

CS-464 Term Project Final Report

Breakout Atari™ Game with Reinforcement Learning

Group 4

Abdullah Arda Aşçı (21702748), Alim Toprak Fırat (21600587),
Atahan Yorgancı (21702349), Tuna Alikasıfoğlu (21702125)

I. INTRODUCTION

Breakout is an arcade game that was published and developed by Atari, where the player is in control of a paddle on the x -axis and tries to break down each brick by handling a ball. At each time the ball touches to a brick, that particular brick is destroyed and the player gains a point, the game is won if all the bricks are destroyed. A sample position from the game is provided in Figure 1.

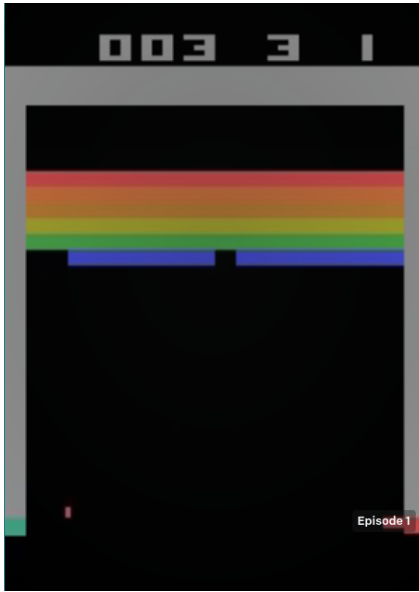


Fig. 1: Sample Position from the Breakout Atari Game [1]

In order to obtain a model that can play the game and get better results than a random agent with the use of Supervised or Unsupervised Learning, a large amount of data needs to be generated. Reinforcement learning doesn't need any data. In

reinforcement learning, an agent who has a set of actions is trained in an environment using rewards or punishments considering the decisions of the agent in different situations.

In reinforcement learning, there is a *environment*, which have different states changes at each time step, $s_t \in S$, and a learning *agent* who can observe the environment and then learns from the outcomes of it's actions, $a_t \in A$. In reinforcement learning problems, the agents' decisions influence their own actions in proceeding steps, agents don't have any information of which actions to take in what situation and they try to maximize reward (or minimize punishment) by trying out different combinations in the set of actions.

One of the main challenges of RL is to build a model that can learn from the noisy and delayed environmental data. The latency between the incoming data forces training model to give a delayed input reaction and the delays add up to each other in a cascading manner. To overcome the delay problem significantly, rather than using computer vision to gather frames and analyze them one by one, a breakout game developed in such a way that all in-game data can be accessible to the model.

Another challenge is the incoming state data being highly correlated with the previously acquired ones, rather than being independent. To build a model which can perform well withing the high-dimensional correlated state space, a Deep Q-Network (DQN) is developed which combines reinforcement learning with deep neural networks. The many layered approach made it possible to build up more abstract representations of the correlated spatial data.

To be able to train our model, the action set, state space and rewards were defined for the Breakout game. The positive reward are defined as unit positive reward for every break of a brick. No negative rewards are defined. However, in order to stop the repetitive behavior, such as bouncing the ball of the wall and getting it back without breaking any bricks and exploiting the positive rewards, a random restart protocol is implemented, so that the model is pushed to explore rather than exploit.

II. PROBLEM DESCRIPTION

The problem that we have is the utilization of deep reinforcement learning, with the specific model of convolutional neural networks, trained with Deep Q-Learning, whose input is raw pixel values, varying between $\{0, 1, \dots, 255\}$, of a grayscale version of the image of the current state, and output is the most suitable action from the action space. The question that we would like to address is whether the trained agent can outperform a random agent, which takes random actions at each time step, and if so whether it can outperform a human agent that plays the game in the same settings.

III. METHODS

As our project is to create a model with reinforcement learning, specifically DQN, we do not have a predefined dataset. With this in mind, our methodology can be separated into three parts.

A. Reinforcement Learning

Reinforcement Learning (RL) is a subset of Machine Learning (ML) where the aim is to teach a model, called agent, via its interactions with the surrounding environment. This method of learning requires no set of labeled or unlabeled data to be collected before the learning actually starts. Instead the agent, typically a neural network, is used for predicting the optimal action to be taken at each step based on its observation and a reward is determined by the environment which is used for training the agent. In RL, this environment is modeled as a Markov Decision Process (MDP).

1) *Markov Decision Process*: MDP is a mathematical framework based on Markov Chains for decision making processes with inherent randomness. In Markov Decision Processes, we define:

- S : State Space (finite set),
- A : Action Space (finite set),
- A_s : Set of actions available at state s ,
- $P_{ss'}^a = P_a(s, s') = P(s_{t+1} = s' \mid s_t = s, a_t = a)$ is the probability of action a in state s leading to state s' ,
- $R_a(s, s')$: Reward received after transitioning from s to s' .

2) *Markov Property*: For RL agents to work with MDPs we need the environment to be fully observable, meaning that state s must capture all the characteristics of the environment. In more technical terms, any state s must satisfy the Markov Property which is defined as.

$$P[s_{t+1} \mid s_t] = P[s_{t+1} \mid s_1, \dots, s_t] \quad (1)$$

This property essentially enables the environment to be memoryless which is required for Markov Chains and more importantly its extension MDPs.

3) *Policy*: The objective in an MDP is to optimize the *policy* of the decision making algorithm. Here, we define the function $\pi(s)$ that outputs the action chosen by the decision maker based on the current state s . This optimization mainly done by maximizing the cumulative reward function. This function can be expressed as:

$$G_t = \sum_{t=0}^{\infty} \gamma^t R_a(s_t, s_{t+1}), \quad (2)$$

where $0 \leq \gamma \leq 1$ is the discount factor. The equation above introduces the concept of *discount factors*. This parameter is quite important as it is one of the hyperparameters of RL training loops. It is useful for avoiding cyclic behavior and infinite returns, and representing an exponentially increasing uncertainty for the future time steps.

Moreover, the policy that maximizes the function given above is regarded as the *optimal policy* and denoted as $\pi^*(s)$. It should be noted that this optimal policy is not necessarily unique.

4) *State-Value Function*: The state-value function, or just value function, is denoted by $V_\pi(s)$. It is the expected return starting from state s and following policy π of an MDP. In most applications, it is used to evaluate how good being in a state is. It is mathematically expressed as:

$$V_\pi(s) = E_\pi[G_t \mid s_t = s] \quad (3)$$

As can be seen in the equation above, calculating the cumulative reward function is required to find the value function of a state. This can be decomposed into a recursive function as the current reward plus the discounted value function of the successor by utilizing the Bellman Equation.

$$V_\pi(s) = E_\pi[R_{t+1} + \gamma V_\pi(s_{t+1}) \mid s_t = s] \quad (4)$$

With this done, we now define the state-value function that is produced by the optimal policy π^* as the *optimal state-value function*. In mathematical terms:

$$V_*(s) = \max_\pi V_\pi(s) \quad (5)$$

5) *Action Value Function*: The action value function, also called the Q-function, is the expected return starting from state s , taking action a , and then following policy π . This is expressed as:

$$Q_\pi(s, a) = E_\pi[G_t \mid s_t = s, a_t = a] \quad (6)$$

Similar to the state-action function, we can also decompose this into a recursive function:

$$Q_\pi(s, a) = E_\pi[G_{t+1} \mid s_t = s, a_t = a] \quad (7)$$

where $G_{t+1} = R_{t+1} + \gamma Q_\pi(s_{t+1}, a_{t+1})$, and define the optimal action-value function as:

$$Q_*(s, a) = \max_\pi Q_\pi(s, a) \quad (8)$$

6) *Finding the Optimal Policy*: In all MDPs, three conditions are satisfied:

- An optimal policy π^* exists (not necessarily unique),
- Optimal policy achieves the optimal state-value function

$$V_{\pi^*}(s) = V_*(s) \quad (9)$$

- Optimal policy achieves the optimal action-value function

$$Q_{\pi^*}(s, a) = Q_*(s, a) \quad (10)$$

These three assumptions can be used to show that finding the optimal policy π^* to solve the MDP can be done by maximizing over $Q_*(s, a)$ with:

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } \arg \max_a Q_*(s, a) \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

This equation essentially means that finding the optimal policy can be done by following the path given by $Q_*(s, a)$ assuming that the optimal Q function is known.

B. Neural Networks

IV. RESULTS

See Figure 2 in section B.

The results of applying the methods to the data set. Include the list of questions your experiments are designed to answer Details of the experiments; observations.

V. DISCUSSION

Interpretation and discussion of the results.

VI. CONCLUSION

All in all, we successfully implemented a deep reinforcement learning architecture using convolutional neural networks, trained with Deep Q-Learning, whose input is raw pixel values, varying between $\{0, 1, \dots, 255\}$, of a grayscale version of the image of the current state, and output is the most suitable action from the action space. As we demonstrated in our final presentation demo, our trained agent outperforms the random agent with less than an hour of training. As demonstrated, our own group members also played the game using the keyboard agent, although none of the group members can perform more than 20 in any of the trials. In this context, the trained agent outperforms every group member with the score of 28. Therefore, the answer to our initial question is yes! The trained agent with the specified model can outperform both random and human agents.

Throughout the project, we learned the basic approaches of reinforcement learning. We learned that by using CNNs and DQNs, it is possible to learn control policies directly from high-dimensional sensory input using reinforcement learning. The resultant agent is satisfactory to respond our question

with “yes”. In addition to basic theoretical approach, we also had the chance to interact with a practical issue that we were not expecting to face. While we were trying to adjust the exploration versus exploitation hyperparameter, ϵ , we faced with an issue that we later learned it is called the credit assignment problem. The credit assignment problem means that it can happen that we choose an action, and we only win or lose hundreds of actions later, leaving us with no idea of as to which of our actions led to this win or lose, thus, making it difficult to learn from our actions [2]. We learned that solely increasing the training duration does increase the performance of the agent. We learned that in order to overcome this credit assignment pitfall, using random restarts actually increases the performance of the agent, by adjusting the exploration versus exploitation problem in the direction of exploration.

As the future work of the project, performance of the agent can be increased by changing the deep neural network architecture and increasing the training duration. During our project, we have encountered with more complex approaches like double DQNs, Never Give-Up (NGU), etc. Experimentation on these approaches can be expressed as potential future work. Furthermore, these approaches can be generalized as *DeepMind’s Agent57* to provide trained agents for 57 atari games [3].

REFERENCES

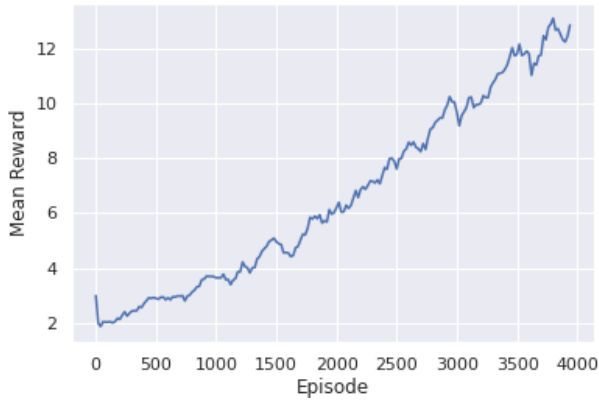
- [1] OpenAI. (Oct. 2019). Gym, [Online]. Available: <https://gym.openai.com/envs/Breakout-v0/>. [Accessed: Mar. 3, 2021].
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, *Playing atari with deep reinforcement learning*, 2013. arXiv: [1312.5602](https://arxiv.org/abs/1312.5602) [cs.LG].
- [3] A. P. Badia, B. Piot, S. Kapturowski, P. Sprechmann, A. Vitvitskyi, Z. D. Guo, and C. Blundell, “Agent57: Outperforming the atari human benchmark,” *CoRR*, vol. abs/2003.13350, 2020. arXiv: [2003.13350](https://arxiv.org/abs/2003.13350). [Online]. Available: <https://arxiv.org/abs/2003.13350>.

APPENDIX A
CONTRIBUTION

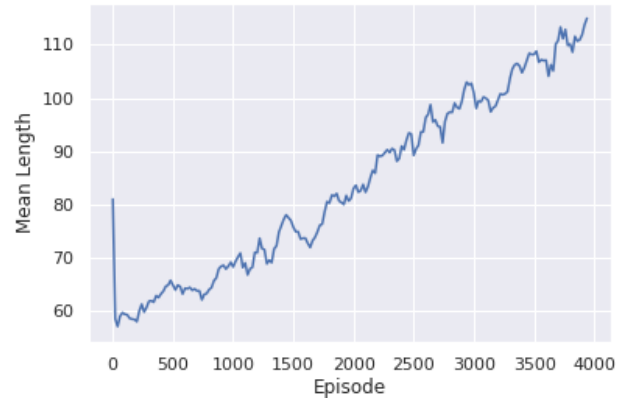
TABLE I: Task Sharing

Student	Task
Abdullah Arda Aşçı	Setting OpenAI, and developing wrappers for gym environment. Training & testing with different hyperparameters.
Atahan Yorgancı	Background research about RL and DQN. Implementation of artificial agent. Cloud based training training & testing with different hyperparameters.
Alim Toprak Fırat	Background research about RL and Q-Learning. Training & testing with different hyperparameters.
Tuna Alikaşifoğlu	Developing CLI for running, and training using gym environment. Implementation of DQN. Training & testing with different

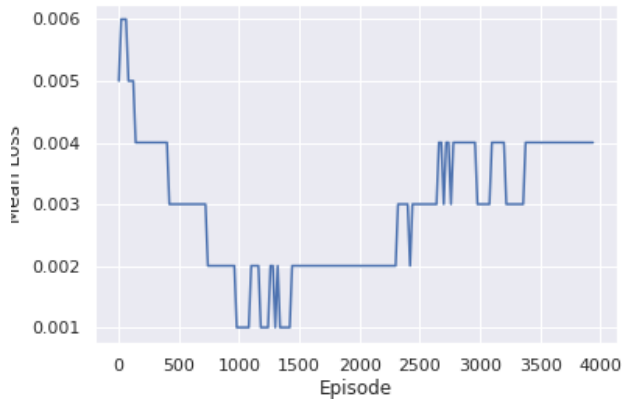
APPENDIX B RESULT PLOTS



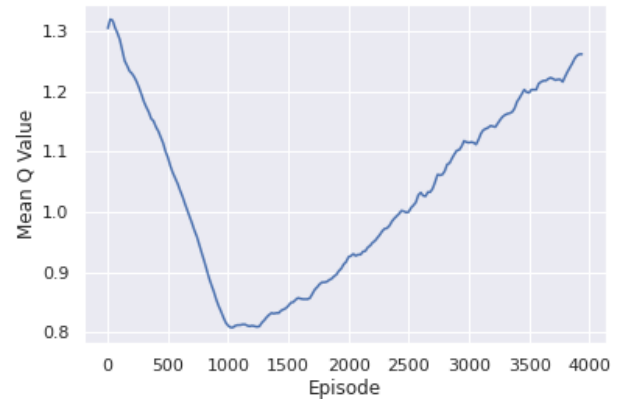
(a) Mean Reward



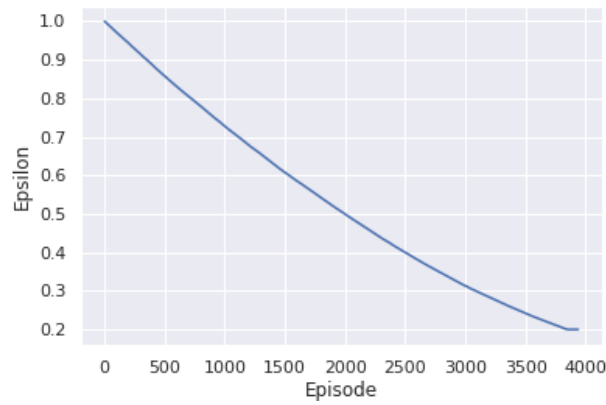
(b) Mean Length



(c) Mean Loss



(d) Mean Q Value



(e) Epsilon

Fig. 2: Change per Episode