

# Reinforcement Learning for multi-mobile manipulators cooperation with obstacle avoidance

\*

1<sup>st</sup> Hotae Lee

Department of Mechanical & Aerospace engineering  
Seoul National University  
Seoul, South Korea  
hotae319@snu.ac.kr

**Abstract**—Multi-robot’s cooperation to carry an object to a target problem is a challenging problem because path planning should be computed for each other differently. This problem consists of several sub-problems like navigation, cooperation, and obstacle avoidance. In this paper, I apply DQN(Deep Q-network) to solving this problem. To optimize the learning method for navigation problem, this work uses curriculum learning and develops revised mini-batch sampling method which always includes a success state. Moreover, a centralized multi-robot problem has a curse of dimensionality, which means the size of action space increases exponentially. This work uses the method which reduces the action space through allocating one step for one agent. It reduces the computation time for network updates and predictions. Simulation results show that revised sampling 'discounted reward sum is more stable than previous random sampling's thing. In addition, curriculum learning yields faster and more stable convergence compared to normal learning. Finally, reduction of action space lowers the computation time, however, it makes the reward sum graphs fluctuate unstably.

**Index Terms**—Navigation, Cooperation constraint, Obstacle avoidance, DQN, Curriculum learning, Sampling method, Reduced action space

## I. INTRODUCTION

There have been rapid developments of cooperative multi-robot systems in recent years. Search and Rescue, collision avoidance, and formation control are main challenges. Also, as described in Fig. 1, mobile manipulators cooperation to carry the objects is very promising for many applications [1]. Unlike a single robot, multi-robots which grasp the same object require many complex considerations in order to keep objects not to drop and avoid the obstacle. Each robot has different strategy and planning because they have each different states such as pose and distance from obstacles. Not only how the object’s trajectory should be set, but also how each robot should move are significant issues. A lot of previous research about robot cooperation focused on tracking the given trajectory of CoM(Center of Mass)of the object or real-time solving optimization problem while the robots are moving [2]. But, the former requires a priori knowledge of trajectory and it is not flexible for various situations. The latter requires an

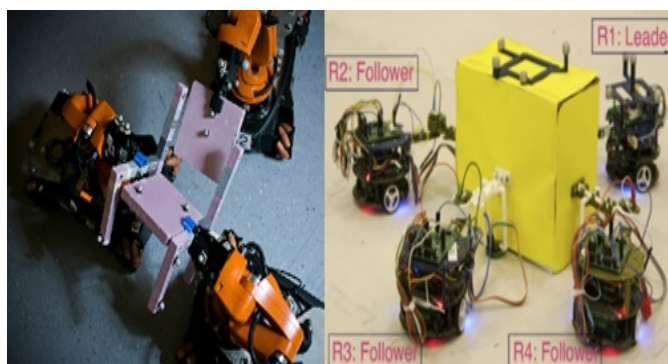


Fig. 1. Mobile manipulators cooperation to carry the objects

exact environment’s information and this requires a lot of computation time to judge the situation and decide its movements in real-time. In recent years, reinforcement learning provides a good solution for this sequential decision-making problem. This can help us already know the optimal policy over every instances. So, simulation based reinforcement learning can provide how to decide its action in various cases. The robots which are trained from off-line simulation in advance can determine its action much faster.

In fact, many researches about applying reinforcement learning to obstacle avoidance problem have already been studied. Q-learning algorithm can be used to the decision of when one agent chooses to rotate or move backward [3]. It works well if they can detect the obstacle without priori knowledge of obstacle position and planning. In addition, navigation problems also have been studied recently [4], [5], [6]. They apply some learning algorithms such as Q-learning, SARSA to the task which reaches the target. To navigate toward the target, the agent has to know the information about the distance from the target position. They get the positive reward when they get closer to the target position. They need more computation time and large iterations. The reward sum is inclined to increase, but, it has unstable fluctuations in common. In the research of Google DeepMind that is one

of the best reinforcement learning team, they applied many revised algorithms like A3C with LSTM to single agent's navigation in complex environments [7]. They also made a single agent navigate in a human-level control using SLAM and auxiliary depth predictions.

However, there are few attempts for multi-robots. If we consider each agent in view of centralized manner, the size of action space becomes larger. So, it is not good for convergence of approximate function of state-value function(Q-function) because there are a lot of state-action pairs and computation complexity is too large. Issue of exploration is also significant for solving problem. Moreover, cooperation imposes constraints on robot's actions. Therefore, some decentralized controls with reinforcement learning are also studied to avoid the difficulties mentioned above. However, I try to apply DQN to multi-robot cooperation in centralized manner. Instead, I suggest some technical methods in order to resolve above issues. The main contribution of this work is an application of reinforcement learning to controlling of multi-robots cooperation, not only a single navigation. Also, it can decrease the computation time compared to the traditional real-time path planning or decision-making. For this achievement, this work includes curriculum learning, revised sampling method specialized for navigation problem, and action space reduction.

The rest of this paper is organized as follows. In Sec. II, we formulated this problems with how to manage objectives in mathematical description and MDP(Markov Decision Process). Which algorithm and approach I apply to this problem and how I revise them with my ideas are established in Sec.III. Simulation set-up such as environment condition and episode is presented in Sec.IV. All simulation results are presented in Sec.V, and Sec.VI concludes the paper.

## II. PROBLEM FORMULATION

### A. Sequential decision-making to reach the target

Many multi-robots problem can be described as a Dec-POMDP(Partially-Observable Markov Decision Process). But, I assume this problem is formulated as a centralized fully observable Markov decision process. The main task is to carry the objects to the target position without dropping and obstacle collision. Then, my objective is to find the decision-making strategy to perform the task well. It is assumed that the object is large or heavy. So, they need to cooperate for carrying the object. I assume two mobile manipulators collaboratively carry a rigid object. Each robot has a manipulator and omni-directional wheels. In order to focus on the navigation problem, I consider their grasping controls and grasping planning could be guaranteed and it is not included in the decision-making process. It only decides how to move for each robots in case of cooperation. Each robot can move toward four directions (UP, DOWN, LEFT, RIGHT) and stop in 2-D area and change the directions arbitrarily. In this work, the center point of them has to reach the target position. When this is accomplished, I regard it as a task success. To carry the object cooperatively, the required condition is the distance between two robots. Because I assume required

manipulator's motion is guaranteed, the cooperation constraint is only that the distance is lower than the given value. In addition, they also need to avoid obstacles while they are going toward the target. This is also described as keeping the distance over certain value. Therefore, this sequential decision-making problem is formulated as navigation problem with the cooperation constraint and the collision avoidance constraint.

### B. Markov Decision Process Description

This sequential decision-making problem can be formulated as a MDP(Markov Decision Process). So, this problem has a tuple of  $\langle S, A, P, r, \gamma \rangle$ .  $S$  is the state space,  $A$  is the action space,  $P$  is the state-transition model,  $r$  is the reward(stage-wise cost), and  $\gamma$  is a discount factor.

- *State Space* :

State is constructed as the robot1's position, robot2's position, target position, obstacle's positions. So,  $s = [x_1, y_1, x_2, y_2, p_{target}, p_{obs1}, p_{obs2}] \in R^{10}$  can be described.

- *Action Space* :

Action space is expressed with  $U = \prod_i^n U_i = U_1 \times U_2$  and  $U_i = [\text{UP}, \text{DOWN}, \text{LEFT}, \text{RIGHT}, \text{STOP}]$  where UP = [0,1], DOWN = [0,-1], LEFT = [-1,0], RIGHT = [1,0], STOP = [0,0]. One step means that every robot decides each actions. It can express every situation because it includes stop motion in the individual action space.

- *State-transition model* :

In this work, only position of each robot changes for simplicity. I do not consider the dynamic obstacles and random targets. State-transition model is simple in this problem. The position is determined by velocity, time gap, noise and moving directions decided by actions. It can be expressed mathematically as below.

$$s_{t+1} = s_t + action \times Vel \times \Delta t + noise \quad (1)$$

- *Reward function* :

A reward function is specified to award the robots for reaching the target and penalize the robots for dropping the objects due to large distance or colliding with obstacles. Detailed value and condition will be described in simulation set-up section.

- *Discount factor* :

It is specified as 0.98 in this work.

In this formulation, each robot has to choose its action every moment in order to maximize the discounted sum of stage-wise cost. But, it is not easy to choose its action because we do not know the exact cost function. So, we apply simulation based learning method to this problem. Two robots run on the simulation environment and get the reward of simulator. The reward is included in the simulation and it is determined by human 's intuition.

## III. APPROACH

I introduce four approaches which represent this work. They are DQN, Revised sampling with success state, Curriculum learning, Reduced action space.

**Algorithm 1: deep Q-learning with experience replay.**  
Initialize replay memory  $D$  to capacity  $N$   
Initialize action-value function  $Q$  with random weights  $\theta$   
Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$   
**For** episode = 1,  $M$  **do**  
  Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$   
  **For**  $t = 1, T$  **do**  
    With probability  $\varepsilon$  select a random action  $a_t$   
    otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$   
    Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$   
    Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$   
    Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$   
    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$   
    Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$   
    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$   
    Every  $C$  steps reset  $\hat{Q} = Q$   
  **End For**  
**End For**

Fig. 2. DQN 2015 algorithm

### A. Deep Q-network

Google DeepMind team showed great achievements in the reinforcement learning in Atari2600 games with DQN(Deep Q-network) 2013 and DQN 2015. They applied Deep learning technique into Q-learning algorithm. Moreover, they introduced replay memory to remove the samples'correlations and the separation of target network from main network. DQN has a great power to approximate Q-value accurately and to find optimal policy with greedy policy. It is based on fundamental Q-learning, but it approximates Q function through DNN(Deep Neural Network) with replay sample and Double Q-network. It can tackle the huge state space and discrete actions through neural network. This problem also has a huge state space which is 10 dimensional vector and discrete actions consisted of 5 actions. DQN can approximate the Q-value of each state and yield the optimal policy in case of its state. In particular, I use DQN2015 learning method in this work. DQN2015 algorithm can be showed in Fig. 2 [8].

### B. Revised sampling with adding success state

One of DQN 2015's gist is random sampling from experience replay memory. This memory stores experience tuple as (state, action, reward, next state, whether episode terminates) every time-step with a limited storage size. But, only minibatches which is randomly sampled from replay memory are used to update the Q-network due to removal of correlation. However, replay memory has not uniform distribution of state-action-reward tuples in this problem. Navigation problems have fewer positive rewards than negative rewards or non-reward in a whole episode. This is because a crucial reward is to reach a target and it occurs only once or never in an episode. So, replay memory has a lot of tuples with obstacle collision or dropping the object(cooperation failure), whereas there are few tuples with reaching target in a replay memory. Even if random sampling method gets rid of correlation of states, randomly

sampling in replay memory mostly does not include success states like reaching a target. It can make the network update weights in an appropriate way. It could be different from our real Q-value and it yields divergence. Therefore, we change some sampling method a little. If we train the neural network with a mini-batch which samples from the replay memory, we add one success state to the mini-batch necessarily. Then, network would update with considerations for success state. It can update the network more properly.

### C. Curriculum learning

Navigation problem is a chronic problem of reinforcement learning because it has a difficulty of sufficient exploration. Other MDP problems usually have some explicit rewards frequently during training and they also just want to increase their total reward(cost function) ultimately. For example, Atari game steadily gets the score reward and obstacle avoidance also get penalty reward often. But, the navigation problem gets only one strong reward after they reach the target and they do not need to increase the reward after they achieve it. The episode is terminated if the agent reaches the target. So, the issue of exploration is more important than other problems. This tendency also increases when the environment is huge and state space is large. If they cannot explore the space enough, the problem's solution can go to the local optima or very slow convergence. When they do not explore near the target, they have little experience about high rewards and Q-network updates with improper weights. Also, the convergence of Q-value toward optimal Q-value is not guaranteed in a fundamental assumption of Q-learning because every state-action pairs are not visited infinitely.

To resolve the issue of exploration and slow convergence of neural network, various method are suggested. Some researchers developed curriculum learning, where the goal is to design a sequence of source tasks for an agent to train on, such that final performance or learning speed is improved [9]. It is similar to providing the agents with the guidance of human-education. The researcher provides some easy missions for agents in advance because they want to update the network with great initialized weights and the agents can explore some suitable state earlier. After they complete some easy missions, I can provide more complicated missions. Like human's curriculum, the researcher increases the level of difficulty gradually.

In this experiment, I first offer a state with a close distance toward the target. I save the weights of the network. After that, a state with the further distance is offered and started to train with previous weights. Then, I start the original mission with pre-trained Q-network. It will have a great effect to stabilize convergence and reduce hours of computation.

### D. Reduced Action space

When multi-robots perform on some tasks, the size of action space inevitably increases rapidly. It is proportional to the product of each robot's action space sizes. So, it increases exponentially. If we use n-robots in this experiment, I have to

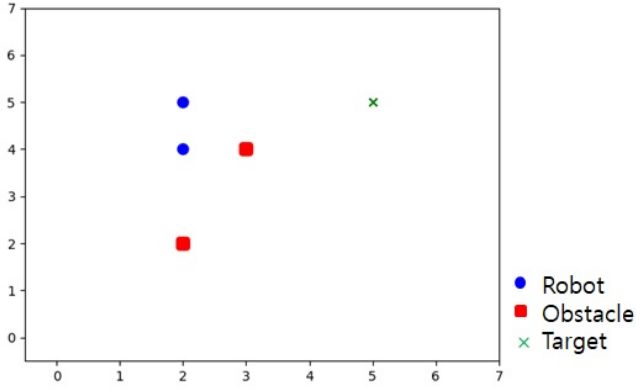


Fig. 3. 2-Dimensional map environment

consider the action space size of  $5^n$ . It can be so huge that I have to spend a lot of time to converges the network because it is same as the output layer size. I need to use large size of layers and need to compute more times. So, I suggest new method which can create Q-value in a different way. It is the method which reduces the size of action space. Previous action space  $U = \prod_i^n U_i$ , however, we select  $U = \sum_i^n U_i$  instead. This reduction of action space size changes the meaning of each time step. In previous action space, each time step means the decision of every robots'action. But, in this action space, each time step means the decision of one action among every robots'actions. Even if steps need to proceed twice(n-times), this will be much faster when the neural network updates and predicts.

#### IV. SIMULATION SETUP

##### A. Environment

This work is modeled in a 2-D planar environment which is a  $7 \times 7$  square. Two robots can move inside the map and both position x and position y are included in  $[0,7]$ . This work assumes the robot 's sensor detecting and localization are sufficient to know the position of each robot and target directly. So, this work can include the position value in the state. Target position and two obstacles positions are fixed in this work, but, we can change them if we have the long training-time. Robot1's position and robot2's position can change while the simulation runs. Fig.3 shows this simulation environment. This environment is made by myself using python matplotlib.

##### B. State-action space discretization

Since DQN method cannot be applied to continuous action values, I have to discretize the action space. So, I divide it by direction and add stop motion to it. Also, I should consider state space because it can have infinite states unless we discretize the state space. The condition that approximation of q function would be accurate is the small number of states. The more states there are, the harder the network approximates the Q-function. Also, Q-learning is unlikely to converge because states cannot be visited sufficiently. Moreover, even if I

consider only 4 elements among 10 elements, the complexity of computation can be so high when I divide the space finely. So, this work determines that the state space is divided by 0.5 intervals. It means that this work decides the velocity as fixed value(0.5) and time gap is 1. At each state, robot1 and robot2 have 5 actions each other. So, total action space has 25 different action selections.

##### C. Reward function

Reward is provided every time the robots do some action in the simulation. The robots try to estimate the real cost function from this reward. The reward function has been chosen by human intuition. It should be able to award the agent for approaching the target and penalize the agent for dropping the object or colliding with obstacles. The reward function is set up as follows in this work.

- $Reward = 1000$   
(when the center of robots reaches the target )
- $Reward = -2 \times distance$   
(when the distance between two robots is higher than criteria)
- $Reward = -1$   
(when they collide with the obstacles)
- $Reward = 100 \times (\frac{1}{d_{t+1}} - \frac{1}{d_t})$   
(from potential field difference)

##### D. Episode and step

The simulation runs on python with tensorflow library. I proceed this simulation based on curriculum learning. So, I divide training process into two steps as below. Easy task has a starting point of  $[2,5,3,6]$  and hard task has a starting point of  $[2,4,2,5]$ .

- 1) Easy task : 500 episodes (max 2000 steps per episode)
- 2) Hard task : 3000 episodes(max 10000 steps per episode)

One episode iterates within maximum step number 2000 or 10000 steps. Reaching target terminates the episode. Mini-batch size of 50 and one epoch corresponds to weight updates as same number of episode. Also, I use annealing-e in an e-greedy policy, which means e decreases from 1 to 0.1 before episode gets to  $0.1 \times MAX_{EPISODE}$  0.1 after it gets to  $0.1 \times MAX_{EPISODE}$ .

##### E. Neural Network structure

This work used MLP(Multi-Layer PerceptronS) instead of CNN because the computation capacity cannot cover the Image input. This work uses a neural network with two hidden layers of width (16,32) and activation layer is ReLU function. We have two types of how to construct neural network about Q-value. First thing is to put both states and actions into the input layers. The second thing is to put only states into the input layers and output layer has the same number of values as the number of actions. This work chooses the latter one.

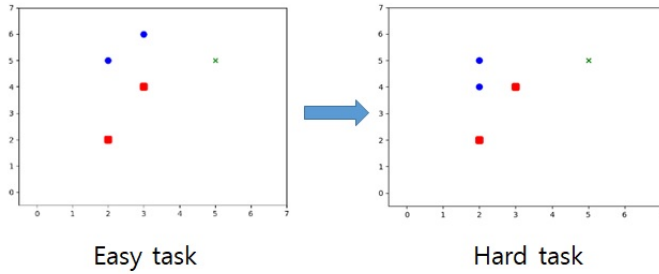


Fig. 4. Two types of task(Left : easy, Right : hard)

The network updates its weights to minimize the difference between the predicted output and the Q-value. The minimizer is Adam-optimizer and the period of copying target network is 100 steps for easy task and 1000 steps for hard task.

## V. EXPERIMENT RESULTS

This work performs on a LG-gram laptop with an i5-6200U CPU, using a python and tensorflow library on jupyter notebook. It takes a lot of time to perform this task because of this low performance capacity. This restricts the diverse experiments. This work usually compares the change of the discounted reward sum or average of maximum Q-value.

### A. Penalty reward vs constrained action

This problem is constrained maximization problem, which has the cooperation constraint, obstacle avoidance constraint, map size constraint. This work already decides to manage cooperation and obstacle avoidance constraint through negative reward. However, I try to tackle the map size constraint through two different methods. The first method is to penalize the agent with a negative reward(-50) if its action violates the constraint which makes robots outside the map. The second method is to enforce the agent not to be able to violate the constraint. The latter is decided inside the environment setting itself. As a result, the former does not show any convergence while the latter shows faster convergence.

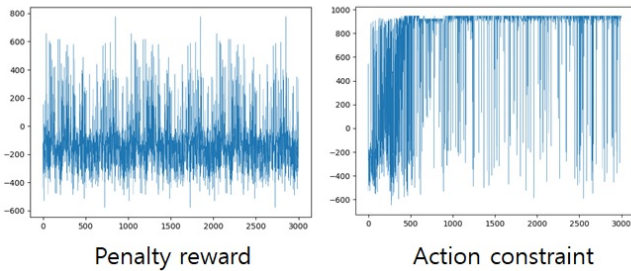


Fig. 5. Discounted reward sum versus episode (Left : penalty reward, Right : constrained action)

### B. Random sampling vs Revised sampling

Replay memory has to include diverse tuples, so we add success case into mini-batch by force to cover its rarity. This work runs with no curriculum learning and 3000 episode with maximum 10000 steps. This work compares discounted reward sum and average of maximum Q-value.

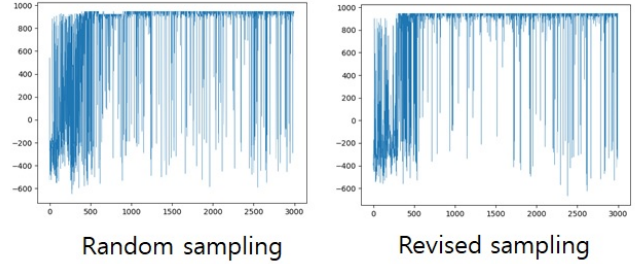


Fig. 6. Discounted reward sum versus episode (Left : random sampling, Right : revised sampling)

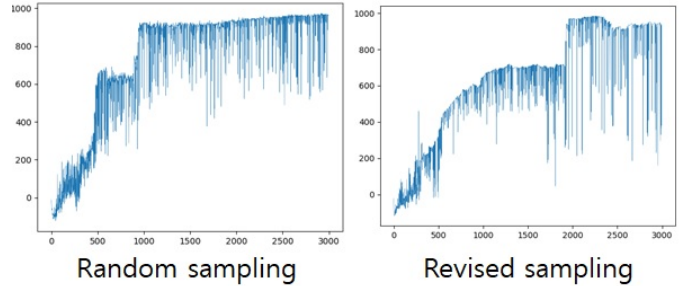


Fig. 7. Average of maximum Q value sum versus episode (Left : random sampling, Right : revised sampling)

In this work, two sampling method cannot show remarkable difference between two things because it converges fast enough. However, I can find revised sampling is more stable as you can see in Fig.6. Average of maximum Q-value looks a little different. Random sampling looks more stable than revised sampling. It needs other examples or experiments for longer episodes. However, it shows that random sampling's early stage increases rapidly with lucky inclusion of many success cases whereas revised sampling's early stage steadily increases. It is supposed that low exploration stage of early episode can be stable from revised sampling.

### C. Curriculum learning vs Normal learning

This work try to compare curriculum learning with normal learning. This work struggled to show the obvious difference between two learnings and the great effect of curriculum offer. So, very hard task was tried in 6000 episode with maximum 10000 steps. But, it took a lot of time and it cannot be managed by my laptop. Therefore, I changed it to easier task. Nevertheless, it can be compared well. It shows enough effect of curriculum learning.

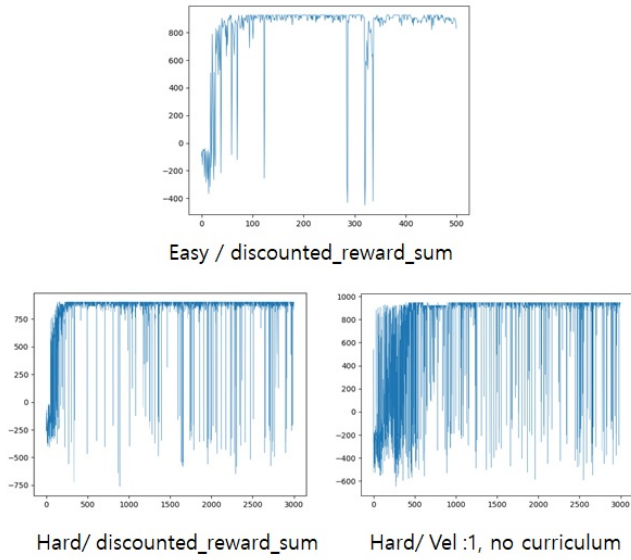


Fig. 8. Average of maximum Q value sum versus episode (Top: discounted reward sum of easy task, Left : discounted reward sum of hard task, Right : discounted reward sum of hard task without curriculum)

Fig.8 shows that curriculum learning converges much faster and becomes more stable than normal learning. Furthermore, normal learning even raises velocity because the case of velocity 0.5 did not converge until 2000 episode with 10000 steps. In spite of this penalty, curriculum learning achieves greater performance.

#### D. Action space reduction

As mentioned before, this work reduces actions space size  $5^2$  to  $5+5$ . It can help the computation time decrease and it can help me tackle this problem in centralized manner within available operations.

As expected, its running time is shorter than previous method 's thing. It takes 556s for easy task and 3297s for hard task while previous method takes 820s for easy task and 5897s for hard task. It is showed that reduced action space learning requires more episodes to reach the optimal reward sum. However, if we consider the meaning of one step, it is not slow convergence. But, its stability becomes much worse until enough exploration as you can see in Fig.9. The hopeful point is that fluctuation decreases and becomes quite stable after enough exploration proceeds(after 400 episodes for easy task and 2500 episodes for hard task).

## VI. CONCLUSION

### A. Analysis

After training from DQN, two robots can navigate toward the target well as you can see in Fig.10 and Fig.11. Before training, two robots move meaninglessly far away from the target. However, two robots become to be able to go toward the target and keep distance under the criteria with obstacle

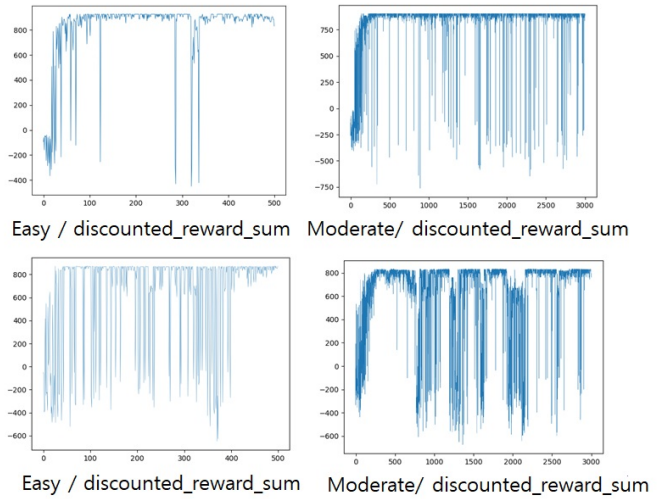


Fig. 9. Comparison of total action space and reduce action space about discounted reward sum (Top: total action space with curriculum learning, Bottom : reduced action space with curriculum learning)

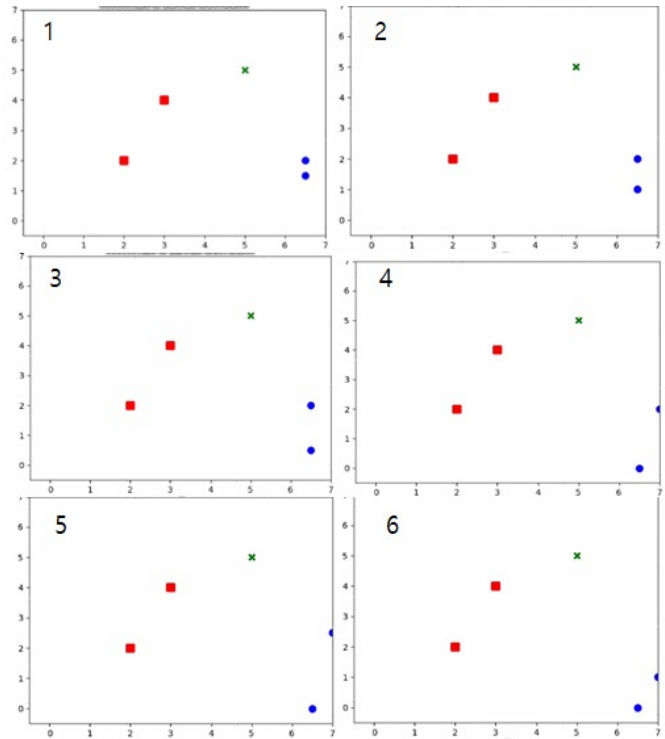


Fig. 10. Simple snapshots of robot movement before training

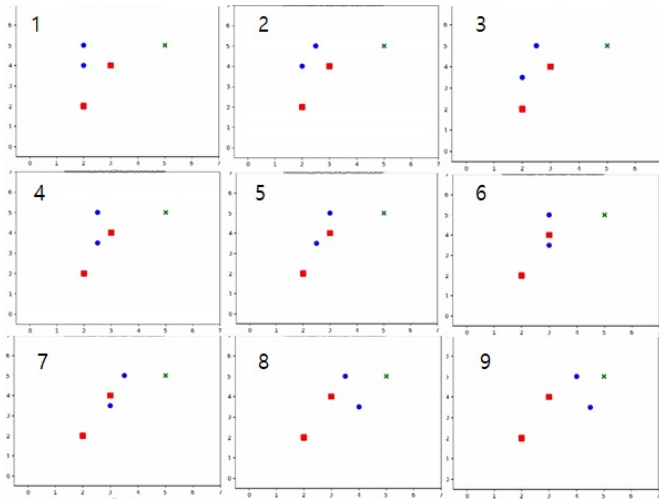


Fig. 11. Simple snapshots of robot movement after training

avoidance. I can find the achievements of task from the DQN training.

This work contributes to the applications of DQN into multi-robot cooperation. The random sampling from replay memory is not always good for navigation problem, so revised replay memory sampling is suggested and it stabilizes when sufficient explorations are not satisfied. Also, it shows that curriculum learning is very powerful in fast convergence and stability of navigation problem. New attempt to consider the curse of dimensionality is the reduction of action space, which lowers computation time, but increases the instability and fluctuation.

The obstacle avoidance problem is well solved by Q-learning as mentioned before. Moreover, dynamic obstacle problem has also been solved. I guess cooperation problem is similar to this dynamic obstacles problem because it is about changing distance constrained problem. However, both action selection influences on the distance and it makes difference from dynamic obstacles. It is hard to choose best action for obstacle avoidance if I also consider reaching the target.

Navigation problem has a characteristic that its rewards are often sparsely distributed in the environment and only one strong reward is located densely. It can yield some fluctuations in overall even if it converges a fixed value. It can be showed in all graphs.

### B. Future work

New curriculum learning will be studied. Two suggestions exists. First thing is that multi-robot's cooperation are trained after path planning is created through learning of a single agent earlier. Second thing is that human provides the guidance of traditional navigation solution, which pre-calculated optimal solution such as LQR. Decentralized decision-making can be also studied. Only one agent knows all state information, while others just know the distance from each other. More complicated manipulators kinematics can be included in a state model. So, I will consider the end-to-end learning of carrying

the object task. Because DQN is not stable for policy and it has clear limitations of computation, other learning method like A3C has to be tried and compared to DQN. It is necessary to elaborate the environment for fancy performance. This work will be practical if it considers more actions like 8 directions. Moreover, velocity or obstacles position can change.

### ACKNOWLEDGMENT

This work is researched in graduate class named "Stochastic Control and Reinforcement Learning".

### REFERENCES

- [1] J. Alonso-Mora, S. Baker, and D. Rus, "Multi-robot navigation in formation via sequential convex programming," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 4634–4641.
- [2] —, "Multi-robot formation control and object transport in dynamic environments via constrained optimization," *The International Journal of Robotics Research*, vol. 36, no. 9, pp. 1000–1021, 2017.
- [3] B.-Q. Huang, G.-Y. Cao, and M. Guo, "Reinforcement learning neural network to the problem of autonomous mobile robot obstacle avoidance," in *Machine Learning and Cybernetics, 2005. Proceedings of 2005 International Conference on*, vol. 1. IEEE, 2005, pp. 85–89.
- [4] J. Muhammad and I. Ö. Bucak, "An improved q-learning algorithm for an autonomous mobile robot navigation problem," in *Technological Advances in Electrical, Electronics and Computer Engineering (TAECE), 2013 International Conference on*. IEEE, 2013, pp. 239–243.
- [5] L. Khriji, F. Touati, K. Benhmed, and A. Al-Yahmedi, "Mobile robot navigation based on q-learning technique," *International Journal of Advanced Robotic Systems*, vol. 8, no. 1, p. 4, 2011.
- [6] N. ALTUNTAŞ, E. Imal, N. Emanet, and C. N. Öztürk, "Reinforcement learning-based mobile robot navigation," *Turkish Journal of Electrical Engineering & Computer Sciences*, vol. 24, no. 3, pp. 1747–1767, 2016.
- [7] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu, *et al.*, "Learning to navigate in complex environments," *arXiv preprint arXiv:1611.03673*, 2016.
- [8] 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, p. 529, 2015.
- [9] A. Graves, M. G. Bellemare, J. Menick, R. Munos, and K. Kavukcuoglu, "Automated curriculum learning for neural networks," *arXiv preprint arXiv:1704.03003*, 2017.