

# Robotics Project : Deep RL Arm Manipulation

Shotaro Oyama

**Abstract**—The goal of this project is to create a DQN agent and define reward functions to teach a robotic arm to carry out two primary objectives: First, have any part of the robot arm touch the object of interest, with at least a 90% accuracy. Second, have only the gripper base of the robot arm touch the object, with at least a 80% accuracy.

**Index Terms**—Deep RL Arm Manipulation, Reinforcement Learning, DQN, Gazebo, Robotics Nanodegree, Udacity

## 1 BACKGROUND

DEEP Reinforcement Learning for Robotics is a paradigm shift. The basic idea is to let a robot figure out, through trial and error, the best way to achieve the goal. The learning framework is characterized by an agent learning to interact with its environment. At each time step, the agent receives the environment's state, and the agent must choose an appropriate action in response. One time step later, the agent receives a reward and a new state. All agents have the goal to maximize expected cumulative reward, or the expected sum of rewards attained over all time steps. In this report, with training DQN agent by the information from camera and sensors on the robotic arm, the arm become to be manipulated as the way how it is expected.

## 2 REWARDS

Rewards consists o of the calculation from the result on each time step, and the incident when terminating the trial.

### 2.1 Terminal step

The reward values provided on each incident are as Table 1.

TABLE 1  
Reward on Incident

Incident	Reward Value
Collision on Arm (Target)	+100.0
Collision on Arm (Non-Target)	-10.0
Ground Contact	-100.0
Episode Timeout	-100.0

### 2.2 Each time step

The reward needs to be defined so that the it will encourage the agent to accomplish its task. So the reward is defined based on the distance between gripper and the object. From the result of several trial, the rewards include the Smoothed Moving Average, or SMMA, and the distance itself. By using SMMA, it shapes the reward function and gives the agent gradual feedback to let it know that its doing better as it gets closer to the object of interest, that helps eliminate

fluctuations and allows for agent to follow prevailing trend to maximize its rewards, allowing the agent to learn faster. The reason why distance itself was also used as rewards, is that it could make the learning faster. Also, the distance was adjusted a little bit in order for the Arm to get more precise target position. The Reward function is defined as follows:

```

distDelta=lastGoalDistance-distGoal
alpha=0.2 f
avgGoalDelta=
    (avgGoalDelta*alpha)+(distDelta*(1-alpha))
rewardDist= -pow((1+distGoal),3)/20
rewardHistory=avgGoalDelta*4+rewardDist

```

## 3 HYPERPARAMETERS

After tweaking, the same hyperparameters was used in both tasks. The highlight is as follows.

TABLE 2  
HyperParameter

Hyperparameter	Value
INPUT WIDTH	64
INPUT HEIGHT	64
OPTIMIZER	RMSprop
LEARNING RATE	0.01f
REPLAY MEMORY	20000
BATCH SIZE	64
USE LSTM	true
LSTM SIZE	256
VELOCITY CONTROL	Off

- INPUT WIDTH & HEIGHT was decreased from 512 to 64, to reduce complexity of neural network
- "RMSprop" OPTIMIZER was selected in spite of ADAM, because of better results
- LEARNING RATE was set to 0.01 after tweaking several values especially for achieving the objective of task 2
- LSTM MEMORY, BATCH SIZE, and LSTM SIZE were set as much as possible, taking computer resource into account. It seemed make the Arm working more softly.
- VELOCITY CONTROL was turned off because it added more complexity to the arm, resulted poorer result.

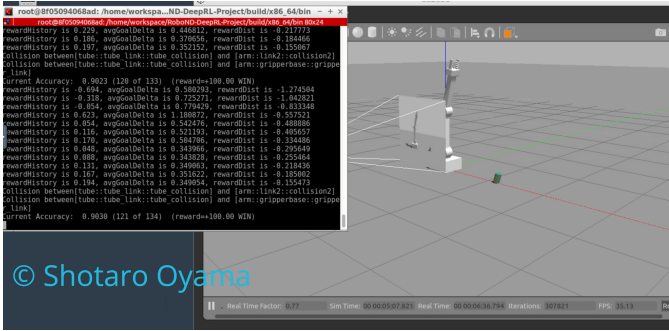


Fig. 1. Task1 result

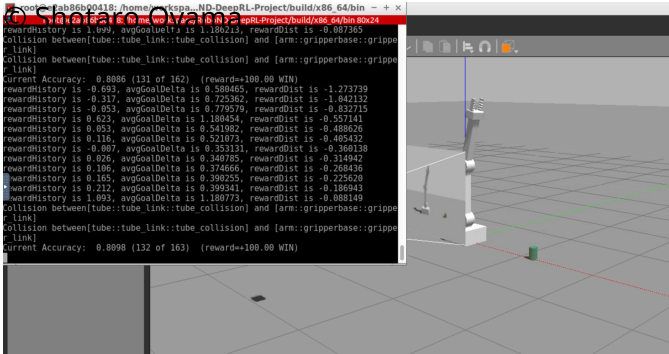


Fig. 2. Task2 result

## 4 RESULTS

Both tasks met the required accuracy. (Table. 2, Fig. 1 and Fig. 2)

TABLE 3  
Result Summary

Task	Accuracy
Task1	90.3% (121 of 134)
Task2	80% (132 of 163)

At Task 1, the arm was able to achieve the objective, almost from the beginning. However, at Task 2, the arm had failed many times in the first 10 episodes. Task 2 required more elaborated behavior of the arm, so it seemed taking time to learn that complicated movement.

## 5 FUTURE WORK

VELOCITY CONTROL may have possibility to improve the result of this project. At this time it was turned off because of its complexity, but if it is handled correctly, it may be able to realize more complicated robotic arm behavior.