# Arm Manipulation with Deep RL

Esraa Magdy

**Abstract**—This paper presents a Deep learning based approach to control the motion of a robotic arm using images from a single camera sensor. This approach utilizes a Deep Q Network (DQN) to train the robotic arm in a simulation environment.

**Index Terms**—Robot, IEEEtran, Udacity, LATEX, RL.

✦

## 1 REWARD FUNCTIONS

There are four rewards that the arm can have during any episode:

1- A reward for hitting the ground : This is a simple reward loss minus a factor that represents how early the robot hits the ground to encourage the robot to not hit it early.

$$rewardHistory = REWARD\_LOSS * 10-$$

$$((maxEpisodeLengthepisodeFrames)/maxEpisodeLength)$$

2- A reward loss for when the episode terminates. 3- A reward win when the right part of the arm touches the object, This was the same for both objective functions. It's simply a reward win times a multiplier plus a factor that is inversely proportional to how many frames have already done before wining. this factor to encourage the robot to finish fast.

$$rewardHistory = REWARD\_WIN * 100$$

$$+(maxEpisodeLength - episodeFrames) * 100$$

4- the last reward function is the interim reward function: For task 1, it was a simple smoothed average of the delta distance between the robot and the object times a multiplier as this average is usually a small number so it was better to add a multiplier. the delta distance will be a positive number if the robot is moving towards the object but a negative number if it's moving away from it. so this interim reward serves as a reward win and a reward loss at the same time based on in which direction the arm is moving.

$$avgGoalDelta = (avgGoalDelta * alpha)$$

$$+(distDelta * (1 - alpha))$$

Different values for alpha were experimented but at the end the value for alpha was 0.3, that means that the recent values for delta distance were more important than the history. the reward function was defined as :

$$rewardHistory = tanh(avgGoalDelta) * 100$$

For task2, Multiple rewards were experimented but finally a simple reward function was able to make it. The reward function used here was first presented in [1]. And states that:

$$if delta distance > 0$$

$$rewardHistory = REWARD\_WIN * 10 + avgGoalDelta$$

$$if delta distance < 0$$

$$rewardHistory = REWARD\_LOSS * 10$$

$$else reward History = 0.$$

Other reward functions could achieve a similar results but maybe with different hyper parameters . It depends on how fast the robot discovers the target object and learns that touching it results in a great reward. As long as the reward function is logically right, the robot will be able to achieve the objectives but it might take long time. More complex reward functions might produce the same results faster.

## 2 HYPERPARAMETERS

A starter guess for the hyperparameters was based on fruit and catch examples.The hyperparameters were the same for task one and task2.

```
#define INPUT_WIDTH    64
#define INPUT_HEIGHT   64
#define OPTIMIZER "RMSprop"
#define LEARNING_RATE 0.01f
#define REPLAY_MEMORY 10000
#define BATCH_SIZE 128
#define USE_LSTM true
#define LSTM_SIZE 256
#define NUM_ACTIONS 6
```

Fig. 1. HyperParams

The width and the height was set to 64 as Gazebo reported that this was the size of the image. beside a square image would reduce the complexity for the calculations. The Optimizer used was RMSprop. Other optimizers like SGD or Adam could have been used as well. The learning rate was set to 0.01. Usually increasing the batch size have the same effect as decreasing the learning rate, so instead of decreasing the learning rate below 0.01, the batch size

was increased to 128. the agent was learning slowly but it manged to converge smoothly after some time. Increasing the LSTM size was a good practice to help it remember the recent good moves.

## 3 RESULTS

### 3.1 Task1

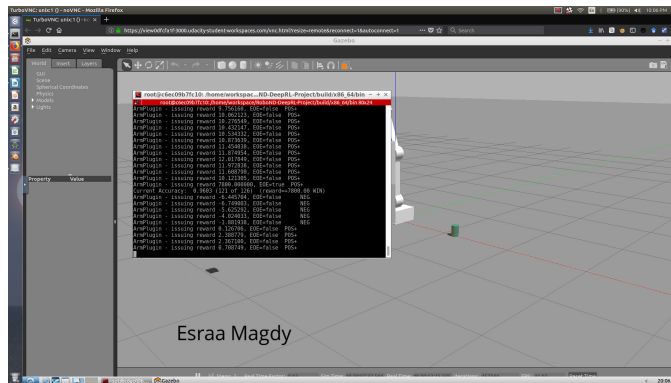For task1, Velocity control was first used on the Udacity workspace environment and achieved 96% accuracy.



Fig. 2. Velocity Control

Then, position control was experimented. Position control movements are smoother than velocity control movements but it took more time to reach a good accuracy.
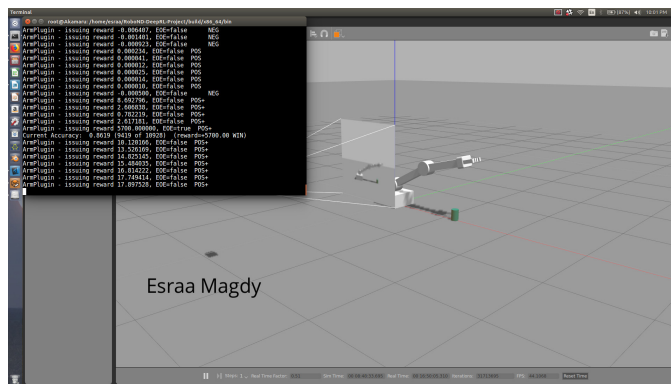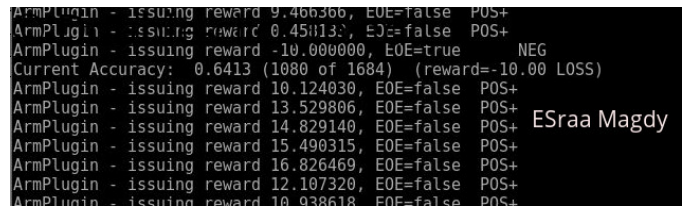


Fig. 3. Position Control

### 3.2 Task2

At first the arm couldn't get beyond 65%.That was odd but might be a result to the exploration and trying random values even after converging to the good move. Then After some tweaks with the reward function and decreasing EPS_End to 0.01 the arm reached the required accuracy but the convergence was really slow. For task2, position control was better in that case as gripper base collision depends on reaching a precise position.

## 4 FUTURE WORK

Further improvements could concentrate on how to implement a complex reward function that would make the robot
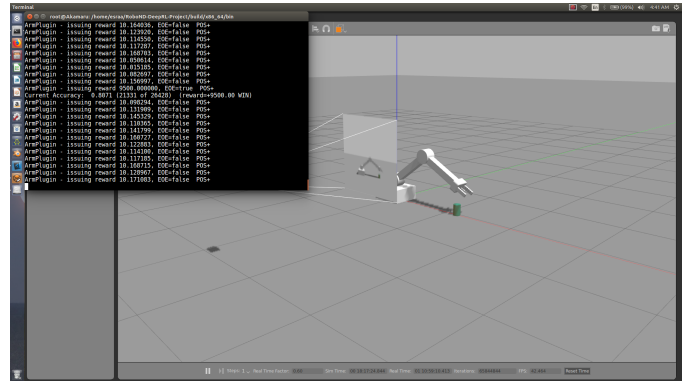


Fig. 4. Results



Fig. 5. Task2 Results

touch the object faster to increase the convergence rate. Also, Other optimizers could be tested to see the effect of each optimizer on the convergence rate.

## 5 REFERENCES

[1]https://arxiv.org/pdf/1511.03791.pdf