

## 7.5 Implementation

To reduce the state space complexity without sacrificing learning, the observations were converted to grayscale, as the colours of the environment do not contribute to learning the optimal policy. Furthermore, the image resolution was reduced to 84x84, reducing the state space without sacrificing functionality.

Lastly, frame stacking was introduced to allow the network to understand direction and velocity of the ball, joining 4 frames into a single frame. Additionally, to improve learning, all rewards were considered of equal value to encourage the model to remove all blocks rather than maximize rewards. To increase training speed, auto-shoot was implemented when starting a new episode or when starting a new life.

The architecture for the neural network consisted of 3 convolutional layers and 2 fully connected linear layers, with ReLu activation functions between each layer, based on [5]. An epsilon-greedy policy was chosen, linearly decaying $\epsilon$ during a number of frames. The sum-tree data structure was constructed following the code provided in [8]. Our expectation based on previous literature, is to observe the best performance in DQN rather than DDQN, with prioritized experience replay alongside DDQN between both results. DDQN can mitigate the overestimation of some states, while PER should select samples which are more informative to the learning process. Interestingly, DQN with this implementation of PER was not added to benchmarks, and could provide an improvement over DQN.

# 8    Analyse the results quantitatively and qualitatively.

The vast number of hyperparameters can greatly affect the performance of the algorithms. As a competitive starting point, the hyperparameters based on the DQN paper were used on most experiments, with the specific hyperparameters in PER based on values suggested by the PER paper. An additional experiment using slightly tuned hyperparameters for DDQN based on the DDQN paper was also evaluated. The different set of hyperparameters are detailed in Appendix B. All experiments were trained up to 4 million frames for fair comparisons.



**Figure 7:** Mean rewards over 100 episodes for DQN, DDQN, Tuned DDQN and DQN with PER. Note the different lengths correspond to the number of episodes performed by each configuration (except Tuned-DDQN which was stopped early as no improvements were noted)

To evaluate the performance during training, the mean reward over the most recent 100 episodes and the scores per episodes were plotted. Note this evaluation considers all rewards equally. A checkpoint for a model is saved when a best mean reward is achieved. These checkpoints are then evaluated, choosing only greedy actions during inference, and a video of the agent interacting with environment is stored.

DQN with PER achieves the best results, reaching higher rewards in less episodes. DQN followed these results, achieving similar results but requiring more episodes. DDQN and its tuned equivalent achieved similar results by the end but at much slower rate of learning. Curiously, no difference was noted when training DDQN on tuned parameters and DQN hyperparameters, as seen in Figure 7. In terms of time complexity, DQN with PER encumbered a significant cost, almost doubling training time compared to other models (6 hours vs 3-4 hours on Hyperion).

Overall, all models achieved similar rewards, although higher rewards in certain episodes are achieved only in DQN and DQN with PER, as seen in Figure 13. During inference (rewards now have different values), DDQN achieves a best score of 356, DQN a score of 366 and DQN with PER a score of 407. During inference some models struggled to achieve improved results as they got stuck in a loop, with the ball bouncing between a gap of bricks and back to the paddle, repeated continuously in a cycle.

**Figure 8:** Rewards per episode for DQN, Double-DQN, Tuned Double-DQN and DQN with PER.

An interesting observation is the development of strategies. All models developed a tendency to create a hole reaching the ceiling, allowing the ball to repeatedly bounce between the ceiling and the top rows, achieving points rapidly. This long term goal is encouraged by the $\gamma$ value of 0.99. Another strategy to reach the ball on average fas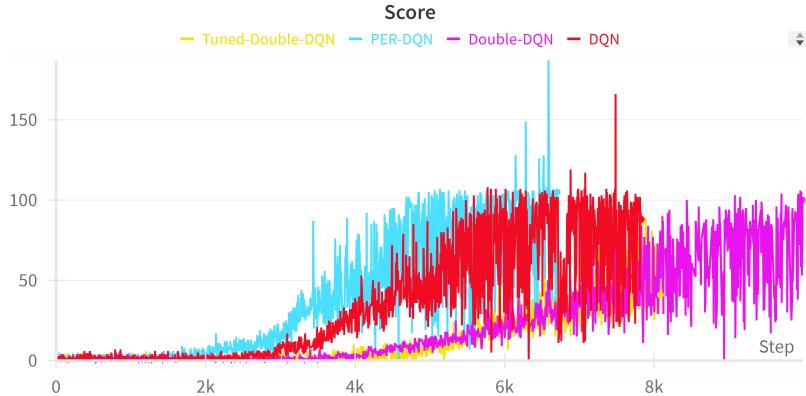ter is to move to the middle of the screen, allowing the agent to move to either directions at a constant speed. Examples of the strategies can be found in the videos uploaded alongside this report and the first strategy can be found in Appendix C.

# 9 Individual Component: Atlantis

For my individual part of the coursework, I decided to implement and analyze the Atari game Atlantis [9]. This is an Atari game featured in the Gymnasium collection where the objective is to defend the underwater city of Atlantis from aerial attacks by shooting to the enemy planes which fly closer to the surface of the water. The player controls three fixed gun turrets located on the center and the sides of the cityscape, firing at enemy aircrafts that pass overhead in varying patterns and speeds. Players can perform 4 actions: NOOP (no movement), FIRE (fire from THE central turret), LEFTFIRE, and RIGHTFIRE. Each destroyed enemy plane -after being shot- increases the player's score, with big planes being worth 600 points and small ones 200. A terminal state is reached when Atlantis is completely destroyed, marked by the loss of all defensive turrets and structures. The Atari game is impossible to win -planes kept coming faster and the city is eventually destroyed-,as such, the RL training stops after 130,000 points are obtained in more than 25 episodes.

This is a complex environment for which each state is represented by the game screen which has dimensions of 210x160x3, presented in RGB. Preprocessing is applied making it a greyscale and 84x84 pixels, which effectively reduces the state space without sacrificing functionality. Also, compression every four frames into one is applied in order to guide the agent with the idea of movement. Also rewards for the agent are divided by 100, which is more manageable for using neural networks. The architecture used is a neural network with 3 convolutional layers to process the visual input of the the current state and 2 fully connected layers for further processing, all using the ReLU activation function. The output of this network is the number of actions.

Three different algorithms where tried: vanilla DQN, Prioritized Experiences Replay (PER) with DQN, and Double DQN (DDQN). They took all less than 8 hours in a NVIDIA 1650 GTX, with PER taking almost twice per episode than vanilla DQN, which make sense given the extra complexity of the samples introduced by the sampling of PER. The networks were trained for 5000 episodes using the hyperparameters on the Appendix B.

# 10 Analysis of the results.

The three networks converged to an strategy that won the game in less than 1300, effectively breaking the training loop (5000 episodes). All of them showcase unstable learning as measured in per episode rewards

**Figure 9:** Rewards for DQN with PER, vanilla DQN, and DDQN. *



**Figure 10:** Mean rewards over 100 episodes for DQN with PER, vanilla DQN, and DDQN. *

(*) Episodes are half of the x-axis because of Weights and Biases logging.

in Figure 9, but general convergence can be observed in Figure 10 by taking the mean over 100 episodes. Both figures differ in their ranges because the peak raw rewards are usually not sustained indefinitely as we saw in tabular Q-learning 12 because here the state space is much bigger and it's likely that the agent hasn't seen this space before, or at least is not processing it in the same way. Comparatively, the PER DQN Network converges faster than the vanilla DQN, which aligns with the background theory that says that sampling more relevant experiences yields faster learning. The DDQN network converges the fastest, outperforming significantly to vanilla DQN and PER DQN; this is also shown in the DDQN paper [6]. The key reason for this may lie in the rapid nature of Atlantis which displays continuous firing of targets that appearing large volumes and diverse positions. DQN tends to overestimate Q-values in highly dynamic environments where the agent must adapt to frequent changes in input, whereas DDQN selects actions that are more consistently beneficial, influencing the long-term return because the game will last for longer.

Thanks to the use of convolutional deep neural networks and the reinforcement learning framework, we are able to train agents that outperform the best humans playing Atlantis by a huge margin, getting 20 times higher score. Humans can come up with heuristics to play beforehand like "whenever a plane appears some distance away of the turret range it should fire" and can also reason about them "because the time the bullet takes to arrive the target will already be there". Despite AI agents being unable to engage with this sort of reasoning, they generate sufficiently good representations to play the game, which coupled with their extremely low reaction times -they can provide an informed response to every four frames of video- it's enough to play super-humanly.

# 11   Contributions

We both contributed evenly to the task having a collaborative approach over the whole project.

We both chose the Task 1 environment, and discussed the implementation of the different modules. He did the backbone of the code (structure of code and plotting) and I added some key modules (reward functions). Then I leaded the parameter finetuning section. We both discussed several environments for Task 2 and converged on Breakout. He did the backbone structure of the code and tweaked the hyperparameters. I implemented the first version of Prioritized Experience Replay. The code was all discussed, checked, reviewed and debugged. The report was evenly written, and peer-reviewed (between us).

# References

[1] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.

[2] OpenAI Gymnasium, "Breakout environment." OpenAI Gymnasium for Atari, 2024.

[3] P. Nash, "A brief history of reinforcement learning in game play." `https://www.projectnash.com/a-brief-history-of-reinforcement-learning-in-game-play/`, 2021. Accessed: 2024-05-10.

[4] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," 2013.

[5] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[6] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," 2015.

[7] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," 2016.

[8] Howuhh, "Prioritized experience replay." `https://github.com/Howuhh/prioritized_experience_replay`, 2016. Implementation with proportional prioritization.

[9] OpenAI Gymnasium, "Atlantis environment." OpenAI Gymnasium for Atari, 2024.

# A    Appendix: Q values map



**(a)** Limited environment. Person yet not collected.



**(b)** Limited environment. Person already collected.



**(c)** Terminal environment. Person yet not collected.



**(d)** Terminal environment. Person already collected.

**Figure 12:** Maximum Q value of an action at each state of an optimal agent. $\gamma = 0.85$

# B   Appendix: Task 2 Hyperparamaters

| Hyperparameter DQN, DDQN and PER | Value |
|---|---|
| Discount Factor $\gamma$ | 0.99 |
| Minibatch size | 64 |
| Experience replay buffer size | 100000 |
| Target Network update frequency | 10000 |
| Network Update frequency | 4 |
| Learning rate | 0.0001 |
| Initial exploration epsilon | 1 |
| Final exploration epsilon | 0.01 |
| Epsilon linear decay frames | 500000 |
| Replay buffer start size | 5000 |
| No op max | 30 |
| PER $\alpha$ | 0.6 |
| PER $\beta$ | 0.4 |

**Table 2:** First set of hyperparameters for DQN, DDQN and PER variants

# C   Appendix: Breakout strategy



**Figure 13:** s trategy developed during Breakout, where the agent creates a hole in the edges of the environment in order to bounce the ball between the ceiling and the top row of blocks.

# D   Code

## D.1   Task 1

### D.1.1   Maze_env.py

```python
import numpy as np
from matplotlib import pyplot as plt


class Maze_env:
    """
    Represents a maze navigation environment for reinforcement learning tasks.
    It manages the maze layout, positions of start, target, and coin, and rewards/transitions.
    Functionality includes:
    - '__init__(start, target, coin, maze)': Initializes the environment.
    - 'plot_env()': Visualizes the maze with important positions highlighted.
    - 'plot_env_position(position, timestep)': Visualizes  maze with agent's position at
        specific timestep.
    - 'create_r_matrix()': Generates a reward matrix based on the maze layout.
    - 'reward(state, action)': Calculates the reward for an action taken from a state.
    - 'transition(state, action)': Determines the new state after an action.
    - 'done()': Checks if the target has been reached, ending the episode.
    - 'create_q_matrix()': Initializes a Q-learning matrix for action selection.
    """

    def __init__(self, start, target, coin, maze, reward_type):
        self.maze = maze
        self.target = target
        self.start = start
        self.coin = coin
        self.reward_type = reward_type
        self.position = 0
        self.R = self.create_r_matrix(self.reward_type)
        print(f"Shape of the R matrix is {self.R.shape}")
        self.Q = self.create_q_matrix()
        print(f"Shape of the Q matrix is {self.Q.shape}")
        self.coin_collected = False
        self.terminate = False

    def plot_env(self):
        cmap = plt.cm.colors.ListedColormap(
            ["white", "orange", "red", "blue", "yellow"]
        )
        maze_plot = self.maze.copy()
        maze_plot[self.target] = 2
        maze_plot[self.start] = 3
        maze_plot[self.coin] = 4
        plt.imshow(maze_plot, cmap=cmap)
        plt.show()

    def plot_env_position(self, position, timestep):
        cmap = plt.cm.colors.ListedColormap(
            ["white", "orange", "red", "blue", "yellow"]
        )
        maze_plot = self.maze.copy()
        maze_plot[self.target] = 2
        maze_plot[position] = 3
        maze_plot[self.coin] = 4
        plt.imshow(maze_plot, cmap=cmap)
        plt.savefig(f"img/plot_{timestep:06d}.png", dpi=300)
        plt.show()
        plt.close()

    def create_r_matrix(self, reward_type):
        """
        This synthesizes the reward and transition functions.
        reward_type (str): The type of reward to use.
        Options are "terminal_movement" and "free_movement".
        """
        actions = [(-1, 0), (1, 0), (0, -1), (0, 1)]
        num_states = self.maze.shape[0] * self.maze.shape[1] * 2  # times coin state
        coin_states = 2  # 0 for no coin collected, 1 for coin collected
        R = np.full(
            (self.maze.shape[0], self.maze.shape[1], coin_states, len(actions)), np.nan)

        if reward_type == "terminal_movement":
            print("Reward type: Terminal Movement")
            # actions beyond limits get -10 (and terminate)
            # actions to a 0 -10 (and terminate)
            # action to coin get 200
            # action to target get 100
            # allowed actions get -1 (for the time)
```
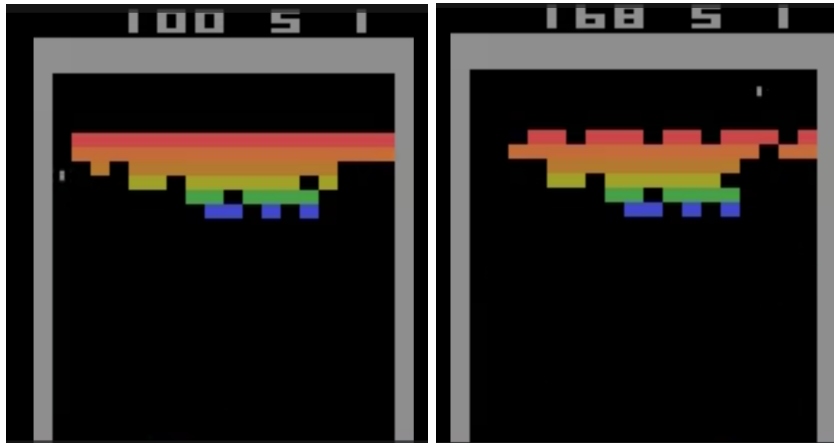
15

```python
                    for i in range(self.maze.shape[0]):
                        for j in range(self.maze.shape[1]):
                            for coin_state in range(coin_states):
                                for action_index, action in enumerate(actions):
                                    new_i, new_j = i + action[0], j + action[1]
                                    if new_i >= 0 and new_i < self.maze.shape[0] and new_j >= 0 and
                                        new_j < self.maze.shape[1]:
                                        # Actions to a wall (1 in the maze) get -1
                                        if self.maze[new_i, new_j] == 1:
                                            R[i, j, coin_state, action_index] = -1 # for the fire
                                        elif self.maze[new_i, new_j] == 0:
                                            R[i, j, coin_state, action_index] = -0.005  # for an
                                                allowed action
                                            if (new_i, new_j) == self.coin and not coin_state:
                                                # print("Assigning coin")
                                                R[i, j, coin_state, action_index] = 10  # coin
                                            elif (new_i, new_j) == self.target:
                                                # print("Assigning target")
                                                R[i, j, coin_state, action_index] = 1  # target
                                    else:
                                        R[i, j, coin_state, action_index] = -1  # actions beyond the
                                            limits are forbidden

                return R
                # then add the transition function so that if reward smaller than -1, then
                    terminate.

        if reward_type == "limited_movement":
            print("Reward type: Limited Movement")
            # actions beyond limits get None (can't move)
            # actions to a 0 (get -10)
            # action to coin get 200
            # action to target get 100
            # allowed actions get -1 (for the time)

            for i in range(self.maze.shape[0]):
                for j in range(self.maze.shape[1]):
                    for coin_state in range(coin_states):
                        for action_index, action in enumerate(actions):
                            new_i, new_j = i + action[0], j + action[1]

                            if new_i >= 0 and new_i < self.maze.shape[0] and new_j >= 0 and
                                new_j < self.maze.shape[1]: # inside of maze
                                # Actions to a wall (1 in the maze) get None
                                if self.maze[new_i, new_j] == 1:
                                    R[i, j, coin_state, action_index] = None # for the fire
                                elif self.maze[new_i, new_j] == 0:
                                    R[i, j, coin_state, action_index] = -0.005  # for an
                                        allowed action
                                    if (new_i, new_j) == self.coin and not coin_state:
                                        R[i, j, coin_state, action_index] = 10  # coin
                                    elif (new_i, new_j) == self.target:
                                        R[i, j, coin_state, action_index] = 1  # target
                            else:
                                R[i, j, coin_state, action_index] = None  # actions beyond the
                                    limits are forbidden
            return R

    def transition_R(self, state, action, reward_type):
        initial_state = state
        x, y = initial_state
        new_x = x
        new_y = y
        if action == 0:  # up
            new_x -= 1
        elif action == 1:  # down
            new_x += 1
        elif action == 2:  # left
            new_y -= 1
        elif action == 3:  # right
            new_y += 1

        if reward_type == "terminal_movement":
            if new_x >= 0 and new_x < self.maze.shape[0] and new_y >= 0 and new_y < self.maze.
                shape[1]:
                if self.R[x, y, int(self.coin_collected), action] == -1: # fire
                    # print("Fire")
                    self.terminate = True
                    return state
                elif (new_x,new_y) == self.coin and not self.coin_collected: # coin
                    # print("Coin")
                    self.coin_collected = True
                    # print(self.coin_collected)
                    return new_x,new_y
                elif (new_x,new_y) == self.target: # target
                    # print("Target")
                    self.terminate = True
```

```
158                         return new_x,new_y
159                     elif self.R[x, y, int(self.coin_collected), action] == -0.005: # normal action
160                         # print("Allowed")
161                         return new_x,new_y
162                 else:
163                     self.terminate = True  # walls
164                     # print("Out of bounds")
165                     return state
166
167         if reward_type == "limited_movement": # should not attempt to access fire or wall
168             if new_x >= 0 and new_x < self.maze.shape[0] and new_y >= 0 and new_y < self.maze.
                 shape[1]:
169                 if (new_x,new_y) == self.coin and not self.coin_collected: # coin
170                     self.coin_collected = True
171                     return new_x,new_y
172                 elif (new_x,new_y) == self.target: # target
173                     self.terminate = True
174                     return new_x,new_y
175                 elif self.R[x, y, int(self.coin_collected), action] == -0.005: # normal action
176                     return new_x,new_y
177
178     def done(self):
179         return self.terminate
180
181     def coin_reached(self):
182         return self.coin_collected
183
184     def create_q_matrix(self):
185         Q = np.zeros_like(self.R)
186         return Q
187
188 if __name__ == "__main__":
189     maze = np.array(
190         [
191             [1,  0,  1,  1,  1,  1,  1,  0,  0,  1],
192             [0,  1,  1,  1,  1,  0,  1,  0,  1,  1],
193             [0,  0,  0,  0,  1,  0,  0,  0,  0,  0],
194             [0,  1,  1,  0,  0,  0,  0,  0,  1,  0],
195             [1,  0,  0,  0,  0,  0,  0,  1,  0,  0],
196             [0,  0,  0,  0,  1,  1,  0,  1,  0,  0],
197             [1,  1,  1,  0,  0,  1,  0,  0,  0,  1],
198             [0,  1,  0,  0,  0,  0,  0,  0,  1,  1],
199             [1,  1,  1,  1,  0,  1,  0,  1,  1,  0],
200             [0,  1,  0,  0,  1,  1,  1,  0,  0,  1],
201         ]
202     )
203     env = Maze_env((2, 0), (0, 8), (7, 5), maze, reward_type="terminal_movement")
204     env.plot_env()
205     print("Info")
206     print(env.R[0, 7, 0])
207     print(env.R[0, 7, 1])
208     print(env.R[7, 5, 0])
209     print(env.R[7, 5, 1])
210     print(env.R[7, 4, 0])
211     print(env.R[7, 4, 1])
```

**Listing 1:** Maze environment code

### D.1.2   agent.py

```
1
2  import numpy as np
3  import cv2
4  import os
5  from matplotlib import pyplot as plt
6  from maze_env import Maze_env
7  from tqdm.auto import tqdm
8
9
10 class Q_learning:
11     """
12     Implements the Q-learning algorithm for reinforcement learning tasks within a predefined
            environment.
13     This class is responsible for learning optimal action-selection policies to maximize
            rewards over episodes of interactions with the environment.
14
15     - `__init__(alpha, gamma, epsilon, episodes, steps, env, states)`: Initializes the
            learning parameters, environment, and states.
16     - `plot_rewards()`: Plots the rewards accumulated over each episode, visualizing the
            learning progress.
17     - `show_Q_spec(coord)`: Displays Q-values for a specific coordinate/state.
18     - `greedy_policy(state)`: Selects an action based on a greedy policy (highest Q-value)
            with an epsilon chance of random action for exploration.
19     - `softmax_policy(state, temperature)`: Selects an action based on the softmax of Q-values
            , factoring in the temperature for exploration-exploitation balance.
```

```
20      - 'train()': Conducts the learning process over a specified number of episodes and steps
            per episode, updating Q-values based on the received rewards.
21      - 'create_video()': Generates a video from saved images of the agent's journey through the
            maze, illustrating the learned policy in action.
22      - 'test(limit)': Evaluates the learned policy by navigating the environment for a given
            number of steps, visualizing the path taken and summarizing the rewards.
23
24      The class utilizes epsilon-greedy and softmax policies for action selection, balancing the
            exploration of the state space with the exploitation of known rewards.
25      """
26
27      def __init__(self, alpha, gamma, epsilon, episodes, steps, env, policy):
28          self.alpha = alpha
29          self.gamma = gamma
30          self.epsilon = epsilon
31          self.temperature = 50.0
32          self.policy = policy
33          self.R = env.R
34          self.R_mod = self.R
35          self.Q = env.Q
36          self.episodes = episodes
37          self.steps = steps
38          self.start = env.start
39          self.target = env.target
40          self.coin = env.coin
41          self.env = env
42          self.episodes_rewards = []
43          self.max_list_size = 10
44          self.list_rewards = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
45          self.threshold = 2
46          self.window_size = 20
47          self.current_average = 0
48
49          # print("Initial Q matrix shape is '{}'".format(self.Q.shape))
50          # print("Initial Q matrix values are '{}'".format(self.Q))
51
52      def plot_rewards(self):
53          plt.plot(self.episodes_rewards)
54          plt.show()
55
56      def show_Q_spec(self, coord):
57          i, j = coord
58          print(self.Q[i, j, int(self.env.coin_reached()), :])
59
60      def greedy_policy(self, state):
61          i, j = state
62          available_actions = np.where(~np.isnan(self.R_mod[i, j, int(self.env.coin_reached())])
                )[0]
63          # print(available_actions)
64          q_values = [self.Q[i, j, int(self.env.coin_reached()), a] for a in available_actions]
65          best_actions = available_actions[np.where(q_values == np.max(q_values))[0]]
66          # print(best_actions)
67
68          # available_actions = np.array([0, 1, 2, 3])
69          # q_values = [self.Q[state, a] for a in available_actions]
70          # best_actions = available_actions[np.where(q_values == np.max(q_values))[0]]
71
72          if np.random.uniform() < self.epsilon:
73              # a = np.random.choice(4)
74              a = np.random.choice(available_actions)
75          else:
76              #                    a = np.argmax(self.Q[s,:])
77              a = np.random.choice(best_actions)
78          # print(a)
79          return a
80
81      def softmax_policy(self, state):
82          i, j = state
83          available_actions = np.where(~np.isnan(self.R_mod[i, j, int(self.env.coin_reached())])
                )[0]
84          # print(f"Available actions: {available_actions}")
85          q_values = [self.Q[i, j, int(self.env.coin_reached()), a] for a in available_actions]
86          max_q_value = np.max(q_values)
87          exp_values = np.exp((q_values - max_q_value) / self.temperature)
88          action_probs = exp_values / np.sum(exp_values)
89          # print(f"Actions Probability: {action_probs}")
90          # Sample an action based on the probabilities
91          selected_action_index = np.random.choice(len(action_probs), p=action_probs)
92          selected_action = available_actions[selected_action_index]
93          # print(f"Selected Action: {selected_action}")
94
95          return selected_action
96
97      def train(self):
98          print("Target is '{}'".format(self.target))
99          print("Starting state is '{}'".format(self.start))
100
```

```python
            for episode in tqdm(range(self.episodes), desc= f"Training agent on {self.episodes}
                episodes", unit="episode", total=self.episodes):
                # print("New episode")
                s = self.start
                episode_reward = 0
                self.env.coin_collected = False
                self.env.terminate = False
                # print("New episode")
                for timestep in range(self.steps):
                    # print(self.env.coin_reached())
                    i, j = s
                    # Epsilon-greedy action choice
                    if self.policy == "greedy":
                        a = self.greedy_policy(s)
                    elif self.policy == "softmax":
                        a = self.softmax_policy(s)
                    else:
                        raise ValueError("Policy must be 'greedy' or 'softmax'")
                    #a = self.softmax_policy(s, self.temperature)
                    # Environment updating
                    # r = env.reward(s, a)
                    # print(self.R_mod[i,j,int(self.env.coin_collected)])
                    r = self.R_mod[i, j, int(self.env.coin_reached()), a]
                    # print(r)
                    # if self.env.coin_reached():
                    #     print()
                    #     print("Coin collected")
                    # print(self.env.coin_reached())
                    # print("Reward")
                    episode_reward += r
                    new_state = self.env.transition_R((i, j), a, self.env.reward_type)
                    # print(f"Action is {a}, and state is {(i, j)}")
                    new_i, new_j = new_state

                    if r == 10: # picked up coin for first time
                        # print("COOOOOOOOOOOOOOOOOOIIIIIIIINNNNNNNN")
                        # Current q in state o and next in state 1 for coin_collected
                        self.Q[i, j, 0, a] = self.Q[i, j, 0, a] + self.alpha * (r + self.gamma *
                            np.max(
                            self.Q[new_i, new_j, 1, :]) - self.Q[i, j, 0, a])
                    else:
                        self.Q[i, j, int(self.env.coin_reached()), a] = self.Q[i, j, int(
                            self.env.coin_reached()), a] + self.alpha * (r + self.gamma * np.max(
                            self.Q[new_i, new_j, int(self.env.coin_reached()), :]) - self.Q[
                                                                            i, j, int(self.env.
                                                                                coin_reached()),
                                                                                a])

                    if self.env.done():
                        # print("Death")
                        # print(self.env.terminate)
                        break
                    s = new_state

                self.episodes_rewards.append(episode_reward)

                self.list_rewards.append(episode_reward)
                if len(self.list_rewards) > self.max_list_size:
                    self.list_rewards.pop(0)
                window = self.list_rewards[-self.window_size:]
                window_average = sum(window) / self.window_size
                self.current_average = window_average

                if episode == self.episodes - 1:
                    if self.policy == "greedy":

                        print(
                            "Episode {} finished. Episode Reward {}. Timesteps {}. Average {}.
                                Epsilon {}".format(
                                episode,
                                episode_reward,
                                timestep,
                                window_average,
                                self.epsilon,
                            )
                        )
                    else:
                        print(
                            "Episode {} finished. Episode Reward {}. Timesteps {}. Average {}.
                                Temp {}".format(
                                episode,
                                episode_reward,
                                timestep,
                                window_average,
                                self.temperature,
                            )
                        )
```

19

```python
182                lambda_rate = 0.1
183                mimimum_epsilon = 0.00
184                initial_epsilon = 1
185                self.epsilon = mimimum_epsilon + (initial_epsilon − mimimum_epsilon) * np.exp(−
                       lambda_rate * episode)
186
187                lambda_rate_temp = 0.1
188                minimum_temperature = 0.001
189                initial_temperature = 5
190                self.temperature = minimum_temperature + (initial_temperature −
                       minimum_temperature) * np.exp(−lambda_rate_temp * episode)
191
192        def create_video(self):
193            image_folder = "img"   # Directory containing your saved plot images
194            video_name = "video_agent.mp4"
195
196            images = [
197                img
198                for img in os.listdir(image_folder)
199                if img.endswith((".jpg", ".jpeg", ".png"))
200            ]
201            frame = cv2.imread(os.path.join(image_folder, images[0]))
202            height, width, layers = frame.shape
203
204            video = cv2.VideoWriter(
205                video_name, cv2.VideoWriter_fourcc(*"mp4v"), 1, (width, height)
206            )
207
208            for image in images:
209                video.write(cv2.imread(os.path.join(image_folder, image)))
210
211            cv2.destroyAllWindows()
212            video.release()
213
214        def test(self, limit=40):
215            s = self.start
216            print("Starting state is '{}'".format(s))
217            episode_reward = 0
218            env.coin_collected = False
219            env.terminate = False
220            for timestep in range(limit):
221                i, j = s
222                self.env.plot_env_position(s, timestep)
223                a = np.argmax(self.Q[i, j, int(self.env.coin_reached())])
224
225                # Environment updating
226                r = self.R_mod[i, j, int(self.env.coin_reached()), a]
227                print(f"Step {timestep}. Action is {a}. State is {(i, j)}. Q value of {self.Q[i, j
                       , int(self.env.coin_reached()), a]}. And reward {r}")
228                episode_reward += r
229                new_state = self.env.transition_R((i, j), a, self.env.reward_type)
230                new_i, new_j = new_state
231
232                if env.done():
233                    self.env.plot_env_position(new_state, timestep+1)
234                    break
235                s = new_state
236            # print('Episode Reward {}.Q matrix values:\n{}'.format(episode_reward, self.Q.round
                   (1)))
237            self.create_video()
238
239
240 if __name__ == "__main__":
241     maze = np.array(
242         [
243             [1, 0, 1, 1, 1, 1, 1, 0, 0, 1],
244             [0, 1, 1, 1, 1, 0, 1, 0, 1, 1],
245             [0, 0, 0, 0, 1, 0, 0, 0, 0, 0],
246             [0, 1, 1, 0, 0, 0, 0, 0, 1, 0],
247             [1, 0, 0, 0, 0, 0, 0, 1, 0, 0],
248             [0, 0, 0, 0, 1, 1, 0, 1, 0, 0],
249             [1, 1, 1, 0, 0, 1, 0, 0, 0, 1],
250             [0, 1, 0, 0, 0, 0, 0, 0, 1, 1],
251             [1, 1, 1, 1, 0, 1, 0, 1, 1, 0],
252             [0, 1, 0, 0, 1, 1, 1, 0, 0, 1],
253         ]
254     )
255     env = Maze_env(start=(2, 0), target=(0, 8), coin=(7, 5), maze=maze, reward_type="
           limited_movement")
256
257     q_learning = Q_learning(alpha=0.9, gamma=0.85, epsilon=1, episodes=250, steps=200, env=env
           , policy="softmax")
258     print("INFO. State is (ROW, COLUMN IS_COIN). Action is [up, down, left, right]")
259     print(f" R values for state (0, 7, 0) {q_learning.R_mod[0, 7, 0]}")
260     print(f" R values for state (0, 7, 1) {q_learning.R_mod[0, 7, 1]}")
261     print(f" R values for state (7, 5, 0) {q_learning.R_mod[7, 5, 0]}")
262     print(f" R values for state (7, 5, 1) {q_learning.R_mod[7, 5, 1]}")
```

```
263    print(f" R values for state (7, 4, 0) {q_learning.R_mod[7, 4, 0]}")
264    print(f" R values for state (7, 4, 1) {q_learning.R_mod[7, 4, 1]}")
265
266    q_learning.train()
267    q_learning.plot_rewards()
268    # q_learning.test()
269    print(f" Q values for state (3, 3, 0) {q_learning.Q[3, 3, 0]}")
270    print(f" Q values for state (7, 4, 0) {q_learning.Q[7, 4, 0]}")
```

**Listing 2:** Q-Learning agent code

## D.2 Task 2 code

### D.2.1 buffer.py

```
1    from collections import namedtuple, deque
2    import numpy as np
3    import random
4    import torch
5
6
7    class ReplayMemory(object):
8        def __init__(self, capacity, use_per=False, alpha=0.6, epsilon=0.001):
9            self.use_per = use_per
10           self.capacity = capacity
11           self.memory = deque([], maxlen=capacity)
12           self.epsilon = epsilon
13           self.count = 0
14           if self.use_per:
15               self.alpha = alpha
16               self.sum_tree = SumTree(capacity)
17               self.max_priority = 1.0
18
19       def push(self, transition):
20           self.memory.append(transition)
21           if self.use_per:
22               self.sum_tree.add(self.max_priority, self.count)
23               self.count = (self.count + 1) % self.capacity
24
25       def sample(self, batch_size, beta=0.4):
26           batch = []
27           idxs = []
28           is_weights = []
29           if self.use_per:
30               total_priority = self.sum_tree.total
31               # print(f"Total Priority: {total_priority}")
32               segment = total_priority / batch_size
33               for i in range(batch_size):
34                   # Guard against sampling error: https:// github.com / rlcode / per / issues
                         / 4
35                   while True:
36                       s = random.uniform(segment * i, segment * (i + 1))
37                       tree_idx, priority, idx = self.sum_tree.get(s)
38                       if idx is not None:
39                           break
40                       else:
41                           print("Attempted to sample unitialised memory")
42                   sampling_probability = priority / total_priority
43                   is_weight = (len(self.memory) * sampling_probability) ** -beta
44                   is_weights.append(is_weight)
45                   # print(idx)
46                   batch.append(self.memory[idx])
47                   idxs.append(tree_idx)
48               max_weight = max(is_weights)
49               is_weights = [w / max_weight for w in is_weights]
50           else:
51               batch = random.sample(self.memory, batch_size)
52               is_weights = [1.0] * batch_size
53               idxs = None
54           return batch, idxs, is_weights
55
56       def update_priority(self, idxs, priorities):
57           if self.use_per:
58               for idx, priority in zip(idxs, priorities):
59                   adjusted_priority = (priority + self.epsilon) ** self.alpha
60                   self.max_priority = max(self.max_priority, adjusted_priority)
61                   self.sum_tree.update(idx, adjusted_priority)
62           else:
63               raise ValueError("Not using PER")
64
65       def __len__(self):
66           return len(self.memory)
67
68
```

```
69  # https://github.com/Howuhh/prioritized_experience_replay/blob/main/memory/tree.py
70  class SumTree:
71      """This will be binary tree stored as a list (self.tree), where:
72       − the experiences priorities are the leaves, stored in the second half of the list
73       − the remaining positions (first half) are the binary sums of children nodes
74       − the root tree (the first element) is the sum of all the elements"""
75      def __init__(self, size):
76          self.nodes = [0] * (2 * size − 1)
77          self.data = [None] * size
78
79          self.size = size
80          self.count = 0
81          self.real_size = 0
82
83      @property
84      def total(self):
85          return self.nodes[0]
86
87      def update(self, data_idx, value):
88          idx = data_idx + self.size − 1  # child index in tree array
89          change = value − self.nodes[idx]
90          self.nodes[idx] = value
91          parent = (idx − 1) // 2
92          while parent >= 0:
93              self.nodes[parent] += change
94              parent = (parent − 1) // 2
95
96      def add(self, value, data):
97          self.data[self.count] = data
98          self.update(self.count, value)
99          self.count = (self.count + 1) % self.size
100         self.real_size = min(self.size, self.real_size + 1)
101
102     def get(self, cumsum):
103         assert cumsum <= self.total
104
105         idx = 0
106         while 2 * idx + 1 < len(self.nodes):
107             left, right = 2*idx + 1, 2*idx + 2
108             if cumsum <= self.nodes[left]:
109                 idx = left
110             else:
111                 idx = right
112                 cumsum = cumsum − self.nodes[left]
113
114         data_idx = idx − self.size + 1
115         return data_idx, self.nodes[idx], self.data[data_idx]
```

**Listing 3:** Buffe which allows Prioritized Sampling

## D.2.2  dqn.py

```
1   import gymnasium as gym
2   from gymnasium.utils.save_video import save_video
3
4   import math
5   import random
6   import matplotlib
7   import matplotlib.pyplot as plt
8   from collections import namedtuple, deque
9   from itertools import count
10  import torch
11  import torch.nn as nn
12  import torch.optim as optim
13  import torch.nn.functional as F
14  import numpy as np
15  import os
16  from buffer import ReplayMemory
17  from logger import Logger
18
19  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
20  print(f"Device is {device}")
21
22  Transition = namedtuple('Transition',
23                          ('state', 'action', 'reward', 'next_state', 'done'))
24
25  os.environ['https_proxy'] = "http://hpc−proxy00.city.ac.uk:3128"
26
27
28  class DQN(nn.Module):
29
30      def __init__(self, n_observations, n_actions, hidden_units=512):
31          super(DQN, self).__init__()
32          self.layer1 = nn.Linear(n_observations, hidden_units)
33          self.layer2 = nn.Linear(hidden_units, hidden_units)
```

```python
34            self.layer3 = nn.Linear(hidden_units, n_actions)
35
36        def forward(self, x):
37            x = F.relu(self.layer1(x))
38            x = F.relu(self.layer2(x))
39            return self.layer3(x)
40
41
42    class DQNCNN(nn.Module):  # DQN/DDQN
43        def __init__(self, input_shape, n_actions, hidden_units=512):
44            super(DQNCNN, self).__init__()
45            self.conv = nn.Sequential(
46                nn.Conv2d(input_shape[0], 32, kernel_size=8, stride=4),
47                nn.ReLU(),
48                nn.Conv2d(32, 64, kernel_size=4, stride=2),
49                nn.ReLU(),
50                nn.Conv2d(64, 64, kernel_size=3, stride=1),
51                nn.ReLU()
52            )
53            conv_out_size = self.get_conv_out_size(input_shape)
54
55            self.value = nn.Sequential(
56                nn.Linear(conv_out_size, hidden_units),
57                nn.ReLU(),
58                nn.Linear(hidden_units, n_actions)
59            )
60
61        def get_conv_out_size(self, shape):
62            conv_size = self.conv(torch.zeros(1, *shape))
63            return int(np.prod(conv_size.size()))
64
65        def forward(self, x):
66            conv_out = self.conv(x).view(x.size()[0], -1)
67            return self.value(conv_out)
68
69
70    class Agent:
71        def __init__(self, env, per=False, double = False, logger = None):
72
73            self.logger = logger
74            self.GAMMA = 0.99
75            self.LR = 1e-4
76            self.ALPHA = 1
77            self.update_frequency = 4
78            self.update_target_frequency = 10000 # 20k for tuned ddqn
79            self.batch_size = 64
80            self.per = per
81            self.double_dqn = double
82
83            self.replay = ReplayMemory(100000, use_per=self.per)
84            if self.per:
85                self.alpha = self.replay.alpha
86                self.sum_tree = self.replay.sum_tree
87                self.max_priority = self.replay.max_priority
88            self.memory = self.replay.memory
89
90            self.max_episodes = 5000
91            self.number_episodes = 0
92            self.max_timesteps = 2000
93            self.number_timesteps = 0
94            self.epsilon = 1
95
96            # Get number of actions from gym action space
97            self.env = env
98            self.n_actions = 4
99            self.number_lives = 5
100           # self.n_actions = env.action_space.shape[0]
101           # num_bins = 61  # Number of bins for each action dimension
102           # self.n_actions = num_bins ** self.n_actions
103           print(self.n_actions)
104           print(f"Number actions: {self.n_actions}")
105           seed = None
106           self.random_state = np.random.RandomState() if seed is None else np.random.RandomState(seed)

107
108           # Get the number of state observations
109           self.state, self.info = env.reset()
110           print(f"State shape: {self.state.shape}")
111           # self.n_observations = len(self.state)
112           self.n_observations = self.state.shape
113           self.policy_net = DQNCNN(self.n_observations, self.n_actions, hidden_units=512).to(device)
114           self.target_net = DQNCNN(self.n_observations, self.n_actions, hidden_units=512).to(device)
115           self.target_net.load_state_dict(self.policy_net.state_dict())
116           self.optimizer = optim.AdamW(self.policy_net.parameters(), lr=self.LR, amsgrad=True)
117           print(self.n_observations)
```

```python
118            print(env.observation_space.shape)
119            # self.policy_net = DQN(env.observation_space.shape, self.n_actions).to(device)
120            # self.target_net = DQN(env.observation_space.shape, self.n_actions).to(device)
121
122            self.video = []
123
124    # def discrete2cont_action(self, action):
125    #     # Map the discrete action index to continuous torques
126    #     num_bins = 61
127    #     action_indices = np.unravel_index(action, (num_bins, num_bins, num_bins))
128    #     torque_min = -1.0
129    #     torque_max = 1.0
130    #     torques = [torque_min + (torque_max - torque_min) * idx / (num_bins - 1) for idx in
               action_indices]
131    #     return np.array(torques)
132
133    def has_sufficient_experience(self):
134        """True if agent has enough experience to train on a batch of samples; False otherwise
               ."""
135        # return len(self.memory) >= self.batch_size
136        if len(self.memory) == 5000:
137            print("Sufficient experience recently obtained!!!")
138        return len(self.memory) >= 5000
139
140    def has_full_experience(self):
141        """True if agent has enough experience to train on a batch of samples; False otherwise
               ."""
142        # return len(self.memory) >= self.batch_size
143        if len(self.memory) == 100000:
144            return len(self.memory) >= 100000
145
146    def save(self, filepath):
147        checkpoint = {
148            "q-network-state": self.policy_net.state_dict(),
149            "optimizer-state": self.optimizer.state_dict(),
150        }
151        torch.save(checkpoint, filepath)
152
153    def choose_action(self, state):
154        # print(state.shape)
155        # need to reshape state array and convert to tensor
156        state_tensor = (torch.from_numpy(np.array(state)).unsqueeze(dim=0).to(device)).float()
157        # choose uniform at random if agent has insufficient experience
158        if not self.has_sufficient_experience():
159            action = self.uniform_random_policy(state_tensor)
160        else:
161            # print("Sufficient experience")
162            action = self.epsilon_greedy_policy(state_tensor, self.epsilon)
163        return action
164
165    def epsilon_greedy_policy(self, state, epsilon):
166        """With probability epsilon explore randomly; otherwise exploit knowledge optimally.
               """
167        if self.random_state.random() < epsilon:
168            action = self.uniform_random_policy(state)
169        else:
170            action = self.greedy_policy(state)
171        return action
172
173    def uniform_random_policy(self, state):
174        """Choose an action uniformly at random."""
175        # random_vector = np.random.(low=-1, high=1, size=self.n_actions)
176        # return random_vector
177        return self.random_state.randint(self.n_actions)
178
179    def greedy_policy(self, state):
180        # print(state.shape)
181        # print(state.dtype)
182        """Choose an action that maximizes the action_values given the current state."""
183        action = (self.policy_net(state)
184                    .argmax()
185                    .cpu()  # action_values might reside on the GPU!
186                    .item())
187        return action
188
189    def select_greedy_actions(self, states, q_network):
190        _, actions = q_network(states).max(dim=1, keepdim=True)
191        # print(actions)
192        return actions
193
194    def evaluate_selected_actions(self, states, actions, rewards, dones, gamma, q_network):
195        """Compute the Q-values by evaluating the actions given the current states and Q-
               network."""
196        next_q_values = q_network(states).gather(dim=1, index=actions)
197        q_values = rewards + (gamma * next_q_values * (1 - dones))
198        return q_values
199
```

```python
200        def q_learning_update(self, states, rewards, dones, gamma, q_network):
201            """Q-Learning update with explicitly decoupled action selection and evaluation steps.
               """
202            actions = self.select_greedy_actions(states, q_network)
203            q_values = self.evaluate_selected_actions(states, actions, rewards, dones, gamma,
                   q_network)
204            return q_values
205
206        def double_q_learning_update(self, states, rewards, dones, gamma, q_network1, q_network2):
207            """Q-Learning update with explicitly decoupled action selection and evaluation steps.
               """
208            actions = self.select_greedy_actions(states, q_network1)
209            q_values = self.evaluate_selected_actions(states, actions, rewards, dones, gamma,
                   q_network2)
210            return q_values
211
212        def learn(self, experiences, is_weights, idxs):
213            """Update the agent's state based on a collection of recent experiences."""
214            states, actions, rewards, next_states, dones = (torch.Tensor(np.array(vs)).to(device)
                   for vs in zip(*experiences))
215
216            actions = (actions.long()).unsqueeze(dim=1)
217            rewards = rewards.unsqueeze(dim=1)
218            dones = dones.unsqueeze(dim=1)
219
220            if self.double_dqn:
221                target_q_values = self.double_q_learning_update(next_states, rewards, dones, self.
                       GAMMA, self.policy_net,
222                                                                 self.target_net)
223            else:
224                target_q_values = self.q_learning_update(next_states, rewards, dones, self.GAMMA, self
                       .target_net)
225            online_q_values = (self.policy_net(states).gather(dim=1, index=actions))
226            losses = F.mse_loss(online_q_values, target_q_values, reduction='none')
227            td_errors = torch.sqrt(losses)  # used for PER
228            is_weights_tensor = torch.tensor(np.array(is_weights), dtype=torch.float32, device=
                   device)
229            weighted_losses = losses * is_weights_tensor  # Apply IS weights
230            loss = weighted_losses.mean()
231            # updates the parameters of the online network
232            self.optimizer.zero_grad()
233            loss.backward()
234            self.optimizer.step()
235
236            if self.replay.use_per:
237                self.replay.update_priority(idxs, td_errors.cpu().detach().numpy()) #necessary?
238
239
240        def step(self, state, action, reward, next_state, done):
241            experience = Transition(state, action, reward, next_state, done)
242            self.replay.push(experience)
243            if not done:
244                self.number_timesteps += 1
245                # every so often the agent should learn from experiences
246                if self.number_timesteps % self.update_frequency == 0 and self.
                       has_sufficient_experience():
247
248                    batch, idxs, is_weights = self.replay.sample(self.batch_size)
249                    self.learn(experiences=batch, is_weights=is_weights, idxs=idxs)
250
251                if self.number_timesteps % self.update_target_frequency == 0:
252                    self.target_net.load_state_dict(self.policy_net.state_dict())
253
254        def train_for_at_most(self):
255            """Train agent for a maximum number of timesteps."""
256            state, info = self.env.reset()
257            state, _, _, _, _ = self.env.step(1)
258
259            self.number_lives = 5
260            score = 0
261            done = False
262            episode_timestep = 0
263            # for t in range(self.max_timesteps):
264            while not done:
265                action = self.choose_action(state)
266                # print(f"Action Dis: {action} Timestep: {episode_timestep}")
267                # action_cont = self.discrete2cont_action(action)
268                next_state, reward, done, truncated, info = self.env.step(action)
269                reward = min(1, reward)
270                if info.get("lives") < self.number_lives:
271                    self.number_lives = info.get("lives")
272                    self.step(state, action, reward, next_state, True)
273                    next_state, _, _, _, _ = self.env.step(1)
274
275                else:
276                    self.step(state, action, reward, next_state, done)
277                self.epsilon = np.interp(self.number_timesteps, [0, 500000], [1, 0.01])
```

```
278                     episode_timestep +=1
279                     state = next_state
280                     score += reward
281                     if done or truncated:
282                         print(f"Episode: {self.number_episodes} Timesteps {episode_timestep}  Died :("
                                )
283                         self.number_episodes += 1
284                         self.video = []
285                         break
286             if self.number_episodes % 200 == 0:
287                 print(f"Episode: {self.number_episodes} finished in {episode_timestep} timesteps
                        score: {score}")
288                 with open('prints.txt', 'a') as f:
289                     f.write(f"\nEpisode: {self.number_episodes} finished in {episode_timestep}
                            timesteps score: {score}")
290             return score
291
292         def train(self):
293             scores = []
294             target_score = float("inf")
295             most_recent_scores = deque(maxlen=100)
296             best_score = float("-inf")
297             self.policy_net.train()
298             self.target_net.train()
299             with open('prints.txt', 'w') as f:
300                 f.write("Starting prints")
301             for i in range(self.max_episodes):
302                 score = self.train_for_at_most()
303                 logger.log({'Score': score})
304                 scores.append(score)
305                 most_recent_scores.append(score)
306                 average_score = np.mean(most_recent_scores)
307                 logger.log({'Mean Score 100 Episodes': average_score})
308
309                 if average_score >= target_score or self.number_timesteps >= 4000000: # 3 million
                        episode limit
310                     print(f"\nEnvironment solved in {i:d} episodes!\tAverage Score: {average_score
                            :.2f}")
311                     checkpoint_filepath = f"rl_chk/new-dqn-per-checkpoint{self.number_episodes}.
                            pth"
312                     os.makedirs(os.path.dirname(checkpoint_filepath), exist_ok=True)
313                     self.save(checkpoint_filepath)
314                     break
315                 elif average_score > best_score:
316                     best_score = average_score
317                     plt.plot(average_score)
318                     plt.savefig("rewards.png")
319                     with open('prints.txt', 'a') as f:
320                         f.write("\nSaving checkpoint")
321                     print("Saving checkpoint")
322                     checkpoint_filepath = f"rl_chk/new-dqn-per-checkpoint_4mil.pth"
323                     self.save(checkpoint_filepath)
324                 if (i + 1) % 100 == 0:
325                     plt.plot(scores)
326                     plt.savefig("rewards.png")
327                     with open('prints.txt', 'a') as f:
328                         f.write(f"\n\rEpisode: {i + 1}\tAverage Score: {average_score:.2f} Epsilon
                                : {self.epsilon} N_Frames: {self.number_timesteps}")
329                     print(f"\rEpisode: {i + 1}\tAverage Score: {average_score:.2f} Epsilon: {self.
                            epsilon} N_Frames: {self.number_timesteps}")
330
331             return scores
332
333
334     def Preprocessing_env(env):
335
336         env = gym.wrappers.AtariPreprocessing(env, noop_max=30,
337                                           screen_size=84, terminal_on_life_loss=False,
338                                           grayscale_obs=True, grayscale_newaxis=False, scale_obs=
                                                    False)
339
340         env = gym.wrappers.FrameStack(env, 4)
341         return env
342
343     if "main":
344         # env = gym.make('CartPole-v1', render_mode="rgb_array")
345         env = gym.make("BreakoutNoFrameskip-v4", render_mode="rgb_array")
346         # env = gym.make('Hopper-v4')
347         env = Preprocessing_env(env)
348
349         wandb_logger = Logger(
350             f"PER-DQN-New",
351             project='INM707-Task2')
352         logger = wandb_logger.get_logger()
353
354         dqn = Agent(env, per=True, double=False, logger = logger)
355         scores = dqn.train()
```

```
356        plt.plot(scores)
357        plt.savefig("rewards.png")
358        # plt.show()
```

**Listing 4:** DQN and extensions code

### D.2.3   dqn_inference.py

```
1    import gymnasium as gym
2    from gymnasium.utils.save_video import save_video
3
4    import math
5    import random
6    import matplotlib
7    import matplotlib.pyplot as plt
8    from collections import namedtuple, deque
9    from itertools import count
10   import torch
11   import torch.nn as nn
12   import torch.optim as optim
13   import torch.nn.functional as F
14   import numpy as np
15
16   device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
17
18
19   class DQN(nn.Module):
20
21       def __init__(self, n_observations, n_actions, hidden_units=512):
22           super(DQN, self).__init__()
23           self.layer1 = nn.Linear(n_observations, hidden_units)
24           self.layer2 = nn.Linear(hidden_units, hidden_units)
25           self.layer3 = nn.Linear(hidden_units, n_actions)
26
27       def forward(self, x):
28           x = F.relu(self.layer1(x))
29           x = F.relu(self.layer2(x))
30           return self.layer3(x)
31
32
33   class DQNCNN(nn.Module): # DQN/DDQN
34       def __init__(self, input_shape, n_actions, hidden_units=512):
35           super(DQNCNN, self).__init__()
36           self.conv = nn.Sequential(
37               nn.Conv2d(input_shape[0], 32, kernel_size=8, stride=4),
38               nn.ReLU(),
39               nn.Conv2d(32, 64, kernel_size=4, stride=2),
40               nn.ReLU(),
41               nn.Conv2d(64, 64, kernel_size=3, stride=1),
42               nn.ReLU()
43           )
44           conv_out_size = self.get_conv_out_size(input_shape)
45
46           self.value = nn.Sequential(
47               nn.Linear(conv_out_size, hidden_units),
48               nn.ReLU(),
49               nn.Linear(hidden_units, n_actions)
50           )
51
52       def get_conv_out_size(self, shape):
53           conv_size = self.conv(torch.zeros(1, *shape))
54           return int(np.prod(conv_size.size()))
55
56       def forward(self, x):
57           conv_out = self.conv(x).view(x.size()[0], -1)
58           return self.value(conv_out)
59
60
61   class Agent:
62       def __init__(self, env, per=False, double=False, logger=None):
63           self.logger = logger
64           self.max_timesteps = 5000
65           self.number_timesteps = 0
66           self.number_episodes = 0
67           self.epsilon = 1
68           # Get number of actions from gym action space
69           self.env = env
70           self.n_actions = 4
71           self.number_lives = 5
72           # self.n_actions = env.action_space.shape[0]
73           # num_bins = 61  # Number of bins for each action dimension
74           # self.n_actions = num_bins ** self.n_actions
75           print(self.n_actions)
76           print(f"Number actions: {self.n_actions}")
77           seed = None
```

27

```python
78          self.random_state = np.random.RandomState() if seed is None else np.random.RandomState
                  (seed)
79
80          # Get the number of state observations
81          self.state, self.info = env.reset()
82          print(f"State shape: {self.state.shape}")
83          # self.n_observations = len(self.state)
84          self.n_observations = self.state.shape
85          checkpoint = torch.load(f"rl_chk/dqn-per-checkpoint_4mil.pth")
86          self.policy_net = DQNCNN(self.n_observations, self.n_actions, hidden_units=512).to(
                  device)
87          self.policy_net.load_state_dict(checkpoint['q-network-state'])
88
89          print(self.n_observations)
90          print(env.observation_space.shape)
91          # self.policy_net = DQN(env.observation_space.shape, self.n_actions).to(device)
92          # self.target_net = DQN(env.observation_space.shape, self.n_actions).to(device)
93
94          self.video = []
95      # def discrete2cont_action(self, action):
96      #     # Map the discrete action index to continuous torques
97      #     num_bins = 61
98      #     action_indices = np.unravel_index(action, (num_bins, num_bins, num_bins))
99      #     torque_min = -1.0
100     #     torque_max = 1.0
101     #     torques = [torque_min + (torque_max - torque_min) * idx / (num_bins - 1) for idx in
            action_indices]
102     #     return np.array(torques)
103
104     def choose_action(self, state):
105         # need to reshape state array and convert to tensor
106         state_tensor = (torch.from_numpy(np.array(state)).unsqueeze(dim=0).to(device)).float()
107         action = self.epsilon_greedy_policy(state_tensor, self.epsilon)
108         return action
109
110     def epsilon_greedy_policy(self, state, epsilon):
111         """With probability epsilon explore randomly; otherwise exploit knowledge optimally.
                """
112         action = self.greedy_policy(state)
113         return action
114
115     def uniform_random_policy(self, state):
116         """Choose an action uniformly at random."""
117         # random_vector = np.random.(low=-1, high=1, size=self.n_actions)
118         # return random_vector
119         return self.random_state.randint(self.n_actions)
120
121     def greedy_policy(self, state):
122         # print(state.shape)
123         # print(state.dtype)
124         """Choose an action that maximizes the action_values given the current state."""
125         action = (self.policy_net(state)
126                   .argmax()
127                   .cpu()  # action_values might reside on the GPU!
128                   .item())
129         return action
130
131     def select_greedy_actions(self, states, q_network):
132         _, actions = q_network(states).max(dim=1, keepdim=True)
133         # print(actions)
134         return actions
135
136     def step(self, state, action, reward, next_state, done):
137         if not done:
138             self.number_timesteps += 1
139
140     def train_for_at_most(self):
141         """Train agent for a maximum number of timesteps."""
142         state, self.info = self.env.reset()
143         score = 0
144         done = False
145         episode_timestep = 0
146         state, _, _, _, _ = self.env.step(1)
147         self.policy_net.eval()
148         with torch.no_grad():
149             for t in range(self.max_timesteps):
150         #     while not done:
151
152                 action = self.choose_action(state)
153                 next_state, reward, done, truncated, info = self.env.step(action)
154                 if info.get("lives") < self.number_lives:
155                     self.number_lives = info.get("lives")
156                     next_state, _, _, _, _ = self.env.step(1)
157
158                 self.video.append(self.env.render())
159                 self.step(state, action, reward, next_state, done)
160                 episode_timestep +=1
```

```
161                      state = next_state
162                      score += reward
163                      if done or truncated:
164                          print("GAME-OVER!")
165                          save_video(self.video, "videos", fps=25, name_prefix="video-inference")
166                          self.number_episodes += 1
167                          self.video = []
168                          break
169                  print(f"Episode {self.number_episodes} finished in {episode_timestep} timesteps score: {score}")
170                  if not done:
171                      print("TOO-LONG!")
172                      save_video(self.video, "videos", fps=25, name_prefix="video-inference")
173                      self.number_episodes += 1
174                      self.video = []
175          return score
176
177      def train(self):
178          scores = []
179          target_score = float("inf")
180          most_recent_scores = deque(maxlen=100)
181          score = self.train_for_at_most()
182          scores.append(score)
183          most_recent_scores.append(score)
184          return scores
185
186
187  def Preprocessing_env(env):
188
189      env = gym.wrappers.AtariPreprocessing(env, noop_max=30,
190                                  screen_size=84, terminal_on_life_loss=False,
191                                  grayscale_obs=True, grayscale_newaxis=False, scale_obs=
                                          False)
192
193      env = gym.wrappers.FrameStack(env, 4)
194      return env
195
196  if "main":
197      # env = gym.make('CartPole-v1', render_mode="rgb_array")
198      # env = gym.make('Hopper-v4', render_mode="rgb_array")
199      env = gym.make("BreakoutNoFrameskip-v4", render_mode="rgb_array")
200      env = Preprocessing_env(env)
201      dqn = Agent(env, per=False, double=True)
202      scores = dqn.train()
203      # plt.plot(scores)
204      # plt.savefig("rewards.png")
205      # plt.show()
```

**Listing 5:** DQN inference

### D.2.4   Individual Part

```
1   import gymnasium as gym
2   from gymnasium.utils.save_video import save_video
3
4   import math
5   import random
6   import matplotlib
7   import matplotlib.pyplot as plt
8   from collections import namedtuple, deque
9   from itertools import count
10  import torch
11  import torch.nn as nn
12  import torch.optim as optim
13  import torch.nn.functional as F
14  import numpy as np
15  import os
16  from buffer import ReplayMemory
17  from logger import Logger
18
19  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
20  print(f"Device is {device}")
21
22  Transition = namedtuple('Transition',
23                          ('state', 'action', 'reward', 'next_state', 'done'))
24
25  # os.environ['https_proxy'] = "http://hpc-proxy00.city.ac.uk:3128"
26
27
28  class DQN(nn.Module):
29
30      def __init__(self, n_observations, n_actions, hidden_units=512):
31          super(DQN, self).__init__()
32          self.layer1 = nn.Linear(n_observations, hidden_units)
33          self.layer2 = nn.Linear(hidden_units, hidden_units)
```

```python
34            self.layer3 = nn.Linear(hidden_units, n_actions)

36      def forward(self, x):
37          x = F.relu(self.layer1(x))
38          x = F.relu(self.layer2(x))
39          return self.layer3(x)


42  class DQNCNN(nn.Module):  # DQN/DDQN
43      def __init__(self, input_shape, n_actions, hidden_units=512):
44          super(DQNCNN, self).__init__()
45          self.conv = nn.Sequential(
46              nn.Conv2d(input_shape[0], 32, kernel_size=8, stride=4),
47              nn.ReLU(),
48              nn.Conv2d(32, 64, kernel_size=4, stride=2),
49              nn.ReLU(),
50              nn.Conv2d(64, 64, kernel_size=3, stride=1),
51              nn.ReLU()
52          )
53          conv_out_size = self.get_conv_out_size(input_shape)

55          self.value = nn.Sequential(
56              nn.Linear(conv_out_size, hidden_units),
57              nn.ReLU(),
58              nn.Linear(hidden_units, n_actions)
59          )

61      def get_conv_out_size(self, shape):
62          conv_size = self.conv(torch.zeros(1, *shape))
63          return int(np.prod(conv_size.size()))

65      def forward(self, x):
66          conv_out = self.conv(x).view(x.size()[0], -1)
67          return self.value(conv_out)


70  class Agent:
71      def __init__(self, env, per=False, double = False, logger = None):

73          self.logger = logger
74          self.GAMMA = 0.99
75          self.LR = 1e-4
76          self.ALPHA = 1
77          self.update_frequency = 4
78          self.update_target_frequency = 10000 # 20k for tuned ddqn
79          self.batch_size = 64
80          self.per = per
81          self.double_dqn = double

83          self.replay = ReplayMemory(100000, use_per=self.per)
84          if self.per:
85              self.alpha = self.replay.alpha
86              self.sum_tree = self.replay.sum_tree
87              self.max_priority = self.replay.max_priority
88          self.memory = self.replay.memory

90          self.max_episodes = 5000
91          self.number_episodes = 0
92          self.max_timesteps = 2000
93          self.number_timesteps = 0
94          self.epsilon = 1

96          # Get number of actions from gym action space
97          self.env = env
98          self.n_actions = 4
99          self.number_lives = 5
100         # self.n_actions = env.action_space.shape[0]
101         # num_bins = 61  # Number of bins for each action dimension
102         # self.n_actions = num_bins ** self.n_actions
103         print(self.n_actions)
104         print(f"Number actions: {self.n_actions}")
105         seed = None
106         self.random_state = np.random.RandomState() if seed is None else np.random.RandomState(seed)

108         # Get the number of state observations
109         self.state, self.info = env.reset()
110         print(f"State shape: {self.state.shape}")
111         # self.n_observations = len(self.state)
112         self.n_observations = self.state.shape
113         self.policy_net = DQNCNN(self.n_observations, self.n_actions, hidden_units=512).to(device)
114         self.target_net = DQNCNN(self.n_observations, self.n_actions, hidden_units=512).to(device)
115         self.target_net.load_state_dict(self.policy_net.state_dict())
116         self.optimizer = optim.AdamW(self.policy_net.parameters(), lr=self.LR, amsgrad=True)
117         print(self.n_observations)
```

```python
            print(env.observation_space.shape)
            # self.policy_net = DQN(env.observation_space.shape, self.n_actions).to(device)
            # self.target_net = DQN(env.observation_space.shape, self.n_actions).to(device)

            self.video = []

    # def discrete2cont_action(self, action):
    #     # Map the discrete action index to continuous torques
    #     num_bins = 61
    #     action_indices = np.unravel_index(action, (num_bins, num_bins, num_bins))
    #     torque_min = -1.0
    #     torque_max = 1.0
    #     torques = [torque_min + (torque_max - torque_min) * idx / (num_bins - 1) for idx in
    #         action_indices]
    #     return np.array(torques)

    def has_sufficient_experience(self):
        """True if agent has enough experience to train on a batch of samples; False otherwise
            ."""
        # return len(self.memory) >= self.batch_size
        if len(self.memory) == 5000:
            print("Sufficient experience recently obtained !!!")
        return len(self.memory) >= 5000

    def has_full_experience(self):
        """True if agent has enough experience to train on a batch of samples; False otherwise
            ."""
        # return len(self.memory) >= self.batch_size
        if len(self.memory) == 100000:
            return len(self.memory) >= 100000

    def save(self, filepath):
        checkpoint = {
            "q-network-state": self.policy_net.state_dict(),
            "optimizer-state": self.optimizer.state_dict(),
        }
        torch.save(checkpoint, filepath)

    def choose_action(self, state):
        # print(state.shape)
        # need to reshape state array and convert to tensor
        state_tensor = (torch.from_numpy(np.array(state)).unsqueeze(dim=0).to(device)).float()
        # choose uniform at random if agent has insufficient experience
        if not self.has_sufficient_experience():
            action = self.uniform_random_policy(state_tensor)
        else:
            # print("Sufficient experience")
            action = self.epsilon_greedy_policy(state_tensor, self.epsilon)
        return action

    def epsilon_greedy_policy(self, state, epsilon):
        """With probability epsilon explore randomly; otherwise exploit knowledge optimally.
            """
        if self.random_state.random() < epsilon:
            action = self.uniform_random_policy(state)
        else:
            action = self.greedy_policy(state)
        return action

    def uniform_random_policy(self, state):
        """Choose an action uniformly at random."""
        # random_vector = np.random.(low=-1, high=1, size=self.n_actions)
        # return random_vector
        return self.random_state.randint(self.n_actions)

    def greedy_policy(self, state):
        # print(state.shape)
        # print(state.dtype)
        """Choose an action that maximizes the action_values given the current state."""
        action = (self.policy_net(state)
                    .argmax()
                    .cpu()  # action_values might reside on the GPU!
                    .item())
        return action

    def select_greedy_actions(self, states, q_network):
        _, actions = q_network(states).max(dim=1, keepdim=True)
        # print(actions)
        return actions

    def evaluate_selected_actions(self, states, actions, rewards, dones, gamma, q_network):
        """Compute the Q-values by evaluating the actions given the current states and Q-
            network."""
        next_q_values = q_network(states).gather(dim=1, index=actions)
        q_values = rewards + (gamma * next_q_values * (1 - dones))
        return q_values
```

```python
200        def q_learning_update(self, states,rewards,dones,gamma,q_network):
201            """Q-Learning update with explicitly decoupled action selection and evaluation steps.
                  """
202            actions = self.select_greedy_actions(states, q_network)
203            q_values = self.evaluate_selected_actions(states, actions, rewards, dones, gamma,
                  q_network)
204            return q_values
205
206        def double_q_learning_update(self, states,rewards,dones,gamma,q_network1, q_network2):
207            """Q-Learning update with explicitly decoupled action selection and evaluation steps.
                  """
208            actions = self.select_greedy_actions(states, q_network1)
209            q_values = self.evaluate_selected_actions(states, actions, rewards, dones, gamma,
                  q_network2)
210            return q_values
211
212        def learn(self, experiences, is_weights, idxs):
213            """Update the agent's state based on a collection of recent experiences."""
214            states, actions, rewards, next_states, dones = (torch.Tensor(np.array(vs)).to(device)
                  for vs in zip(*experiences))
215
216            actions = (actions.long()).unsqueeze(dim=1)
217            rewards = rewards.unsqueeze(dim=1)
218            dones = dones.unsqueeze(dim=1)
219
220            if self.double_dqn:
221                target_q_values = self.double_q_learning_update(next_states,rewards, dones,self.
                      GAMMA,self.policy_net,
222                                                                    self.target_net)
223            else:
224                target_q_values = self.q_learning_update(next_states,rewards,dones,self.GAMMA,self
                      .target_net)
225            online_q_values = (self.policy_net(states).gather(dim=1, index=actions))
226            losses = F.mse_loss(online_q_values, target_q_values, reduction='none')
227            td_errors = torch.sqrt(losses)  # used for PER
228            is_weights_tensor = torch.tensor(np.array(is_weights), dtype=torch.float32, device=
                  device)
229            weighted_losses = losses * is_weights_tensor  # Apply IS weights
230            loss = weighted_losses.mean()
231            # updates the parameters of the online network
232            self.optimizer.zero_grad()
233            loss.backward()
234            self.optimizer.step()
235
236            if self.replay.use_per:
237                self.replay.update_priority(idxs, td_errors.cpu().detach().numpy()) #necessary?
238
239
240        def step(self, state, action, reward, next_state, done):
241            experience = Transition(state, action, reward, next_state, done)
242            self.replay.push(experience)
243            if not done:
244                self.number_timesteps += 1
245                # every so often the agent should learn from experiences
246                if self.number_timesteps % self.update_frequency == 0 and self.
                      has_sufficient_experience():
247
248                    batch, idxs, is_weights = self.replay.sample(self.batch_size)
249                    self.learn(experiences=batch, is_weights=is_weights, idxs=idxs)
250
251                if self.number_timesteps % self.update_target_frequency == 0:
252                    self.target_net.load_state_dict(self.policy_net.state_dict())
253
254        def train_for_at_most(self):
255            """Train agent for a maximum number of timesteps."""
256            state, info = self.env.reset()
257            state, _, _, _, _ = self.env.step(1)
258
259            self.number_lives = 5
260            score = 0
261            done = False
262            episode_timestep = 0
263            # for t in range(self.max_timesteps):
264            while not done:
265                action = self.choose_action(state)
266                # print(f"Action Dis: {action} Timestep: {episode_timestep}")
267                # action_cont = self.discrete2cont_action(action)
268                next_state, reward, done, truncated, info = self.env.step(action)
269                reward = min(1, reward)
270                if info.get("lives") < self.number_lives:
271                    self.number_lives = info.get("lives")
272                    self.step(state, action, reward, next_state, True)
273                    next_state, _, _, _, _ = self.env.step(1)
274
275                else:
276                    self.step(state, action, reward, next_state, done)
277                self.epsilon = np.interp(self.number_timesteps, [0, 500000], [1, 0.01])
```

```python
278                    episode_timestep +=1
279                    state = next_state
280                    score += reward
281                    if done or truncated:
282                        print(f"Episode: {self.number_episodes} Timesteps {episode_timestep}  Died :("
                            )
283                        self.number_episodes += 1
284                        self.video = []
285                        break
286            if self.number_episodes % 200 == 0:
287                print(f"Episode: {self.number_episodes} finished in {episode_timestep} timesteps
                        score: {score}")
288                with open('prints.txt', 'a') as f:
289                    f.write(f"\nEpisode: {self.number_episodes} finished in {episode_timestep}
                            timesteps score: {score}")
290            return score
291
292        def train(self):
293            scores = []
294            target_score = float("inf")
295            most_recent_scores = deque(maxlen=100)
296            best_score = float("-inf")
297            self.policy_net.train()
298            self.target_net.train()
299            with open('prints.txt', 'w') as f:
300                f.write("Starting prints")
301            for i in range(self.max_episodes):
302                score = self.train_for_at_most()
303                logger.log({'Score': score})
304                scores.append(score)
305                most_recent_scores.append(score)
306                average_score = np.mean(most_recent_scores)
307                logger.log({'Mean Score 100 Episodes': average_score})
308
309                if average_score >= target_score or self.number_timesteps >= 4000000: # 3 million
                        episode limit
310                    print(f"\nEnvironment solved in {i:d} episodes!\tAverage Score: {average_score
                        :.2f}")
311                    checkpoint_filepath = f"rl_chk/new-dqn-per-checkpoint{self.number_episodes}.
                        pth"
312                    os.makedirs(os.path.dirname(checkpoint_filepath), exist_ok=True)
313                    self.save(checkpoint_filepath)
314                    break
315                elif average_score > best_score:
316                    best_score = average_score
317                    plt.plot(average_score)
318                    plt.savefig("rewards.png")
319                    with open('prints.txt', 'a') as f:
320                        f.write("\nSaving checkpoint")
321                    print("Saving checkpoint")
322                    checkpoint_filepath = f"rl_chk/new-dqn-per-checkpoint_4mil.pth"
323                    self.save(checkpoint_filepath)
324                if (i + 1) % 100 == 0:
325                    plt.plot(scores)
326                    plt.savefig("rewards.png")
327                    with open('prints.txt', 'a') as f:
328                        f.write(f"\n\rEpisode: {i + 1}\tAverage Score: {average_score:.2f} Epsilon
                            : {self.epsilon} N_Frames: {self.number_timesteps}")
329                    print(f"\rEpisode: {i + 1}\tAverage Score: {average_score:.2f} Epsilon: {self.
                        epsilon} N_Frames: {self.number_timesteps}")
330
331            return scores
332
333
334    def Preprocessing_env(env):
335
336        env = gym.wrappers.AtariPreprocessing(env, noop_max=30,
337                                        screen_size=84, terminal_on_life_loss=False,
338                                        grayscale_obs=True, grayscale_newaxis=False, scale_obs=
                                            False)
339
340        env = gym.wrappers.FrameStack(env, 4)
341        return env
342
343    if "main":
344        env = gym.make("AtlantisNoFrameskip-v4", render_mode="rgb_array")
345        env = Preprocessing_env(env)
346
347        wandb_logger = Logger(
348            f"logger-DQN-Atlantis",
349            project='INM707-Task2')
350        logger = wandb_logger.get_logger()
351
352        dqn = Agent(env, per=True, double=False, logger = logger)
353        scores = dqn.train()
354        plt.plot(scores)
355        plt.savefig("rewards.png")
```

```
356        # plt.show()
```

**Listing 6:** Individual Part: dqn.py