

# Non-Reciprocating Sharing Methods in Cooperative Q-Learning Environments

Bryan Cunningham

*Department of Computer Science  
Virginia Polytechnic Institute and State University  
Blacksburg, USA  
bcunn06@vt.edu*

Yong Cao

*Department of Computer Science  
Virginia Polytechnic Institute and State University  
Blacksburg, USA  
yongcao@vt.edu*

**Abstract**—Past research on multi-agent simulation with cooperative reinforcement learning (RL) focuses on developing sharing strategies that are adopted and used by all agents in the environment. In this paper, we target situations where this assumption of a single sharing strategy that is employed by all agents is not valid. We seek to address how agents with no predetermined sharing partners can exploit groups of cooperatively learning agents to improve learning performance when compared to Independent learning. Specifically, we propose 3 intra-agent methods that do not assume a reciprocating sharing relationship and leverage the pre-existing agent interface associated with Q-Learning to expedite learning.

**Keywords**—Multi-Agent Reinforcement Learning; Cooperative Learning; Agent Interaction Protocols; Information Exchanges in Multi-Agent Systems

## I. INTRODUCTION

Traditional reinforcement learning (RL) simulations imbue an agent with no knowledge of the environment in which they are located a priori. By exploring this environment and gaining direct experience with it, the agent will update its policy, or memory, to remember the best action(s) to take in each state. The goal of the agent is to improve its policy to maximize its performance in the environment [1]. Multi-agent-based cooperative RL simulations allow for multiple agents to coexist in the same environment and improve learning performance by exploiting the policy information of the other agents [2]. Enabling multiple agents to exist in the same environment adds complexity to the learning problem. Agents need to balance the tasks of exploring the environment and exploiting their own policy in addition to exploring the usefulness of the other agents' policies and exploiting them. Ideally, this improves the learning performance of the agents when compared to training them separately in single-agent RL.

Past research on multi-agent simulation with cooperative reinforcement learning (RL) focuses on developing sharing strategies that are adopted and used by all agents in the environment. We consider these sharing strategies to be reciprocating because all participating agents have a pre-defined agreement regarding what type of information is shared, when it is shared, and how the participating agent's

policies are subsequently updated. The sharing strategies are specifically designed around manipulating this shared information to improve learning performance. In this paper, we target situations where this assumption of a single sharing strategy that is employed by all agents is not valid. This research generalizes the cooperative learning domain by considering situations in which agents may use differing sharing strategies to cooperatively learn a task within the same environment. We seek to address how agents with no predetermined sharing partners can exploit groups of cooperatively learning agents to improve learning performance when compared to independent learning. Specifically, we propose 3 intra-agent methods that do not assume a reciprocating sharing relationship and leverage the pre-existing agent interface associated with Q-Learning to expedite learning.

This research is directly applicable for physically embodied agents in the multi-agent cooperative robotics field. Consider the following situations: What if a new agent is introduced into the environment that does not know the sharing strategy of the other agents, but it still wants to exploit their knowledge. Alternatively, perhaps the agent cannot use the same sharing strategy due to hardware or computational restrictions. Finally, consider the case where there are 2 groups of agents in the environment that are utilizing their own sharing methods but each want to exploit the knowledge of the other group. The concept of non-reciprocating sharing methods provides a solution to these problems by making agent-group and group-group collaboration possible. To the best of our knowledge, research into non-reciprocating methods has not been addressed within the cooperative Q-learning field for homogeneous agents.

Whereas other cooperative learning research assumes the simulation designer has control of all other agents in the environment for the purposes of encoding an agent interaction strategy, this research will investigate algorithms that do not so assume. The other agents' functions and their sharing strategies are unknown and inaccessible from the point of view of the agent(s) using our proposed methods. We contribute 3 sharing methods: the modified averaging

(MA) method, the modified experience counting (MEC) method, and the hybrid experience counting and averaging (HECA) method. These 3 methods adapt existing sharing strategies to perform in environments where other agents do not reciprocate. We analyze our methods by testing them in 3 unique environments where the other agents employ a variety of standard sharing methods.

Findings indicate that learning performance while using our proposed methods is improved when compared to independent learning. Results show that agents with no predetermined sharing partners can successfully exploit groups of cooperatively learning agents. Additionally, the MEC and HECA methods perform best overall across the presented experiments.

The paper is organized as follows. Section 2 presents the related work and further motivation for this research. Section 3 provides the background information relevant to this research and describes the assumptions behind our work. Section 4 describes the 3 proposed algorithms and section 5 discusses the experiments associated with this research. Section 6 presents the results of our experiments and section 7 provides concluding remarks and presents future work.

## II. RELATED WORK AND MOTIVATION

Previous research on direct communication cooperative learning agents focuses on developing sharing strategies that are adopted by all agents in the environment. Tan [2] introduced cooperative learning methods where agents share sensations, episodic information, or learned policies to achieve improved learning rates. The work presented the concept of policy averaging (PA), in which agents average their learned policies together and exploit the other agents' shared knowledge to expedite learning. Another cooperative learning technique is to have agents make use of a joint policy table [3]. Using this method, several distinct agents explore an environment and update the same policy that is shared by all agents. Agents do not communicate directly with one another but rather to a central entity that stores the shared policy. In our work, we refer to this method as the Centralized sharing method and use it as an upper bound to compare our agents' learning performance against. A similar concept was proposed by Yang et al. with a blackboard-based communication method in which agent communication is relayed through and controlled by this central blackboard architecture [4]. Ribeiro et al. propose a cooperation model that uses a global action policy to ensure proper convergence. The global action policy unifies agents' partial action policies to produce optimal policies for the generic interaction model case [5]. Multi-agent cooperative learning has been shown to improve learning performance for both the homogeneous and heterogeneous cases [6].

The weighted-strategy sharing (WSS) measures the expertness of the cooperating agent's policies and weighs their contributions according to this value. This method has been

shown to be effective when agent experiences differ [7], [8]. The WSS cooperative learning method could be considered a viable non-reciprocating strategy. However, results show that for the case where all agents start in a similar initial location, the agents' learning performance does no better than when using Independent learning. Another technique for determining an agents area of expertise uses several classifiers [9] to extract partial policies containing expertise. Other techniques furthering the identification and exploitation of agents areas of expertise have also been presented [10].

## III. BACKGROUND

In this section we will briefly describe the mechanics necessary to understand the basics of RL in addition to the action-selection (AS) method used in this paper.

### A. Reinforcement Learning

RL is a bottom-up programming methodology that gives agents the ability to extract salient environmental cues online. RL uses Markov decision processes (MDP) to model how an agent interacts with its environment. An MDP is a 4-tuple taking the form  $(S, A, P_{ss'}^a, R_{ss'}^a)$  where  $S$  is the state space,  $A$  is the action set,  $P$  is the transition function where  $P_{ss'}^a$  represents the probability of transitioning from state  $s$  to state  $s'$  via action  $a$ , and  $R$  is the reward function where  $R_{ss'}^a$  represents the expected value of the reward achieved when an agent moves from state  $s$  to state  $s'$  via action  $a$ . As an agent explores its environment, it updates its policy function  $\pi$  that maps each state  $s \in S$  and action  $a \in A(s)$  to  $\pi(s, a)$ . The agent attempts to maximize the expected total sum of rewards gained over time to converge to an optimal policy  $\pi^*$  (multiple may exist).

For this paper, we use a popular form of RL called Q-Learning: a type of temporal-difference (TD) learning where  $Q : S \times A \rightarrow \mathbb{R}$ . TD learning is a learning technique in which an agent will update its previously estimated state values using the differences between its current and former values. This effectively propagates more accurate estimates of the state values as learning continues. Agents define an action-value function for policy  $\pi$  by  $Q^\pi(s, a)$  where the agent takes action  $a$  in state  $s$  and then applies policy  $\pi$ . Q-Learning represents an off-policy form of TD control which, for the one-step case used in this paper, takes the following form:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t(s_t, a_t) \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]. \quad (1)$$

where  $t$  is the time step parameter,  $\alpha$  is the learning rate parameter, and  $\gamma$  is the discount rate parameter. In this paper, the single-agent form of RL is applied to each agent in a multi-agent environment.

## B. Agent Action-Selection

As agents navigate the environment, various strategies for selecting an action in each state can be used, some of which take into account the exploration-exploitation trade-off. Typical strategies include random selection, optimal selection,  $\epsilon$ -greedy selection,  $\epsilon$ -decreasing selection, and softmax selection [1]. In this paper we use softmax AS. Agents use the softmax AS rule with the Boltzmann distribution, as seen in (2), to determine the probability of selecting action  $a_i$  in state  $s$ .

$$P(a_i|s) = \frac{\exp(Q(s, a_i)/\tau)}{\sum_{a \in A(s)} \exp(Q(s, a_k)/\tau)}. \quad (2)$$

where  $\tau$  represents the temperature parameter. As  $\tau \rightarrow \infty$  the probability of selecting actions becomes more random. As  $\tau \rightarrow 0$  the probability of selecting the actions with higher Q-values becomes more likely.

## C. Key Assumptions

We make the following 3 assumptions in this work that persist throughout the paper:

- 1) Agents are assumed to be using Q-Learning. The agents have a public interface that allows other agents to query their Q-tables. This is a minimal method for catalyzing cooperative learning.
- 2) All agents have the same Q-Learning parameter values.
- 3) The agent(s) employing our method knows how many other agents are in the environment and has enough memory to store the Q-tables and E-tables for them.

## IV. THE ALGORITHMS

In this section we describe the architecture of the agent in addition to the 3 algorithms we have developed.

### A. General Agent Structure and Behavior

Traditional agents in a RL simulation have a single policy representing their memory. Our approach gives the agent the memory capacity to embed the policies of the other agents. Note that for this paper only Agent 0, or  $A_0$ , utilizes the agent structure described in this section. The other agents in the environment will retain the standard agent structure with 1 policy,  $\pi_0$ , for memory to store a Q-table and potentially an E-table for counting experience, depending on the sharing method used.

*Definition 1:* An agent's memory,  $M$ , is a set of  $N$  policies  $\{\pi_0, \dots, \pi_{N-1}\}$  where  $\pi_0$  is the agent's own, or self, policy and policies  $\{\pi_1, \dots, \pi_{N-1}\}$  are shared policies. Each shared policy  $\pi_i \in M$  is a policy from each of the other  $(N - 1)$  agents.

Our approach separates the two types of policies in memory and updates both accordingly. The agent updates its self policy when an action is directly taken by the agent and therefore experiences being in state  $s_t$ , taking action  $a_t$ , and arriving at state  $s_{t+1}$  with reward  $r_{t+1}$ . The agent

updates its shared policies after each time step by querying the other agents for their Q-tables.

There are 5 steps taken by the agent at each time step: 1) Action Selection: Take a step with action  $a$  in state  $s$  using policy  $\pi_0$  to make a decision, 2) Update  $\pi_0$  based on the experience gained by taking action  $a$  in state  $s$ , 3) Query Q-values from the other agents' policies and update policies  $\{\pi_1, \dots, \pi_{N-1}\}$  accordingly, 4) Assess/evaluate the best values to use, and 5) Combine the values to create an updated policy,  $\pi_0$  (that is then accessible to others). This interaction cycle is shown in Fig. 1. Our proposed methods operate at the level of steps 4 and 5 and are differentiated by the actions taken in these steps. Note that the 3 methods presented below still perform learning by trial-and-error, as in step 2. These methods simply supplement the agents' knowledge by exploiting the other agents' knowledge.

### B. Algorithm 1: Modified Averaging (MA) Method

The PA method first suggested by [2] is a simple sharing strategy that averages all agent policies together – according to each  $(s, a)$  pair – and assigns the averaged policy to all agents. After all agents have taken their step for time step  $t$  in the episode, the PA method averages the Q-value of each action  $a \in A(s) \forall s \in S$  with all other corresponding actions in their corresponding states across all agent policies. The method then overwrites the Q-values across all policies for each action  $a \in A(s) \forall s \in S$  with the corresponding average for that action. As all agents participate and use the same averaged policy, this is considered a reciprocating sharing strategy.

We propose the MA method, a non-reciprocating version of the PA method. The MA method, as seen in (3), is similar to the PA method except that after averaging all agents' policies, only  $\pi_0$  is assigned the averaged Q-table because the other agents are not participating in the sharing method.

$$\text{MA}(s, a) = \frac{1}{n} \sum_{i=0}^{n-1} Q^{\pi_i}(s, a) \quad (3)$$

for  $s \in S, a \in A(s), \pi_i \in M$

### C. Algorithm 2: Modified Experience Counting (MEC) Method

The experience counting (EC) method, also referred to as the non-trivial update counting method in [11], is a popular sharing strategy that has been shown to outperform the PA method. The EC method makes use of a E-table, or experience table, in addition to a Q-table to keep track of which states and what actions an agent has actually experienced during the simulation. The central idea behind this method is that the most experienced agent with regard to each  $(s, a)$  pair should have the most contribution to the Q-table that is synchronized by all agents after each step.

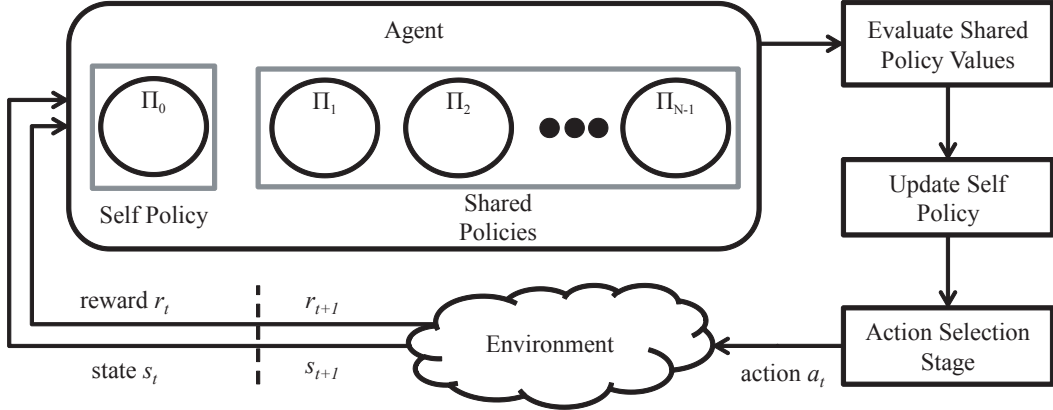


Figure 1. Extended RL agent-environment interaction cycle.

Each  $(s,a)$  visitation results in the corresponding  $(s,a)$  entry of the E-table being updated by increasing the visitation, or experience, count for that entry by 1. After all agents have taken their step for time step  $t$  in the episode, the EC method goes through each  $(s,a)$  pair in the Q-table and allows the agent with the most experience for that  $(s,a)$  pair to contribute its current Q-value for  $(s,a)$  to the resultant Q-table that will overwrite all agent policies. Experience ties are resolved by randomly selecting from the tied experience values, which correspond to Q-values. Similarly, each  $(s,a)$  entry in the E-table for all agents is updated to the highest count experienced out of all agents because the corresponding entry in the Q-table was updated. At the end of each step all agents have the same Q-table and E-table.

We propose the MEC method, as seen in (4), a non-reciprocating version of the EC method. Because an agent employing the MEC method cannot assume the other agents are reciprocating, it cannot rely on receiving others' E-tables for counting experience. Instead, we propose that the agent queries and stores the other agents' Q-tables in  $M$ . Additionally, the agent will create a E-table for each of the other agents in  $M$  as well. The agent can then observe which Q-values have changed since the last time step  $t$  and will associate this with an experience point for that particular  $(s,a)$  entry in the E-table for that agent. Therefore, the essence of the EC method is captured and the agent can reconstruct the  $(s,a)$  pair(s) that the other agents have experienced to determine which agents should contribute to the policy that will be assigned to  $\pi_0$ . The algorithm assigns an experience point when a Q-value for  $(s,a)$  at time step  $t$  does not equal the Q-value for  $(s,a)$  at time step  $t + 1$ .

$$\begin{aligned} \text{MEC}(s,a) &= Q^{\pi_j}(s,a) & (4) \\ \text{where } j &= \arg \max_i E^{\pi_i}(s,a) \\ \text{for } s &\in S, a \in A(s), \pi_i \in M \end{aligned}$$

#### D. Algorithm 3: Hybrid Experience Counting and Averaging (HECA) Method

The HECA method is a hybrid method combining elements from both the MA and MEC methods. The algorithm assigns an experience point when a Q-value for  $(s,a)$  at time step  $t$  does not equal the Q-value for  $(s,a)$  at time step  $t + 1$ . This method is similar to the MEC method except that experience, or E-table, ties are resolved by averaging the Q-values associated with the tied experience values. All Q-values for  $(s,a)$  at time step  $t$  that have equal experience values should be given equal representation in determining the Q-value for the resultant Q-table. Instead of randomly selecting one of the Q-values associated with one of the tied E-table values, this method averages the Q-values associated with the tied E-table values for  $(s,a)$ .

## V. EXPERIMENTS

### A. Maze World Domain

To explore the methods we propose, an  $8 \times 8$  grid-based environment is used to simulate the maze world domain [12]. The homogeneous agents attempt to navigate to the goal location with the intention of reaching it using the fewest number of steps. The set of states  $S$  is composed of each cell position  $(x,y)$  on the map.  $|S| = c \times d$  where  $c$  is the number of cells in the horizontal direction and  $d$  is the number of cells in the vertical direction. In each state, agents can move one step and choose from the following action set:  $A = \{\text{North, East, South, West}\}$ . Agents that attempt to move into a wall are blocked and will remain in the same position. Similarly, if an agent attempts to leave the bounds of the map, the agent will remain in the same position. Agents may occupy the same cell. Agents start each episode in the cell with position  $(0,3)$ . The goal cell is located in position  $(7,3)$ . For every step an agent takes that leads to a non-goal state the agent receives a reward of -1.5. Conversely, actions that lead the agent into a goal state provide the agent with a reward of 20. The maze world map used in testing was

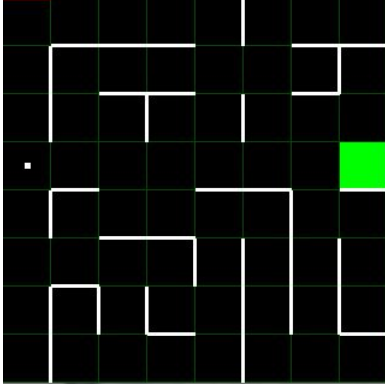


Figure 2. A randomly generated  $8 \times 8$  maze world domain map. The white lines are walls, the white dot is an agent, and the green cell is the goal.

randomly generated with respect to wall placement. See Fig. 2 for an example of the maze world environment.

### B. Evaluation

An episode consists of the steps taken after an agent starts in the initial state and spans until the agent reaches the goal state. A run is defined by a series of episodes. The number of episodes per run is determined by the number of episodes that are necessary to allow the learning performance of the agents to stabilize. Agent policies are carried over from one episode to the next, but not from run to run. In our experiments, we evaluate the learning performance of the agents over 500 runs, where each run consists of 60 episodes. For the simulation parameters we set  $\alpha = 0.15$ ,  $\gamma = 0.99$ , and  $\tau = 0.4$ ; these are default parameters for the maze world domain [11].

For each of the 3 experiments in this paper, we present a graph comparing the average number of steps per episode for  $A_0$  vs. the number of cooperating agents in the simulation for the proposed sharing methods. The average number of steps per episode signifies learning performance. These graphs provide visual results displaying the trends associated with each of the sharing methods. We vary the number of agents up to 8 because, at this point, we near performance convergence for the sharing methods. For each experiment, we include a table displaying the average number of steps per episode for  $A_0$  for the case when the simulation has 8 agents. As aforementioned, the sharing method performance converges when the agent count is 8 and therefore these step averages are good overall indicators of the performance for the sharing methods. We also include the averages' respective 95% confidence intervals (CI) calculated by a t-test. A CI allows us to determine whether or not the differences in learning performance, as indicated by average steps per episode, are statistically significant. Additionally, each experiment is also accompanied by a graph displaying the average number of steps per episode vs. the number of

episodes. These step averages are calculated across all 500 runs for  $A_0$  when there are 7 other agents in the environment simultaneously learning.

To simulate testing environments with varying sharing strategies and to test the effectiveness of our proposed methods, we have set up 3 unique environments that represent a range of the types of sharing strategies that may be encountered.  $A_0$  will employ one of the 3 proposed sharing methods while the other agents,  $[A_1, \dots, A_{N-1}]$  will employ one of the following methods: 1) Independent: No sharing strategy is employed and each agent's Q-table will likely be different from the other agents' Q-tables. At each time step  $t$ , at most 1 Q-value in the Q-table will change. 2) PA: The agents use the PA sharing method. Agents using the PA method will have the same policy  $\pi$  after each time step  $t$  because the method dictates that agents synchronize their Q-tables after averaging the Q-values. Therefore, it is possible for multiple Q-values to change between time steps. 3) EC: The agents use the EC sharing method. Again, agents will have same policy  $\pi$  after each time step  $t$  and it is possible for multiple Q-values to change between time steps. For each of the experiments we compare against 4 reference methods: Centralized, PA, EC, and Independent. The Centralized sharing method represents an upper bound for learning performance as all agents write to the same policy after each step. Conversely, the Independent method serves as a lower bound for performance as any sharing methods that perform worse than this non-sharing method might as well learn independently. Inclusion of the PA and EC reference methods allows us to compare the performances of the reciprocating and non-reciprocating versions of the methods. The 4 reference methods are run in an environment where all agents are using that specific sharing method.

## VI. RESULTS

### A. Experiment 1

In this experiment, we explore how the proposed algorithms perform for  $A_0$  when agents  $[A_1, \dots, A_{N-1}]$  are using the Independent sharing method. Results are shown in Fig. 3. The results indicate that as the number of agents increases, the average number of steps per episode for each method decreases. Intuitively, this makes sense because agents are able to take advantage of more knowledge accumulated by other agents learning in the same environment. Once the sharing method's performance stabilizes, it is clear that the EC and Centralized methods perform best and the Independent method performs worst. Table I shows that, according to t-tests, the differences in performance among all listed methods, excepting the MEC and HECA methods, are statistically significant. Fig. 4 illustrates how all agents ultimately converge to the same policy, but vary according to learning performance rates.

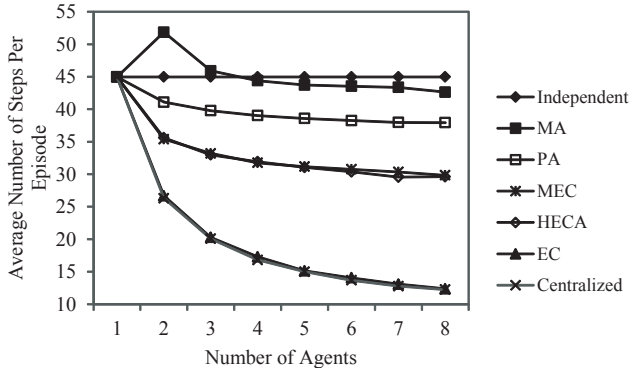


Figure 3. Simulation results for  $A_0$  when the other agents use the Independent sharing method

Table I  
SIMULATION RESULTS FOR  $A_0$  WHEN THE OTHER AGENTS USE THE INDEPENDENT SHARING METHOD

Sharing Method	Average Steps with 95% CI
Independent	$44.98 \pm (0.33)$
MA	$42.64 \pm (0.37)$
PA	$37.94 \pm (0.30)$
MEC	$29.84 \pm (0.21)$
HECA	$29.74 \pm (0.20)$
EC	$12.37 \pm (0.11)$
Centralized	$12.24 \pm (0.11)$

The MA method performs more poorly than the PA method. This is expected as the PA method operates with all agents participating in synchronizing to the same Q-table after each time step  $t$ , thus allowing agents to make an action-selection during their next step using the same Q-table as one other. The MA method operates with  $A_0$  having a Q-table that is not synchronized with the other agent’s Q-tables and learning performance goes down as a result. Q-values can conflict with one another and cause harmful learning effects. Note that there is a negative performance spike for the case when there are 2 agents in the environment, where  $A_0$  is learning with the MA method and  $A_1$  uses Independent learning. The reasoning is as follows: consider the case where  $A_0$  has experienced  $(s,a)$  at time step  $t$  and  $A_1$  has not.  $A_0$  has more beneficial knowledge due to direct experience with  $(s,a)$  and  $A_1$  does not. When determining the resultant Q-value for  $A_0$  using the MA strategy,  $A_1$ ’s Q-value contribution will be just as significant as  $A_0$ ’s Q-value contribution. After each successive time step where no updates occur on the  $(s,a)$  Q-value, the resultant  $(s,a)$  Q-value will continue to be averaged with  $A_1$ ’s outdated Q-value. This guides  $A_0$ ’s Q-value for  $(s,a)$  away from the valid value that  $A_0$  initially set for this state-action Q-value after directly experiencing it. This indicates that for this case and, to a lesser extent, the case where there are 3 agents, that it is more beneficial to use the Independent learning method. However, the MA method outperforms Independent learning as we increase the number of agents. When there are 3 or

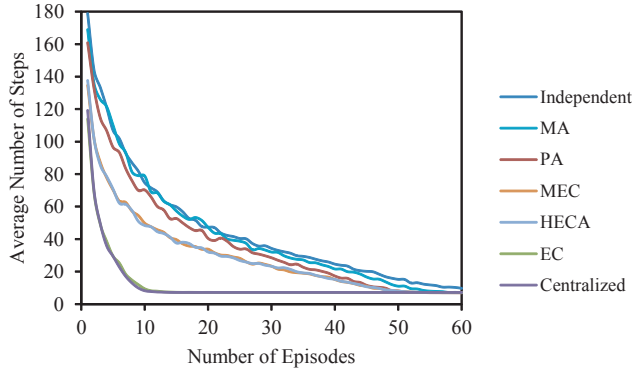


Figure 4. Learning performance for  $A_0$  when the other agents use the Independent sharing method

more agents for  $A_0$  to interact with, this trend holds true.

The performance of the MEC and HECA methods are not significantly different. This result indicates that modifying the MEC method to average the Q-values associated with the tied experience values does not yield any learning performance gains. Both methods outperform the Independent, MA, and PA methods but underperform the EC and Centralized methods. Similar to the reasoning why the reciprocating PA method outperforms the non-reciprocating MA method, the non-reciprocating MEC method cannot synchronize its Q-table and E-table values with the other agents. Therefore, action-selection decisions are made by each agent operating with potentially different Q-tables. Because the Q-tables are not synchronized, the agents learn at a slower pace than they would using the EC method, but still outperform the Independent method by more than 15 average steps per episode.

The results of this experiment indicate that it is possible for an agent to improve learning performance in environments where other agents do not reciprocate with it. This experiment shows that the MEC and HECA methods perform well when the other agent’s policies have different values and do not synchronize to the same policy after each time step  $t$ .

## B. Experiment 2

In this experiment, we explore how the proposed algorithms perform for  $A_0$  when agents  $[A_1, \dots, A_{N-1}]$  are using the PA sharing method. Results are shown in Fig. 5. Table II shows that, according to t-tests, the differences in learning performance among all methods, again excepting the MEC and HECA methods, are statistically significant. For this experiment, the MEC and HECA methods perform well when the other agent’s policies have similar values and synchronize to the same policy after each time step  $t$ . The PA method outperforms the MA method, but to a lesser extent when compared to the results from experiment 1. Fig. 6 shows the learning performance rates associated with this



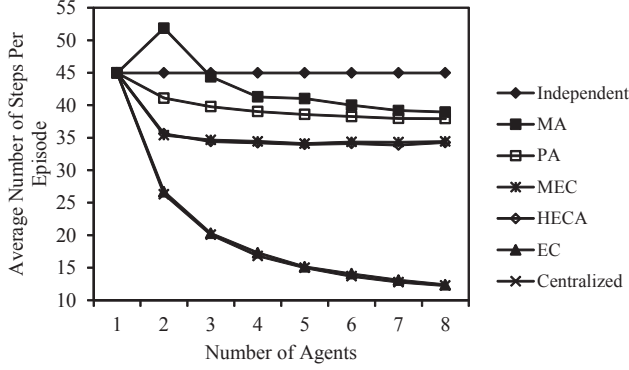


Figure 5. Simulation results for  $A_0$  when the other agents use the PA sharing method

Table II  
SIMULATION RESULTS FOR  $A_0$  WHEN THE OTHER AGENTS USE THE PA SHARING METHOD

Sharing Method	Average Steps with 95% CI
Independent	$44.98 \pm (0.33)$
MA	$38.94 \pm (0.34)$
PA	$37.94 \pm (0.30)$
MEC	$34.42 \pm (0.26)$
HECA	$34.27 \pm (0.26)$
EC	$12.37 \pm (0.11)$
Centralized	$12.24 \pm (0.11)$

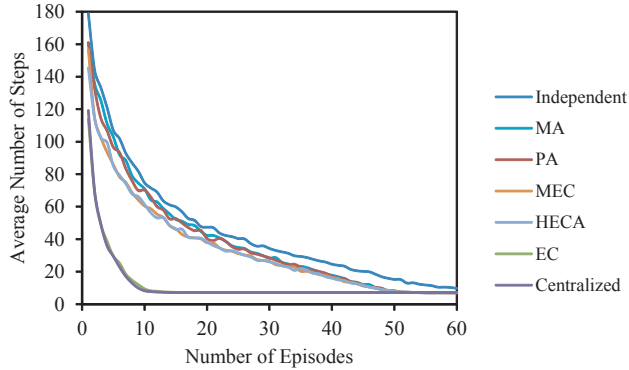


Figure 6. Learning performance for  $A_0$  when the other agents use the PA sharing method

experiment for  $A_0$ .

The experiment confirms that the same learning performance trends occur, even for a different experimental setup. Both experiment 1 and 2 confirm that the proposed methods can suitably adapt to and perform during both types of situations – synchronized and unsynchronized Q-tables – encountered for sharing strategies. Additionally, the methods perform well when faced with single or multiple Q-value changes within the Q-tables per time step.

### C. Experiment 3

In this experiment, we explore how the proposed algorithms perform for  $A_0$  when agents  $[A_1, \dots, A_{N-1}]$  are

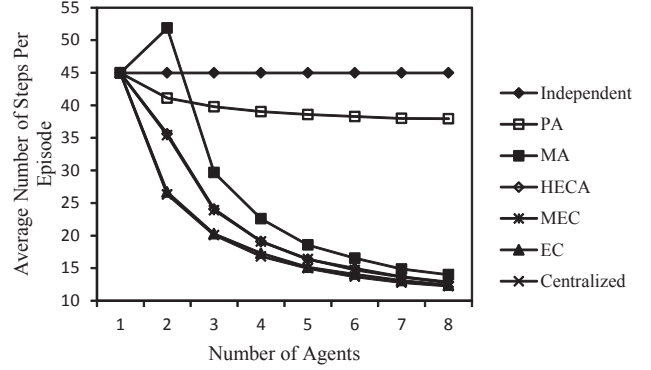


Figure 7. Simulation results for  $A_0$  when the other agents use the EC sharing method

Table III  
SIMULATION RESULTS FOR  $A_0$  WHEN THE OTHER AGENTS USE THE EC SHARING METHOD

Sharing Method	Average Steps with 95% CI
Independent	$44.98 \pm (0.33)$
PA	$37.94 \pm (0.30)$
MA	$13.99 \pm (0.17)$
HECA	$12.86 \pm (0.12)$
MEC	$12.75 \pm (0.12)$
EC	$12.37 \pm (0.11)$
Centralized	$12.24 \pm (0.11)$

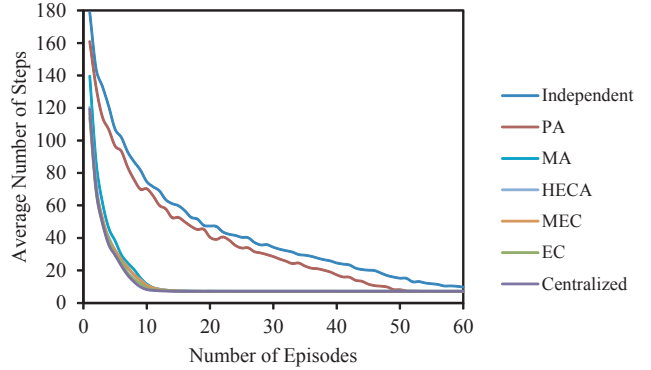


Figure 8. Learning performance for  $A_0$  when the other agents use the EC sharing method

using the EC sharing method. Results are shown in Fig. 7. Table III shows that, according to t-tests, the differences in performance among all listed methods, excepting the MEC and HECA methods, are statistically significant. The learning performance trends for this experiment differs between the PA and MA methods. The MA method significantly outperforms the PA reference method. This is a result caused by the other 7 agents employing the EC method and thus providing access to improved Q-tables for the MA method at each time step. Fig. 8 shows the learning performance rates associated with this experiment for  $A_0$ .

All non-reference methods perform significantly better in this experiment when compared to experiments 1 and

2 because the EC sharing method being used by agents  $[A_1, \dots, A_{N-1}]$  is an effective sharing strategy. Both the MEC and HECA methods perform well when the other agent's policies have similar values and synchronize to the same policy after each time step  $t$ .

## VII. CONCLUSIONS AND FUTURE WORK

Results from the 3 experiments show that it is possible for an agent to improve learning performance in environments where other agents do not reciprocate with it. The 3 proposed methods – MA, MEC, and HECA – all perform better than the Independent learning method. This indicates that by exploiting the pre-existing agent interface, learning performance can be expedited. The proposed methods can suitably adapt to and perform during both types of situations – synchronized and unsynchronized Q-tables – encountered for sharing strategies. Additionally, the methods perform well when faced with single or multiple Q-value changes within the Q-tables per time step. Most notably, the MEC and HECA methods perform best overall across the 3 experiments, with no significant difference in learning performance amongst the 2 methods.

This research opens up the possibility for future work. Currently, our method assumes that agents can communicate with and query the Q-tables of all agents in the environment, regardless of physical distance between the agents. In real situations where agents are robots and take on a physical form, this assumption may not hold. Future work can study how the methods proposed in this paper are affected by a limited communication range and propose modifications for these methods that perform better in such conditions. In this work we also assume that agents employing the proposed methods have enough memory to store the Q-tables and E-tables for all other agents in the environment. Again, if the agents are physically embodied agents, this assumption of a large memory may prove to be unrealistic. Future work can study how the methods proposed in this paper are affected by a limited memory size and propose modifications for these methods that will enable the methods to perform in such environments. An implicit assumption with this work is that agents know how many other agents are in the environment in addition to assuming a reliable communication channel among the agents where Q-tables queries are always satisfied. The information bandwidths associated with our methods can be taken into account. Our method assumes unlimited bandwidth; however, for some RL applications that are used with physical agents, bandwidth is a limiting factor.

## REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Introduction to Reinforcement Learning*, 1st ed. Cambridge, MA, USA: MIT Press, 1998.
- [2] M. Tan, "Multi-agent reinforcement learning: Independent versus cooperative agents," in *ICML*. Morgan Kaufmann, 1993, pp. 330–337.
- [3] H. R. Berenji and D. Vengerov, "Advantages of cooperation between reinforcement learning agents in difficult stochastic problems," in *In Proceedings of 9th IEEE International Conference on Fuzzy Systems*, 2000.
- [4] Y. Yang, Y. Tian, and H. Mei, "Cooperative q learning based on blackboard architecture," in *Proceedings of the 2007 International Conference on Computational Intelligence and Security Workshops*, ser. CISW '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 224–227. [Online]. Available: <http://dx.doi.org/10.1109/CIS.Workshops.2007.36>
- [5] R. Ribeiro, A. P. Borges, and F. Enembreck, "Interaction models for multiagent reinforcement learning," in *CIMCA/IAWTIC/ISE*, M. Mohammadian, Ed. IEEE Computer Society, 2008, pp. 464–469.
- [6] P. Zhang, X. Ma, Z. Pan, X. Li, and K. Xie, "Multi-agent cooperative reinforcement learning in 3d virtual world," in *ICSI (1)*, ser. Lecture Notes in Computer Science, Y. Tan, Y. Shi, and K. C. Tan, Eds., vol. 6145. Springer, 2010, pp. 731–739.
- [7] M. N. Ahmadabadi, M. Asadpour, and E. Nakano, "Cooperative q-learning: the knowledge sharing issue," *Advanced Robotics*, vol. 15, no. 8, pp. 815–832, 2001.
- [8] M. N. Ahmadabadi and M. Asadpour, "Expertness based cooperative q-learning," *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, vol. 32, no. 1, pp. 66–76, 2002.
- [9] M. N. Ahmadabadi, A. Imanipour, B. N. Araabi, M. Asadpour, and R. Siegart, "Knowledge-based extraction of area of expertise for cooperation in learning," in *IROS*. IEEE, 2006, pp. 3700–3705.
- [10] B. N. Araabi, S. Mastoureshgh, and M. N. Ahmadabadi, "A study on expertise of agents and its effects on cooperative q-learning," *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, vol. 37, no. 2, pp. 398–409, 2007.
- [11] L. Torrey and M. E. Taylor, "Help an agent out: Student/teacher learning in sequential decision tasks," in *AAMAS*. (In Press), 2012.
- [12] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: a survey," *Journal of Artificial Intelligence Research*, vol. 4, pp. 237–285, 1996.