# Arm Manipulator with Deep RL
# Project Write-up

Jorge Crespo Cedeño

No Institute Given

## 1 Introduction

This write-up describes the training of a arm manipulator using Deep Reinforcement Learning to touch a tube on the ground. In the first challenge, any part of the arm can touch the tube at least 80% of the times with a minimum of 100 runs. In the second one, only the gripper base can touch the tube 90% of the times with a minimum of 100 runs. In neither of both challenges, the gripper base cannot touch the ground.

## 2 Reward Functions

The DQN agent chooses an action. This action indicates the joint to be actuated upon and whether to increase or decrease the position of that joint. Alternatively, the joint can be controlled by increasing or decreasing the velocity, but position control is used for this project. The position of the selected joint is increased or decreased by 0.15.

The reward, for the first challenge, is given as follows:

- -0.4 when the gripper base touches the ground
- -0.4 when the number of frames captured by the camera exceed 100.
- +4.0 when any part of the arm touches the tube
- -a value which depends on the speed of the arm towards the tube. This is calculated as the smoothed moving average of the difference between the successive goal distances.

The reward, for the second challenge, is given as follows:

- -0.4 when the gripper base touches the ground
- -0.4 when the number of frames captured by the camera exceed 100.
- +4.0 when the gripper base touches the tube
- -0.4 when any part of the arm but the gripper base touches the ground.
- -a value which depends on the speed of the arm towards the tube. This is calculated as the smoothed moving average of the difference between the successive goal distances.

The magnitude of the negative reward is small compared to the positive one, because very small negative rewards, -4.0 for instance, make the agent hesitant to move since it tries to avoid a big penalty when not touching the tube. The positive reward is big compared to the magnitude of the negative reward in order to give a big incentive to the agent to learn from a successful arm movement.

## 3   Hyperparameters

The hyperparameters are the same for both objectives.

**Table 1.** Hyperparameters.

| Hyperparameter | Value | Reason |
|---|---|---|
| INPUT_WIDTH | 64 | The camera frames are 64 pixels wide |
| INPUT_HEIGHT | 64 | The camera frames are 64 pixels high |
| OPTIMIZER | RMSprop | Root mean square prop |
| LEARNING_RATE | 0.025 | A small value to avoid overshooting the minimum |
| REPLAY_MEMORY | 1000 | The amount of past experiences to store |
| BATCH_SIZE | 16 | Number of transitions randomly sampled from replay memory |
| USE_LSTM | true | Use Long Short Term Memory |
| LSTM_SIZE | 128 | Size of LSTM cell |

## 4   Results

The required accuracy for the first challenge was obtained after approximately 10000 runs.

An interesting case is the second challenge, which has almost the same hyperparameters as the first one, with the restriction of getting a positive reward only when the gripper touches the tube, and with addition of one negative reward when an arm part different than the gripper touches the tube. In this second case, the number of runs required to obtained at leat 80% accuracy is 2000 approximately. Althought the learning was stopped soon after reaching 81%, probably a 90% accuracy could have been obtained with less than 10000 runs. A plausible explanation of this case is the addition of a negative reward (when a part different than the gripper was touching the tube). This extra information might be the reason for the DQN agent to learn quicker a more precise function.

### 4.1   First Challenge

Fig. 1 shows the arm touching the tube with the required accuracy.

### 4.2   Second Challenge

Fig. 2 shows only the gripper base of the arm touching the tube with the required accuracy.
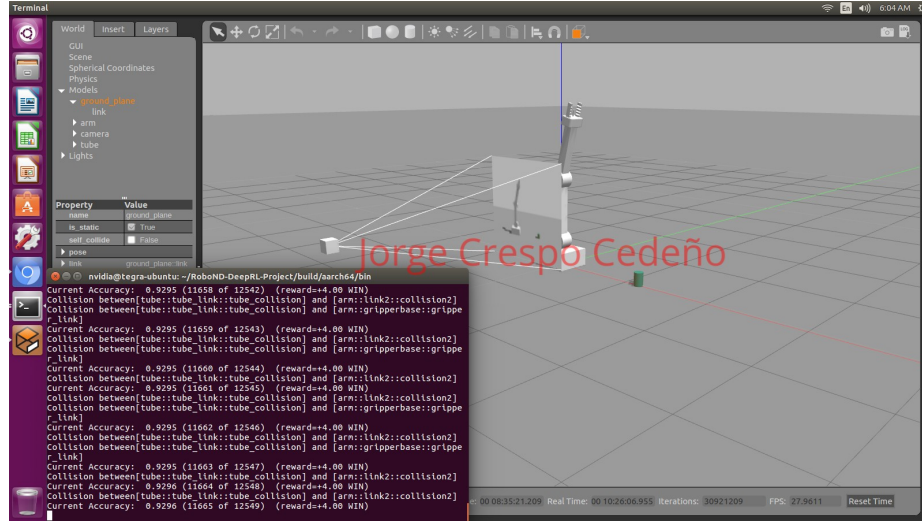
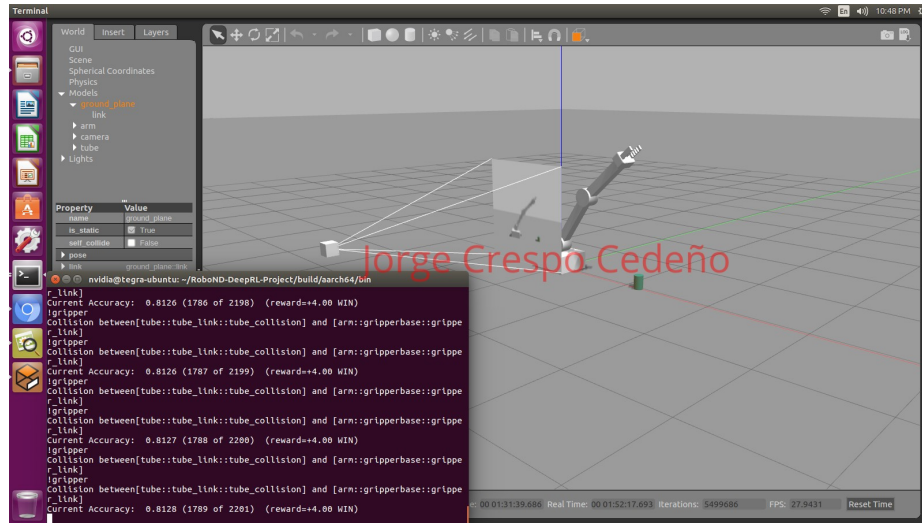**Fig. 1.** Any part of the arm manipulator touching the tube.



**Fig. 2.** Only the gripper of the arm manipulator touching the tube.

## 5    Future Work

In the case of the first challenge, although the reached accuracy was stable, i.e., the agent continue winning after 90% of success rate, the training took several hours. This time can be reduce by experimenting with bigger learning rates, since a small learning rate makes the agent learn slowly. In the case of the second challenge, the 80% required accuracy was reached much faster, in less than two hours approximately, with the same learning rate. In this case, bigger batch sizes can be tried to allow the agent break the correlation of successive experiences. The idea is that the agent learns several ways, i.e., arm configurations, to touch the gripper.