

# Investigation of Independent Reinforcement Learning Algorithms in Multi-Agent Environments

Ken Ming Lee<sup>1\*</sup>, Sriram Ganapathi Subramanian<sup>1</sup>, Mark Crowley<sup>1</sup>

<sup>1</sup>University of Waterloo, Canada

*Submitted to Journal:*  
Frontiers in Artificial Intelligence

*Specialty Section:*  
Machine Learning and Artificial Intelligence

*Article type:*  
Original Research Article

*Manuscript ID:*  
805823

*Received on:*  
31 Oct 2021

*Revised on:*  
11 Aug 2022

*Journal website link:*  
[www.frontiersin.org](http://www.frontiersin.org)

### *Conflict of interest statement*

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

### *Author contribution statement*

KML implemented the algorithms, ran the experiments, analyzed the results and wrote the paper. SS assisted in designing the experiments and writing the paper, while MC provided access to computing resources for running the experiments. Both SS and MC also provided numerous advice and feedback during the entire process, and assisted in polishing the paper.

### *Keywords*

Multi-agent reinforcement learning, reinforcement learning, deep learning, machine learning, artificial intelligence

### *Abstract*

Word count: 166

Independent reinforcement learning algorithms have no theoretical guarantees for finding the best policy in multi-agent settings. However, in practice, prior works have reported good performance with independent algorithms in some domains and bad performance in others. Moreover, a comprehensive study of the strengths and weaknesses of independent algorithms is lacking in the literature. In this paper, we carry out an empirical comparison of the performance of independent algorithms on seven PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. For the cooperative setting, we show that independent algorithms can perform on par with multi-agent algorithms in fully-observable environments, while adding recurrence improves the learning of independent algorithms in partially-observable environments. In the competitive setting, independent algorithms can perform on par or better than multi-agent algorithms, even in more challenging environments. We also show that agents trained via independent algorithms learn to perform well individually, but fail to learn to cooperate with allies and compete with enemies in mixed environments.

### *Contribution to the field*

Independent reinforcement learning algorithms have no theoretical guarantees for finding the best policy in multi-agent settings. However, in practice, prior works have reported good performance with independent algorithms in some domains and bad performance in others. Moreover, a comprehensive study of the strengths and weaknesses of independent algorithms is lacking in the literature. In this paper, we carry out an empirical comparison of the performance of independent algorithms on seven PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. For the cooperative setting, we show that independent algorithms can perform on par with multi-agent algorithms in fully-observable environments, while adding recurrence improves the learning of independent algorithms in partially-observable environments. In the competitive setting, independent algorithms can perform on par or better than multi-agent algorithms, even in more challenging environments. We also show that agents trained via independent algorithms learn to perform well individually, but fail to learn to cooperate with allies and compete with enemies in mixed environments.

### *Funding statement*

This project was made possible through grants from the Canada Wildfire Strategic Network and Discovery Grant Program of the National Research Council of Canada (NSERC).

*Ethics statements*

*Studies involving animal subjects*

Generated Statement: No animal studies are presented in this manuscript.

*Studies involving human subjects*

Generated Statement: No human studies are presented in this manuscript.

*Inclusion of identifiable human data*

Generated Statement: No potentially identifiable human images or data is presented in this study.

*Data availability statement*

Generated Statement: The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

In review

# Investigation of Independent Reinforcement Learning Algorithms in Multi-Agent Environments

Ken Ming Lee<sup>1,\*</sup>, Sriram Ganapathi Subramanian<sup>1</sup> and Mark Crowley<sup>1</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON, Canada

Correspondence\*:

Ken Ming Lee  
kenming.lee@uwaterloo.ca

## 2 ABSTRACT

Independent reinforcement learning algorithms have no theoretical guarantees for finding the best policy in multi-agent settings. However, in practice, prior works have reported good performance with independent algorithms in some domains and bad performance in others. Moreover, a comprehensive study of the strengths and weaknesses of independent algorithms is lacking in the literature. In this paper, we carry out an empirical comparison of the performance of independent algorithms on ~~four~~<sup>seven</sup> PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. ~~We show that For the cooperative setting, we show that independent algorithms can perform on par with multi-agent algorithms in fully-observable environments, independent algorithms while adding recurrence improves the learning of independent algorithms in partially-observable environments. In the competitive setting, independent algorithms can perform on par with or better than multi-agent algorithms in cooperative and competitive settings. For the mixed environments, we even in more challenging environments. We also show that agents trained via independent algorithms learn to perform well individually, but fail to learn to cooperate with allies and compete with enemies . We also show that adding recurrence improves the learning of independent algorithms in cooperative partially observable in mixed environments.~~

Keywords: Multi-Agent Reinforcement Learning, Reinforcement Learning, Deep Learning, Machine Learning, Artificial Intelligence

[Number of words: ~~4484~~<sup>4802</sup>, Number of figures and tables: ~~6~~<sup>12</sup>]

## 0 INTRODUCTION

One of the simplest ways to apply reinforcement learning in multi-agent settings is to assume that all agents are independent of each other. In other words, every other agent is seen as part of the environment from any agent's perspective. Independent algorithms (i.e., single-agent algorithms) face the issue of non-stationarity in the multi-agent domain due to the violation of the Markovian property in a Markov Decision Process (Choi et al., 2000). As a result, independent algorithms lack convergence guarantees, and are not theoretically sound in the multi-agent setting (Tan, 1993). Despite these shortcomings, independent algorithms have the advantage of requiring

28 lower computational resources and being easier to scale to large environments than traditional  
 29 multi-agent algorithms which perform exact opponent modelling of each agent. In practice, prior  
 30 works have reported mixed performance for independent algorithms in different multi-agent domains  
 31 (Zawadzki et al., 2014; Tampuu et al., 2017; Shehama and Leyton-Brown, 2008; Berner et al., 2019; Foerster et al.,  
 32 (Berner et al., 2019; Foerster et al., 2018; Lowe et al., 2017; Rashid et al., 2018; Shoham and Leyton-Brown, 2008  
 33 . However, a study of the strengths and weaknesses of independent algorithms across various categories  
 34 within the multi-agent domain is lacking in the literature.

35 In this paper, we investigate the empirical performance of independent algorithms in multi-agent  
 36 settings, and compare them to various multi-agent algorithms under the Centralized Training and  
 37 Decentralized Execution scheme (Kraemer and Banerjee, 2016; Oliehoek et al., 2008). We evaluate  
 38 these algorithms on 47 multi-agent environments from the PettingZoo library (Terry et al., 2020b),  
 39 which span the 3 main categories of multi-agent environments (i.e., cooperative, competitive,  
 40 and mixed) (Busoniu et al., 2008; Canese et al., 2021; Zhang et al., 2021; Gronauer and Diepold, 2021)  
 41 (Busoniu et al., 2008; Canese et al., 2021; Gronauer and Diepold, 2021; Zhang et al., 2021). We show  
 42 that independent algorithms can perform on par with multi-agent algorithms in the cooperative, fully-  
 43 observable setting, and adding recurrence allows them to perform well compared to multi-agent algorithms  
 44 in partially observable environments. In the competitive setting, we show that parameter sharing alongside  
 45 the addition of agent indicators allow independent algorithms to outperform some perform on par or  
 46 better than multi-agent algorithms, such as Multi-Agent Proximal Policy Optimization (Yu et al., 2021)  
 47 , and Multi-Agent Deep Deterministic Policy Gradient (Lowe et al., 2017), in fully-observable even in  
 48 challenging environments. For the mixed setting, we show that agents of independent algorithms learn to  
 49 perform well individually, but fail in learning to cooperate with allies and compete against enemies.

## 1 BACKGROUND INFORMATION

50 In this section, we provide readers with a brief overview of the various concepts and algorithms that are  
 51 used throughout the paper.

### 52 1.1 Reinforcement Learning

53 In Reinforcement Learning (RL), an agent interacts with the environment by making sequential decisions  
 54 (Sutton and Barto, 2018). At every time step, denoted as  $t$ , the agent observes a state  $s_t$  from the  
 55 environment, and takes an action  $a_t$ . This action is executed in the environment, which returns a reward  
 56  $r_t$  and the next state  $s_{t+1}$  that are determined by the reward function  $R(s_t, a_t)$  and the transition  
 57 probability,  $P(s_{t+1}|s_t, a_t)$ , respectively. Critically,  $R(s_t, a_t)$  and  $P(s_{t+1}|s_t, a_t)$  are part of the environment, and are usually unknown to the agent of a model-free  
 58 RL algorithm. Since the transition probability  $P(s_{t+1}|s_t, a_t)$  conditions the next state  $s_{t+1}$   
 59 purely on the current state  $s_t$  and action  $a_t$ , it satisfies the Markov property (Markov, 1954). This  
 60 interaction between the agent and the environment is called a Markov Decision Process (MDP) (Bellman,  
 61 1957). The objective of an RL agent is to learn a policy  $\pi(a_t|s_t)$ , which maps a state to an action  
 62 that maximizes the expected cumulative reward it receives from the environment. This is written as  $\sum_t \gamma^t r_t$ ,  
 63 where  $\gamma \in [0, 1]$  represents a discount factor on future rewards.

### 65 1.2 Multi-Agent Reinforcement Learning

66 The single-agent MDP framework is extended to the Multi-Agent Reinforcement Learning (MARL)  
 67 setting in the form of stochastic games (Shapley, 1953). In an  $N$ -agent stochastic game, at every time step,

68 each of the  $n$  agents, identified by  $j \in \{1, 2, \dots, n\}$  across all agents, takes an action  $u_t^j$ . The joint action  
 69  $\mathbf{u}_t \triangleq \{u_t^1, \dots, u_t^N\}$  determines the rewards obtained by each agent is written as  $\mathbf{u}_t \triangleq (u_t^1, \dots, u_t^N)$ . Every  
 70 agent receives an agent specific reward through the reward function  $R(s_t, \mathbf{u}_t, j)$ . State transitions of the  
 71 environment are determined by the transition probability  $P(s_{t+1}|s_t, \mathbf{u}_t)P(s_{t+1}|s_t, \mathbf{u}_t)$ , which conditions  
 72 on the state and the joint action at timestep  $t$ .

### 73 1.3 Centralized Training and Decentralized Execution

74 While it is technically possible to learn a centralized controller that maps a state to a distribution over the  
 75 joint action space, the number of possible combinations of actions grows exponentially with the number of  
 76 agents. This makes centralized control intractable for environments with many agents. Therefore, this paper  
 77 is mainly focused on multi-agent algorithms which correspond to a Centralized Training and Decentralized  
 78 Execution (CTDE) scheme (Kraemer and Banerjee, 2016; Oliehoek et al., 2008). A CTDE algorithm has  
 79 two phases. During the control phase, where policies are deployed in the environment, rather than using  
 80 a centralized controller to take actions for all agents, decentralized agents make decisions based on their  
 81 individual observations. During the prediction phase, centralized training is performed, which allows for  
 82 extra information (e.g., the state) to be utilized, as long as this is not required during the control phase.

### 83 1.4 Cooperative, Competitive and Mixed

84 This paper follows the convention of classifying every multi-agent algorithm and environment  
 85 studied into one of three categories – cooperative, competitive, or-and mixed (cooperative-  
 86 competitive) (Busoniu et al., 2008; Canese et al., 2021; Zhang et al., 2021; Gronauer and Diepold, 2021)  
 87 (Busoniu et al., 2008; Canese et al., 2021; Gronauer and Diepold, 2021; Zhang et al., 2021).

88 In the cooperative setting, agents collaborate with each other to achieve a common goal. As a  
 89 result, it is very common for all agents to share the same reward function (i.e., a team goal) (Chang  
 90 et al., 2004). Also known as the multi-agent credit assignment problem, every agent has to deduce its  
 91 own contributions from the team reward (Chang et al., 2004). Algorithms studied in this paper that  
 92 explicitly address the multi-agent credit-assignment problem include QMIX (Rashid et al., 2018) and  
 93 an on-policy algorithm – Counterfactual Multi-Agent Policy Gradients (COMA) (Foerster et al., 2018),  
 94 and an off-policy algorithm, QMIX (Rashid et al., 2018), which addresses the poor sample efficiency of  
 95 on-policy algorithms. Additionally, the CommNet (Sukhbaatar et al., 2016) extension on top of COMA  
 96 is utilized for specific cooperative environments to promote communication between cooperative agents.  
 97 Other multi-agent algorithms that are considered for the cooperative scenario include multi-agent variants  
 98 of single-agent algorithms, such as Multi-Agent Deep Deterministic Policy Gradient (MADDPG) (Lowe  
 99 et al., 2017) and Multi-Agent Proximal Policy Optimization (MAPPO) (Yu et al., 2021).

100 In the competitive setting, agents play a zero-sum game, where one-an agent's gain is another agent's  
 101 loss. In other words,  $\sum_a r(s, u, a) = 0 \forall s, u$ . Algorithms that are studied specifically in this paper include  
 102 Deep Reinforcement Opponent Network (DRON) (He et al., 2016), MADDPG and MAPPO. MADDPG  
 103 and MAPPO learn a separate critic for every agent, which gives the algorithms flexibility to learn different  
 104 behaviours for agents with different reward functions.

105 In a mixed or cooperative-competitive setting, environments are neither zero-sum (competitive) nor  
 106 cooperative, and they do not necessarily need to be general-sum either. A common setting would be  
 107 environments that require every agent to cooperate with some agents, and compete with others (Busoniu  
 108 et al., 2008; Canese et al., 2021; Zhang et al., 2021). MADDPG and MAPPO are used here for the same  
 109 reason as the competitive setting.

## 110 1.5 Independent Algorithms and Non-Stationarity

111 One naive approach for applying single-agent RL to the multi-agent setting would be the use of  
112 independent learners, where each agent treats every other agent as part of the environment, and learns  
113 purely based on individual observations. In a multi-agent setting, the transition probability  $P$  and  
114 reward function  $R$  are conditioned on the joint action  $u$ . Since all agents in the environment are  
115 learning, their policies constantly change. Therefore, from each independent learner's perspective,  
116 the transition probability and reward function appear non-stationary, due to the lack of awareness of  
117 other agents' actions. This violates the Markovian property of an MDP, which then causes independent  
118 algorithms to be slow to adapt to other agents' changing policies, and as a result, face difficulties  
119 in converging to a good policy (Papoudakis et al., 2019; Hernandez-Leal et al., 2017; He et al., 2016)  
120 (He et al., 2016; Hernandez-Leal et al., 2017; Papoudakis et al., 2019).

121 In this paper, we chose to use a popular off-policy algorithm, Deep Q-Network (DQN) (Mnih et al.,  
122 2015), and an on-policy algorithm, Proximal Policy Optimization (PPO) (Schulman et al., 2017). In specific  
123 partially observable environments, Deep Recurrent Q-Network (DRQN) (Hausknecht and Stone, 2015) is  
124 also utilized. Following the work of Gupta et al. (2017), parameter sharing is utilized for all independent  
125 algorithms, such that experiences from all agents are trained simultaneously using a single network. This  
126 allows the training to be more efficient (Gupta et al., 2017). The aforementioned independent algorithms  
127 are tested in all 3 categories of multi-agent environments.

## 2 EXPERIMENTAL SETUP

128 In this section, we introduce the environments used for the experiments, specify the various preprocessing  
129 that were applied, and other relevant implementation details.

### 130 2.1 Environments

131 The experiments were performed on multiple multi-agent environments from the PettingZoo library  
132 (Terry et al., 2020b), which contains the Multi-Agent Particle Environments (MPE) (Lowe et al., 2017;  
133 Mordatch and Abbeel, 2017) and multi-agent variants of the Atari 2600 Arcade Learning Environment  
134 (ALE) (Bellemare et al., 2013; Terry and Black, 2020).

135 For the cooperative setting, experiments were performed on a modified version of the 2-player Space  
136 Invaders (Bellemare et al., 2013; Terry and Black, 2020), and the Simple Reference MPE environment  
137 (Lowe et al., 2017; Mordatch and Abbeel, 2017). In Space Invaders, both agents share the common goal of  
138 shooting down all aliens. To make Space Invaders cooperative, we removed the (positive) reward that is  
139 given to a player whenever the other player gets hit. Additionally, the environment rewards every agent  
140 individually by default. Since a number of cooperative multi-agent algorithms (e.g., QMIX and COMA)  
141 assume that only a team reward is given, we modified the reward function such that a team reward is  
142 given instead (i.e., both agents receive the sum of their individual rewards). This setup exemplifies the  
143 multi-agent credit assignment problem, the effect of which is studied more closely in the Section 3.1.1.  
144 On the other hand, in the Simple Reference environment, two agents are rewarded by how close they are  
145 to their target landmark. However, the target landmark of an agent is only known by the other agent, as a  
146 result communication is required for both agents to navigate successfully to their target landmarks.

147 For the competitive setting, we performed experiments on the three 2-player variant of the original Atari  
148 Pong environment zero-sum competitive games from the Atari suite – Boxing, Pong and Space War. For  
149 the mixed setting, we opted for the Simple Tag MPE environment, which and the Simple Adversary

150 MPE environments. Simple Tag is a Predator-Prey environment (Mordatch and Abbeel, 2017). This  
151 environment consists of 4 agents—that consists of 3 predators and a prey (Mordatch and Abbeel, 2017).  
152 The prey travels faster and has to avoid colliding with the predators, while the 3 predators travel slower  
153 and have to work together to capture the prey. The rewards received by the prey and a predator sum to  
154 0 (i.e., the prey gets a negative reward for collision, while the predators get rewarded positively), and  
155 all predators receive the same reward. The prey is also negatively rewarded if it strays away from the  
156 predefined area (a  $1 \times 1$  unit square). This environment is general-sum because it contains 3 predators and a  
157 single prey. On the other hand, Simple Adversary consists of 2 cooperative agents and an adversary agent.  
158 Both cooperative agents receive equal reward, based on the distance of the closest agent to the target and  
159 the negative distance of the adversary to the target. The adversary is also rewarded based on its distance to  
160 the target, but unlike the cooperative agents, it has to infer the location of the target based on the location  
161 of the cooperative agents, which it observes. Therefore, the cooperative agents have to cooperate in order  
162 to split up to get close to the target while deceiving the adversary away from the target. Similarly, the  
163 Simple Adversary environment is also general-sum.

## 164 2.2 Preprocessing

165 For the MPE environments, no preprocessing was done, and default environment-parameters were used  
166 for all MPE experiments.

167 For the Atari environments, following the recommendations of Machado et al. (2018), we performed the  
168 following preprocessing:

- 169 • Reward clipping to ensure that the rewards at every timestep were clipped between the range of [-1, 1].  
170 • Sticky actions with a probability of 0.25.  
171 • Frame skip of 4.

172 The number of steps per episode was set to a limit of 200 for both all Atari environments, as that yielded  
173 the best results in general. Subsequently Additionally, no-op resets were also performed on the first 130  
174 frames for Space Invaders, and the first 60 frames for Pong.

175 Furthermore, the action spaces for both Atari environments were shrunk to their effective action spaces  
176 in order to improve learning efficiency. For Pong all competitive environments specifically, we also  
177 concatenated a one-hot vector of the agent's index to the observations so that independent algorithms can  
178 differentiate one from the other when parameter sharing is utilized. The effect of this addition is studied  
179 more closely in Section 3.5.

180 All preprocessing were performed using the SuperSuit library (Terry et al., 2020a).

## 181 2.3 Implementation

182 Implementations of all algorithms were based on the following open-sourced libraries/reference  
183 implementations:

- 184 • Implementation of DQN and DRON were based on the Machin library (Li, 2020).  
185 • Implementation of independent PPO was based on Stable Baselines3 (Raffin et al., 2019).  
186 • Implementation of DRQN, QMIX, COMA and CommNet came from a popular public repository by  
187 the name of StarCraft (starry sky6688, 2019).  
188 • Implementation of MADDPG came from the original code implementation (Lowe et al., 2017).

189 • Implementation of MAPPO came from the original code implementation (Yu et al., 2021).  
190 For both DQN and DRON, the underlying DQN implementations included Double DQN (Van Hasselt  
191 et al., 2016), the dueling architecture (Wang et al., 2016) and priority experience replay buffer (Schaul  
192 et al., 2015). On the other hand, the implementation of DRQN did not use any of the aforementioned  
193 add-ons. For PPO and MAPPO, 4 parallel workers were used for all environments with homogeneous state  
194 and action spaces. Default hyperparameters were used for all algorithms, and no hyperparameter tuning  
195 was performed. Details of the hyperparameters used can be found in the Supplementary Material.

196 All experiments were performed across 5 different seeds. Parameter sharing was utilized for all algorithms  
197 throughout the experiments for all environments with homogeneous state and action spaces. For multi-agent  
198 algorithms that perform centralized training (e.g., QMIX, COMA, MADDPG etc.), the global states were  
199 represented by the concatenation of all agents' local observations. We also used the 128-byte Atari RAM  
200 as state inputs, rather than visual observations. This allows the algorithms to focus their learning on control  
201 rather than on both control and perception, improving learning efficiency.

202 For the mixed setting, the observation and action spaces differ between predators and the prey in the  
203 Simple Tag environment, while the observation spaces of the cooperative agents differ from the adversary  
204 agent in the Simple Adversary environment. As a result, in both environments, none of the agents have  
205 their parameters shared. Parameters between the predators/cooperative agents are also not shared to ensure  
206 that no bias is introduced (since they would have more data to learn from compared to opposing agents in  
207 their respective environments).

### 3 EXPERIMENTAL RESULTS AND DISCUSSION

208 In this section, we highlight the experiments performed on ~~the four multi-agent environments (i.e., Simple~~  
209 ~~Reference, Space Invaders, Pong and Simple Tag), and all environments, and~~ provide discussions about  
210 the obtained results.

#### 211 3.1 Cooperative

212 We ran the various algorithms on the Simple Reference environment for 240k episodes ( $6 \times 10^6$  steps).  
213 From ~~figure~~ Figure 1A, it could be observed that all independent algorithms converged to a lower score,  
214 except for DRQN, whose recurrence allowed it to vastly outperform DQN and converge to a score on  
215 par with multi-agent algorithms. However, this trend was not observed when comparing MAPPO to its  
216 recurrent variant (i.e., RMAPPO), as MAPPO performs equally well as RMAPPO. We hypothesize that  
217 since MAPPO's centralized critic learns based on the joint observation and action of both agents, this  
218 minimizes the amount of partial observability of every agent, and allows each agent to learn to communicate  
219 with other agents effectively without recurrence. In contrast, for independent algorithms, such as DQN,  
220 where the interactions between the agents are not explicitly learned (since all other agents are treated as  
221 part of the environment), adding recurrence could help mitigate some resulting partial observability, hence  
222 improving their performance, as described above.

223 Unlike the Simple Reference environment, the Space Invaders environment seemed to favour non-  
224 recurrent variants of algorithms (~~figure~~ Figure 1B). MAPPO vastly outperformed RMAPPO, and similarly  
225 DQN outperformed DRQN. This is also likely the underlying reasoning behind the comparatively poorer  
226 performance of the multi-agent algorithms, such as QMIX, COMA and CommNet, all of which were  
227 implemented with recurrent neural networks under the CTDE scheme.

228 Additionally, since there is no unit collision in the Space Invaders environment (i.e., agents can move past  
229 each other without being blocked), they do not have to coordinate between themselves to achieve a high  
230 score in the environment; a good policy can be learned solely by having agents maximize their individual  
231 rewards. This explains the strong performance that was achieved by DQN. Also, since this is a cooperative  
232 task with both agents having identical goals, learning separate representations for individual agents is not  
233 very important; the learning of both agents assist each other. This is shown in figure Figure 8B in Section  
234 3.5, where the addition of an agent indicator did not yield any performance improvement for DQN on  
235 Space Invaders.

236 Given such circumstances, it is interesting to observe the stronger performance of MAPPO compared to  
237 the independent algorithms. By conditioning on the joint action, MAPPO’s critic has full observability into  
238 the joint action that resulted in the team reward. Therefore, the observed reward is unbiased, which allows  
239 the learning process to be more efficient. In contrast, independent algorithms have to learn from a noisy  
240 team reward signal, where an agent could receive a large positive team reward even when it did nothing.  
241 This relates to the problem of credit assignment in MARL, noted in prior works (Hernandez-Leal et al.,  
242 2019).

### 243 3.1.1 Multi-Agent Credit Assignment Problem in Fully Observable Settings

244 In this section, we attempt to study the effect of using a team reward signal, rather than individual reward  
245 signals on various independent and multi-agent algorithms in a fully observable environment. When team  
246 rewards are the only rewards given, these reward signals are noisy for independent algorithms because  
247 the agent, which treats every other agent as part of the environment, does not know the actions taken  
248 by other agents. This makes it difficult for independent algorithms’ agents to learn how their individual  
249 actions contribute to the team reward signal. We performed the experiments on Space Invaders, in which  
250 the default agents receive individual rewards from the environment. To study the effect of the multi-agent  
251 credit assignment problem, we performed two runs per algorithm, one with team rewards only, and the  
252 other with individual rewards only (i.e., agents are rewarded independently by the environment).

253 For multi-agent algorithms, such as MAPPO (figure Figure 2B) and RMAPPO (figure Figure 2C), having  
254 a team reward does not have a large effect on the performance of the algorithms. This is expected because  
255 these algorithms have critics that learn from the joint action, which allow them to implicitly learn the  
256 estimated contribution of every agent without factorization.

257 On similar lines, regarding independent algorithms, we observe that having team rewards instead of  
258 individual ones do not impact their performance adversely (figure Figure 2A). A plausible explanation  
259 could be that since parameter sharing is utilized and all agents receive the same reward for a given joint  
260 action, this allows the independent algorithms to correlate actions from different agents that produced  
261 similar (high) rewards.

## 262 3.2 Competitive

263 The 2-player Pong environment was used for the competitive setting. All algorithms were first trained  
264 using parameter sharing with the addition of agent indicators (the effect of which is detailed in Section  
265 3.5) for 60k episodes ( $1.2 \times 10^6$  steps), their network parameters were then saved. Since Pong is a  
266 zero-sum game, we evaluated them by putting them For every competitive environment, all algorithms  
267 were trained for a fixed number of steps ( $1.2 \times 10^7$ ). Performance evaluation is performed by pitting  
268 algorithms head-to-head against each other for 3 episodes for all possible combinations. After that,  
269 their positions were swapped, and the entire process was repeated. Swapping their positions is crucial

270 for evaluation, because permutations, during which no training is performed. To ensure fairness, the first  
271 player (playing the right paddle) ordering of algorithms (i.e., Algorithm A playing as player 1 or 2) are  
272 also taken into consideration. This is because in environments such as Pong, the right paddle player is  
273 always the serving player, therefore the first player always has an advantage over the second player (which  
274 plays the left paddle). This advantage is further exacerbated because the winning side always gets to serve  
275 subsequent openings thus having an advantage. At the end of an episode, the agent of an algorithm that  
276 has achieved higher cumulative reward is considered the winner. If both agents achieved the same amount  
277 of cumulative reward (draw scenario), both agents are considered to win. The entire evaluation process  
278 was is repeated across all 5 seeds.

279 From the , and is the basis behind the stacked bar charts shown in figure 4, a similar trend across  
280 the number of games won as the first and second player can be observed (figure 4 A and 4B). DRON  
281 is consistently the best player, closely followed by DQN. Both of these algorithms were also the only  
282 algorithms to have a win rate of greater than 50% for the games they have played (figure 4C). (Figure 3, 4  
283 and 5).

284 An interesting observation that can be made is the strong performance of independent algorithms,  
285 compared to other multi-agent algorithms. Since Pong is fully observable, critics that learn based on  
286 the joint observation of both agents do not necessarily provide any new information. Furthermore, since  
287 Pong is a highly reactive environment, Among all algorithms, DQN and DRON are the best performers in  
288 Boxing and Pong environment, by a large margin. DQN outperforms DRON in the Boxing environment  
289 (Figure 3), while DRON outperforms DQN in the Pong environment (Figure 4). Since both of these  
290 environments are reactive in nature, this meant that an agent can learn a good policy solely by understanding  
291 how to position its react to the situation at hand. For instance, in Pong, this meant learning to position  
292 the paddle according to the trajectory of the ball (towards the agent). While learning on the joint action  
293 could allow agents to learn to better predict the incoming trajectory of the ball, it can be observed that  
294 the additional layer of complexity causes the sample efficiency to decrease and only yields diminishing  
295 returns. Additionally, since both of these environments are fully observable, critics that learn based on the  
296 joint observation of both agents do not necessarily provide any new information.

297 In addition to the above factors, it is possible that parameter sharing benefited agents of independent  
298 algorithms by allowing them to learn better representations of both players, since they were trained to play  
299 as both players simultaneously. Had these algorithms trained without parameter sharing, there would likely  
300 be a larger performance difference between independent algorithms and opponent modelling algorithms  
301 such as DRON. Instead of treating other agents as part of the environment, opponent modelling allows  
302 agents to adapt more quickly to the opponent's changing strategies (He et al., 2016). In the Space Wars  
303 environment, two agents shoot missiles at each other. Critically, since the missiles travel faster than the  
304 agents and the missiles can ricochet off walls, prediction and positioning are key for shooting and dodging  
305 effectively. This makes the environment challenging. In this environment, RMAPPO performed the best,  
306 closely followed by MADDPG and DQN (Figure 5). By conditioning on past trajectories, RMAPPO  
307 was able to better learn the mechanics of the environment (such as ricocheting missiles), allowing it  
308 to have better aim than other algorithms. However, this only yields marginal improvement over the  
309 MADDPG and DQN. Even with the complexity of the environment, DQN's performance was nearly on  
310 par with MADDPG, and performed better than the rest of the minimal improvement DRON has over DQN  
311 suggests that in the Pong environment, an agent's policy may not be significantly affected by changes in  
312 the opponent agent's policy (i.e., individual agents can play the same way regardless of how their opponent  
313 played). multi-agent algorithms, including DRON and MAPPO.

314 **3.3 Mixed**

315 In the Simple Tag (i.e.,

316 In the Simple Tag (i.e., Predator-Prey) environment, the predators are incentivized to cooperate together  
317 to trap the prey, while the prey is incentivized to dodge the predators while staying within a predefined  
318 area. For our method of evaluation, we plot the training curves of the prey (figure 6B), and one of the  
319 predators (figure 6A), since all predators receive the same reward. Since the observation and action spaces  
320 differ between the predators and the prey, none of the agents have their parameters shared. We chose not  
321 to share the parameters of the predators to ensure that bias towards the predators was not introduced (since  
322 they would have 3 times the amount of data to learn from compared to the prey).

323 In the case of DQN, the environment, DQN's prey successfully learned to minimize the number of  
324 collisions with the predators, which can be observed by the strong performance achieved by the prey (figure  
325 Figure 6B). However, similar to PPO, since the predators were trained completely independently (i.e., their  
326 parameters were not shared), they did not manage to learn how to cooperate with one another to capture the  
327 prey (figure Figure 6A). It is interesting to observe that MADDPG converged to a policy similar to DQN,  
328 with the difference being that its predators have learned to cooperate better, thus getting slightly higher  
329 rewards compared to DQN's predators (figure Figure 6A). Subsequently, as a result of the higher rewards  
330 obtained by the predators, MADDPG achieves a slightly lower score for its prey (figure Figure 6B).

331 MAPPO and RMAPPO, on the other hand, learned a different strategy. As we can observe from the  
332 comparatively noisier curves obtained from their predators and preys (figure Figure 6A and 6B), there  
333 is a constant tug-of-war between the prey and the predators - as the predators learn how to cooperate  
334 better, their scores increase, which subsequently causes the prey to learn how to dodge, decreasing the  
335 predators' scores, and vice versa. Since the predators of MAPPO and RMAPPO achieves a much higher  
336 score compared to all other algorithms, this is indicative that the predators have successfully learned to  
337 cooperate to trap the prey.

338 Similar to the Simple Tag environment, the Simple Adversary environment exhibits similar patterns in  
339 terms of relative performances of the various algorithms. DQN and MADDPG converged to a similar  
340 policy, where the stronger performance of the adversary agent implies that the adversary agent was able  
341 to better locate the target location (Figure 7A). In other words, the two cooperative agents were less  
342 successful at deceiving the adversary agent into an incorrect target. On the other hand, MAPPO and  
343 RMAPPO's adversary agent converged to a substantially lower score (Figure 7A), but conversely the  
344 cooperative agents were able to achieve greater score than those of DQN and MADDPG (Figure 7B). This  
345 is indicative that the cooperative agents (of MAPPO and RMAPPO) were more successful in learning to  
346 cooperate by positioning close to the true target while deceiving the adversary away from it.

347 **3.4 Performance of Independent PPO**

348 Throughout all experiments performed, PPO exhibits poor performance in most of the environments  
349 across all 3 settings. While PPO has been shown to work well in a wide variety of  
350 environments (Berner et al., 2019; OpenAI et al., 2018; Schulman et al., 2017; Yu et al., 2021), PPO is  
351 highly sensitive to the choice of implementation and hyperparameters, as reported in prior works  
352 (Andrychowicz et al., 2020; Engstrom et al., 2020). We find that this problem is exacerbated in  
353 multi-agent environments, especially in the mixed settings, since those are the hardest.

### 354 3.5 Importance of Agent Indicator

355 In this section, we list some interesting findings from the addition of agent indicators to independent  
356 algorithms when utilizing parameter sharing.

357 Interestingly, in both cooperative environments, there was no noticeable improvement in the performance  
358 of DQN when an agent indicator was added (figure Figure 8A and 8B). As was previously discussed, in  
359 the case of Space Invaders, since both agents have identical goals and similar representations, there is little  
360 need to distinguish between either agent. On the other hand, due to the partially observable nature of the  
361 Simple Reference environment, DQN performed similarly poorly, regardless of whether agent indicators  
362 were present. In this case, the addition of recurrence would have resulted in a much more significant  
363 difference instead, as was previously shown.

364 Conversely, for the Pong environment, even though it is also fully observable (akin to Space Invaders), the  
365 representation of both agents are not interchangeable. Utilizing parameter sharing without agent indicators,  
366 all algorithms struggled to learn due to the inability to tell which paddle were they controlling at every  
367 timestep. The only exception was RMAPPO (figure Figure 9), which was able to condition on the sequence  
368 of previous observations and actions to infer which paddle was it controlling.

## 4 CONCLUSION

369 In this section, we provide a summary of the findings and discussions from the previous sections.

### 370 4.1 Cooperative

371 In the cooperative setting, for environments where individual agents have full observability such as  
372 Space Invaders, we showed that independent algorithms can perform even better than certain multi-  
373 agent algorithms. Furthermore, we showed that independent algorithms are able to cope well with the  
374 multi-agent credit assignment problem in environments that are fully observable with a relatively small  
375 number of agents, and where every agent has the same task. On the other hand, in the Simple Reference  
376 environment where the need for agents to communicate induces partial observability, adding recurrence  
377 allowed independent algorithms to perform as well as other multi-agent algorithms. We also discussed the  
378 significance of learning on the joint observation and action, rather than individual ones, and showed that  
379 MAPPO performs as well as DRQN in the Simple Reference environment, without the need for an RNN.  
380 Moreover, in Space Invaders, MAPPO was able to consistently achieve the highest score amongst all other  
381 algorithms.

### 382 4.2 Competitive

383 In the Boxing and Pong environment, ~~we saw that~~ DRON and DQN were able to outperform all other  
384 algorithms. We argued that this is due to the ~~fully observable nature of the Pong environment, in addition~~  
~~to the diminishing returns that learning from joint actions could yield.~~ reactive nature of both environment,  
385 which results in diminishing returns for multi-agent algorithms that learn joint actions. On the other hand,  
386 in the more complex Space Wars environment, RMAPPO performed the best as it was able to leverage on  
387 past observations to make better predictions. However, DQN still performed ~~on par or better than~~ other  
388 ~~multi-agent algorithms.~~ Furthermore, we showed that with the ~~use~~ addition of agent indicators, independent  
389 algorithms were able to learn robust policies ~~for both competing agents~~ using parameter sharing ~~in~~ Pong.

### 391 4.3 Mixed

392 In the Predator-Prey environment both mixed environments, we saw that since there were no parameter  
393 sharing to induce cooperation, predators from cooperative agents of independent algorithms were unable  
394 to learn how to cooperate with each other to capture the prey. Conversely, in DQN we saw that its prey was  
395 able to achieve the highest score consistently, showing that the prey has learned to dodge the predators  
396 effectively while staying within the predefined area. Interestingly, we also saw how MADDPG's training  
397 curve for its predators and prey shows resemblance to that of DQN, suggesting that it also faced difficulties  
398 in learning strategies for the predators to coordinate and capture the prey compete with the opposing agent.  
399 This was reflected in DQN's stronger performance as the prey in the Simple Tag environment, and the  
400 adversary in the Simple Adversary environment. Agents of MAPPO and RMAPPO, on the other hand,  
401 were the only algorithms that managed to achieve high scores for their predators, suggesting that their  
402 predators have learned how to collaborate with each other to hunt the prey. The noisiness of their able  
403 to learn to cooperate, leading to higher rewards as predators in Simple Tag, and as cooperative agents in  
404 Simple Adversary. Furthermore, the noisiness of graphs suggest that there is a constant tug-of-war between  
405 the prey and the predators both opposing parties, as one tries to outsmart the other. Interestingly, in both  
406 mixed environments, MADDPG exhibits similar characteristics to DQN, suggesting that its cooperative  
407 agents in both environments also faced difficulties in learning to cooperative.

## 5 FUTURE WORK

408 In this section, we highlight some future work that could potentially bring more insights into having a  
409 broader understanding of dealing with non-stationarity and partial observability for independent algorithms,  
410 both of which are common in the multi-agent setting. In the Space Invaders environment, we observed  
411 that independent algorithms were able to learn well with just a team reward. Future work could be done  
412 to determine if this was only the case for fully observable environments, or under what conditions would  
413 independent algorithms still be able to cope with the multi-agent credit assignment problem. It would also  
414 be interesting to study the performance of non-recurrent variants of multi-agent algorithms such as QMIX  
415 and COMA in fully observable environments. Since the experiments performed in this paper only included  
416 fully-observable competitive and mixed environments, future work can also include a more diverse set of  
417 environments, such as partially observable competitive and mixed environments.

## CONFLICT OF INTEREST STATEMENT

418 The authors declare that the research was conducted in the absence of any commercial or financial  
419 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

420 KML implemented the algorithms, ran the experiments, analyzed the results and wrote the paper. SS  
421 assisted in designing the experiments and writing the paper, while MC provided access to computing  
422 resources for running the experiments. Both SS and MC also provided numerous advice and feedback  
423 during the entire process, and assisted in polishing the paper.

## FUNDING

424 This project was made possible through ~~a~~ grants from the Canada Wildfire Strategic Network and ~~the~~  
425 Discovery Grant Program ~~of~~ from the National Research Council of Canada (NSERC).

## ACKNOWLEDGMENTS

426 The authors would like to thank Compute Canada for providing computing resources for the experiments.

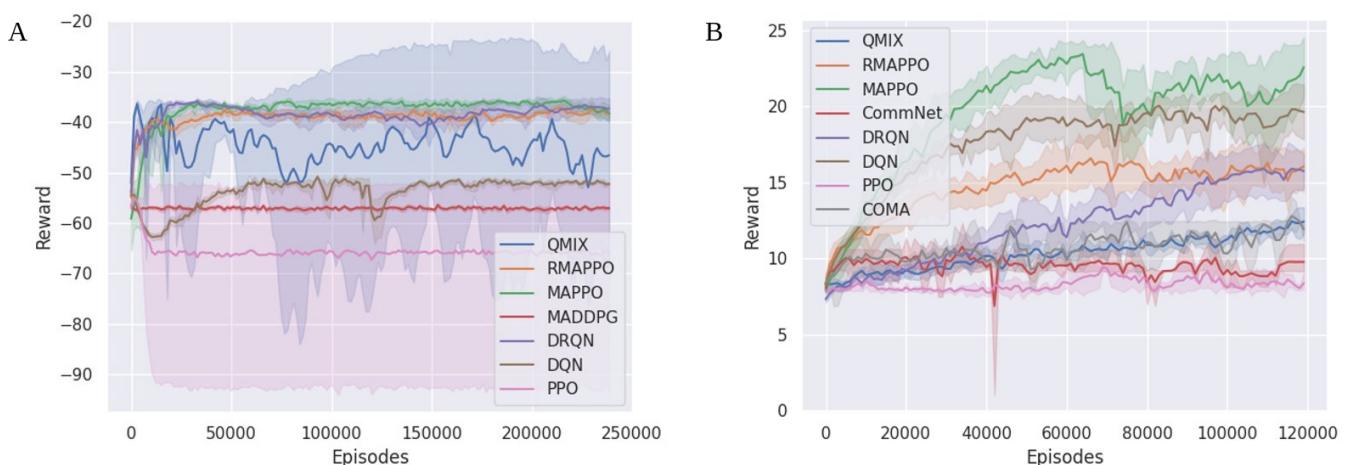
## REFERENCES

- 427 Andrychowicz, M., Raichuk, A., Stańczyk, P., Orsini, M., Girgin, S., Marinier, R., et al. (2020).  
428 What matters in on-policy reinforcement learning? a large-scale empirical study. *arXiv preprint*  
429 *arXiv:2006.05990*
- 430 Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment: An  
431 evaluation platform for general agents. *Journal of Artificial Intelligence Research* 47, 253–279
- 432 Bellman, R. (1957). A markovian decision process. *Journal of Mathematics and Mechanics* 6, 679–684
- 433 Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., et al. (2019). Dota 2 with large  
434 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*
- 435 Busoniu, L., Babuska, R., and De Schutter, B. (2008). A comprehensive survey of multiagent reinforcement  
436 learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 38,  
437 156–172. doi:10.1109/TSMCC.2007.913919
- 438 Canese, L., Cardarilli, G. C., Di Nunzio, L., Fazzolari, R., Giardino, D., Re, M., et al. (2021). Multi-agent  
439 reinforcement learning: A review of challenges and applications. *Applied Sciences* 11. doi:10.3390/  
440 app11114948
- 441 Chang, Y.-h., Ho, T., and Kaelbling, L. (2004). All learning is local: Multi-agent learning in global reward  
442 games. In *Advances in Neural Information Processing Systems*, eds. S. Thrun, L. Saul, and B. Schölkopf  
443 (MIT Press), vol. 16
- 444 Choi, S., Yeung, D.-Y., and Zhang, N. (2000). An environment model for nonstationary reinforcement  
445 learning. In *Advances in Neural Information Processing Systems*, eds. S. Solla, T. Leen, and K. Müller  
446 (MIT Press), vol. 12
- 447 Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D., Janoos, F., Rudolph, L., et al. (2020). Implementation  
448 matters in deep policy gradients: A case study on ppo and trpo. *arXiv preprint arXiv:2005.12729*
- 449 Foerster, J., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. (2018). Counterfactual multi-agent  
450 policy gradients. *Proceedings of the AAAI Conference on Artificial Intelligence* 32
- 451 Gronauer, S. and Diepold, K. (2021). Multi-agent deep reinforcement learning: a survey. *Artificial  
452 Intelligence Review* , 1–49
- 453 Gupta, J. K., Egorov, M., and Kochenderfer, M. (2017). Cooperative multi-agent control using deep  
454 reinforcement learning. In *Autonomous Agents and Multiagent Systems*, eds. G. Sukthankar and J. A.  
455 Rodriguez-Aguilar (Cham: Springer International Publishing), 66–83
- 456 Hausknecht, M. and Stone, P. (2015). Deep recurrent q-learning for partially observable mdps. In *2015  
457 aaai fall symposium series*
- 458 He, H., Boyd-Graber, J., Kwok, K., and Daumé, H. (2016). Opponent modeling in deep reinforcement  
459 learning. In *Proceedings of the 33rd International Conference on International Conference on Machine  
460 Learning - Volume 48 (JMLR.org)*, ICML’16, 1804–1813

- 461 Hernandez-Leal, P., Kaisers, M., Baarslag, T., and de Cote, E. M. (2017). A survey of learning in multiagent  
462 environments: Dealing with non-stationarity. *arXiv preprint arXiv:1707.09183*
- 463 Hernandez-Leal, P., Kartal, B., and Taylor, M. E. (2019). A survey and critique of multiagent deep  
464 reinforcement learning. *Journal of Autonomous Agents and Multi-Agent Systems (JAAMAS)* 33, 750–797
- 465 Kraemer, L. and Banerjee, B. (2016). Multi-agent reinforcement learning as a rehearsal for decentralized  
466 planning. *Neurocomputing* 190, 82–94. doi:<https://doi.org/10.1016/j.neucom.2016.01.031>
- 467 [Dataset] Li, M. (2020). Machin. <https://github.com/iffiX/machin>
- 468 Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, P., and Mordatch, I. (2017). Multi-agent actor-critic for  
469 mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*
- 470 Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M. (2018).  
471 Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents.  
472 *Journal of Artificial Intelligence Research* 61, 523–562
- 473 Markov, A. A. (1954). The theory of algorithms. *Trudy Matematicheskogo Instituta Imeni VA Steklova* 42,  
474 3–375
- 475 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., et al. (2015). Human-level  
476 control through deep reinforcement learning. *nature* 518, 529–533
- 477 Mordatch, I. and Abbeel, P. (2017). Emergence of grounded compositional language in multi-agent  
478 populations. *arXiv preprint arXiv:1703.04908*
- 479 Oliehoek, F. A., Spaan, M. T., and Vlassis, N. (2008). Optimal and approximate q-value functions for  
480 decentralized pomdps. *Journal of Artificial Intelligence Research* 32, 289–353
- 481 OpenAI, Andrychowicz, M., Baker, B., Chociej, M., Józefowicz, R., McGrew, B., et al. (2018). Learning  
482 dexterous in-hand manipulation. *CoRR*
- 483 Papoudakis, G., Christianos, F., Rahman, A., and Albrecht, S. V. (2019). Dealing with non-stationarity in  
484 multi-agent deep reinforcement learning. *arXiv preprint arXiv:1906.04737*
- 485 [Dataset] Raffin, A., Hill, A., Ernestus, M., Gleave, A., Kanervisto, A., and Dormann, N. (2019). Stable  
486 baselines3. <https://github.com/DLR-RM/stable-baselines3>
- 487 Rashid, T., Samvelyan, M., Schroeder, C., Farquhar, G., Foerster, J., and Whiteson, S. (2018). QMIX:  
488 Monotonic value function factorisation for deep multi-agent reinforcement learning. In *Proceedings of  
489 the 35th International Conference on Machine Learning*, eds. J. Dy and A. Krause (PMLR), vol. 80 of  
490 *Proceedings of Machine Learning Research*, 4295–4304
- 491 Schaul, T., Quan, J., Antonoglou, I., and Silver, D. (2015). Prioritized experience replay. *arXiv preprint  
492 arXiv:1511.05952*
- 493 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization  
494 algorithms. *arXiv preprint arXiv:1707.06347*
- 495 Shapley, L. S. (1953). Stochastic games. *Proceedings of the National Academy of Sciences* 39, 1095–1100.  
496 doi:[10.1073/pnas.39.10.1095](https://doi.org/10.1073/pnas.39.10.1095)
- 497 Shoham, Y. and Leyton-Brown, K. (2008). *Multiagent systems: Algorithmic, game-theoretic, and logical  
498 foundations* (Cambridge University Press)
- 499 [Dataset] starry\_sky6688 (2019). Starcraft. [https://github.com/starry-sky6688/  
500 StarCraft](https://github.com/starry-sky6688/StarCraft)
- 501 Sukhbaatar, S., Szlam, A., and Fergus, R. (2016). Learning multiagent communication with  
502 backpropagation. In *Proceedings of the 30th International Conference on Neural Information Processing  
503 Systems* (Red Hook, NY, USA: Curran Associates Inc.), NIPS’16, 2252–2260
- 504 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction* (MIT press)

- 505 Tampuu, A., Matiisen, T., Kodelja, D., Kuzovkin, I., Korjus, K., Aru, J., et al. (2017). Multiagent  
 506 cooperation and competition with deep reinforcement learning. *PLOS ONE* 12, 1–15. doi:10.1371/journal.pone.0172395
- 508 Tan, M. (1993). Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings*  
 509 *of the tenth international conference on machine learning*. 330–337
- 510 Terry, J. K. and Black, B. (2020). Multiplayer support for the arcade learning environment. *arXiv preprint*  
 511 *arXiv:2009.09341*
- 512 Terry, J. K., Black, B., and Hari, A. (2020a). Supersuit: Simple microwrappers for reinforcement learning  
 513 environments. *arXiv preprint arXiv:2008.08932*
- 514 Terry, J. K., Black, B., Jayakumar, M., Hari, A., Sullivan, R., Santos, L., et al. (2020b). Pettingzoo: Gym  
 515 for multi-agent reinforcement learning. *arXiv preprint arXiv:2009.14471*
- 516 Van Hasselt, H., Guez, A., and Silver, D. (2016). Deep reinforcement learning with double q-learning. In  
 517 *Proceedings of the AAAI conference on artificial intelligence*. vol. 30
- 518 Wang, Z., Schaul, T., Hessel, M., Hasselt, H., Lanctot, M., and Freitas, N. (2016). Dueling network  
 519 architectures for deep reinforcement learning. In *International conference on machine learning (PMLR)*,  
 520 1995–2003
- 521 Yu, C., Velu, A., Vinitsky, E., Wang, Y., Bayen, A., and Wu, Y. (2021). The surprising effectiveness of  
 522 mappo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*
- 523 Zawadzki, E., Lipson, A., and Leyton-Brown, K. (2014). Empirically evaluating multiagent learning  
 524 algorithms. *arXiv preprint arXiv:1401.8074*
- 525 Zhang, K., Yang, Z., and Başar, T. (2021). Multi-agent reinforcement learning: A selective overview of  
 526 theories and algorithms. *Handbook of Reinforcement Learning and Control* , 321–384

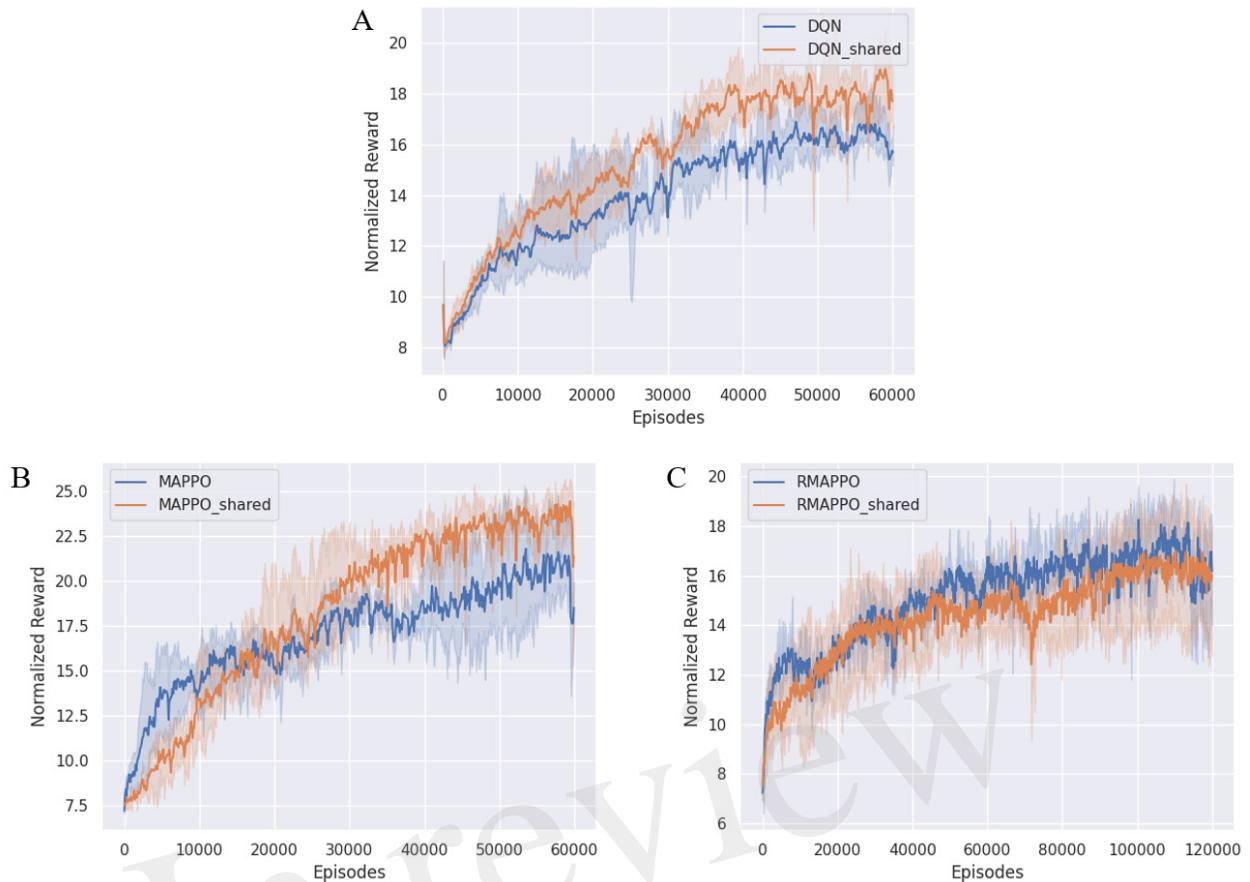
## FIGURE CAPTIONS



**Figure 1.** Training curves of various algorithms in two cooperative environments. For every algorithm, the solid line represents the mean reward per episode, while the shaded region represents the 95% confidence interval around the mean. **(A)** shows training curve for Simple Reference environment, **(B)** shows training curve for Space Invaders environment.

**Table 1.** Final scores (mean and standard deviation) of algorithms obtained over the last 100 episodes across all 5 seeds in the Space Invaders and Simple Reference environments. The highest score in both environments is bolded.

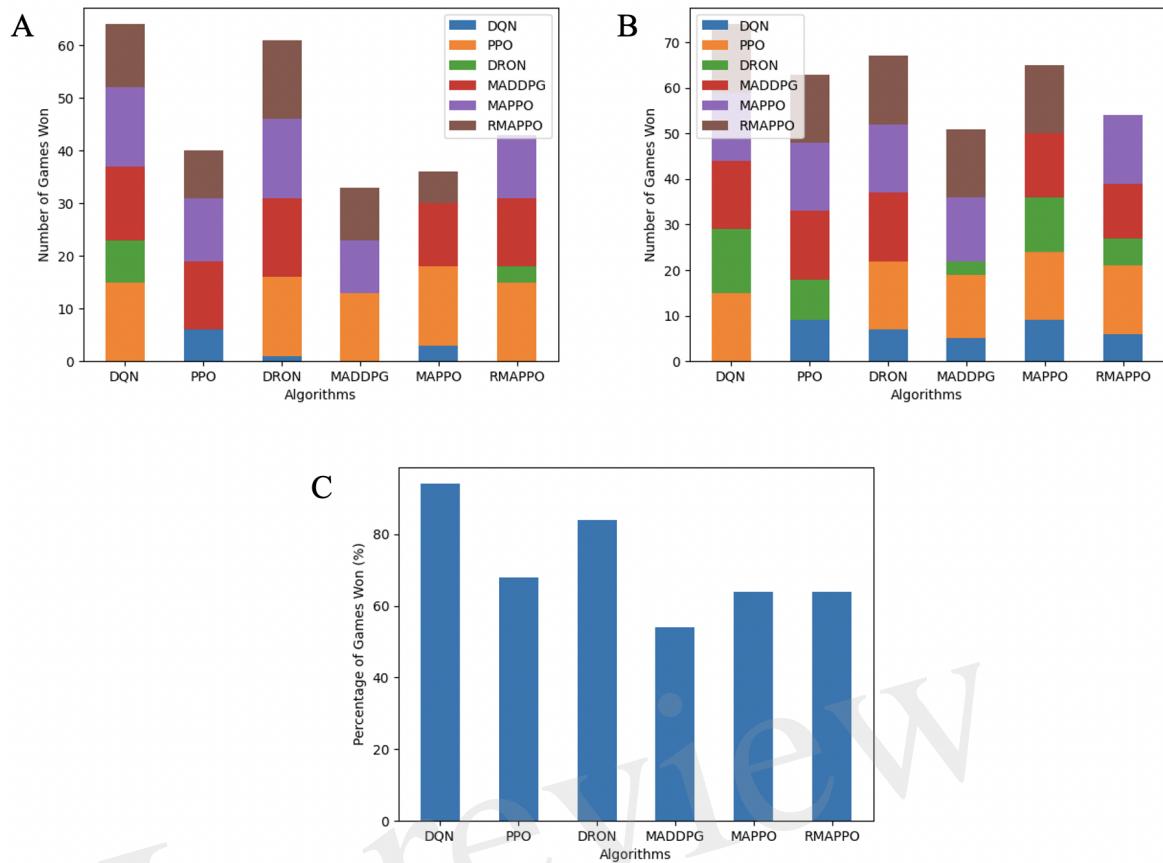
Algorithm	Space Invaders	Algorithm	Simple Reference
QMIX	$12.5 \pm 3.28$	QMIX	$-46.2 \pm 28.4$
RMAPPO	$16.2 \pm 3.31$	RMAPPO	<b><math>-36.8 \pm 11.8</math></b>
MAPPO	<b><math>22.5 \pm 3.45</math></b>	MAPPO	$-38.0 \pm 14.0$
CommNet	$9.78 \pm 0.98$	MADDPG	$-58.2 \pm 18.5$
COMA	$12.2 \pm 1.61$	DRQN	$-36.6 \pm 13.3$
DRQN	$15.7 \pm 3.70$	DQN	$-52.8 \pm 19.1$
DQN	$19.8 \pm 3.52$	PPO	$-66.0 \pm 33.0$
PPO	$7.92 \pm 2.69$		



**Figure 2.** Training curves of various algorithms in Space Invaders, comparing when individual rewards are given (blue) to when team rewards are given (orange). (A) shows training curve of DQN, (B) shows training curve of MAPPO, (C) shows training curve of RMAPPO.

**Table 2.** Overall winrate percentage of various algorithms across Boxing, Pong and Space War environments. The highest winrate percentage for every environment is bolded.

Algorithm	Boxing	Pong	Space War
DQN	<b>94%</b>	88%	70%
PPO	<b>68%</b>	36%	<b>52%</b>
DRON	<b>84%</b>	<b>90%</b>	68%
MADDPG	54%	28%	72%
MAPPO	<b>64%</b>	40%	66%
RMAPPO	<b>64%</b>	24%	<b>76%</b>



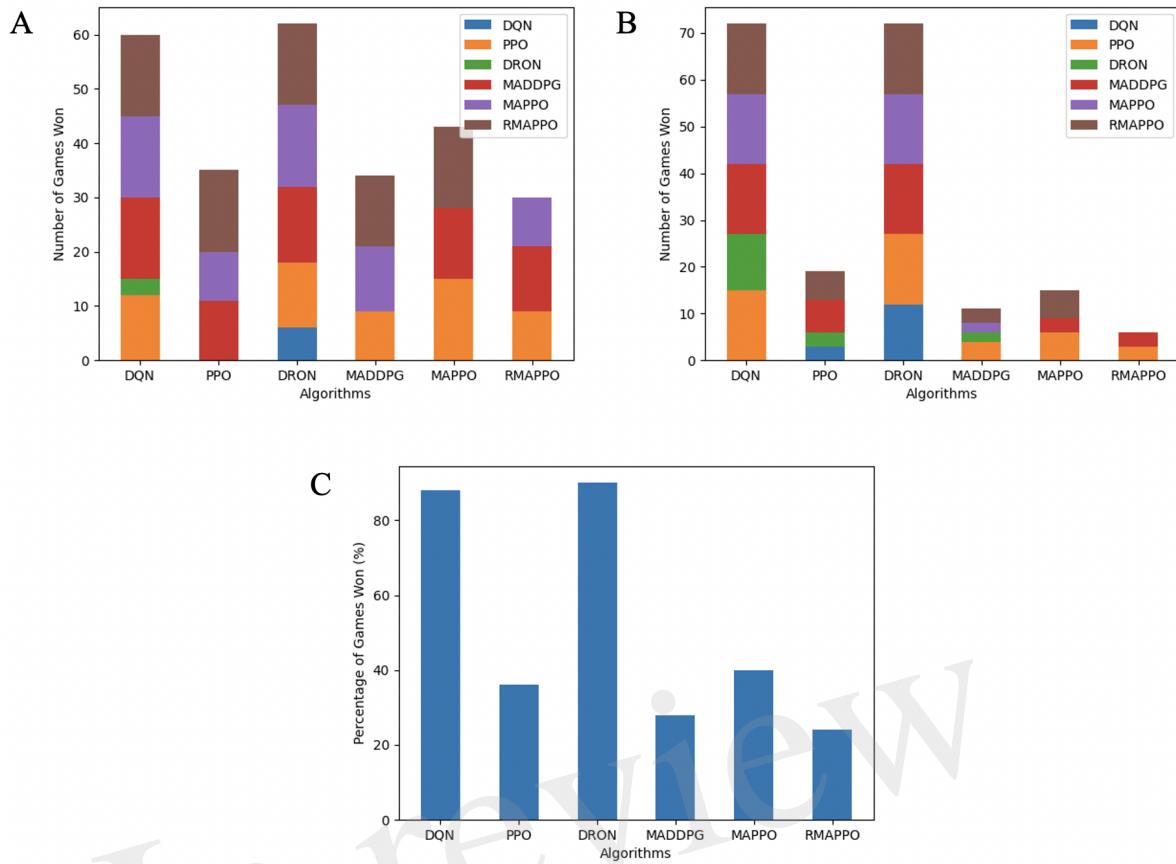
**Figure 3.** Performance of various algorithms when playing against other algorithms in the Boxing environment. (A) shows the number of games won as the first player, (B) shows the number of games won as the second player, (C) shows the overall win rate percentage.

**Table 3.** Final scores (mean and standard deviation) of algorithms obtained over the last 100 episodes across all 5 seeds in the Simple Tag and Simple Adversary environments.

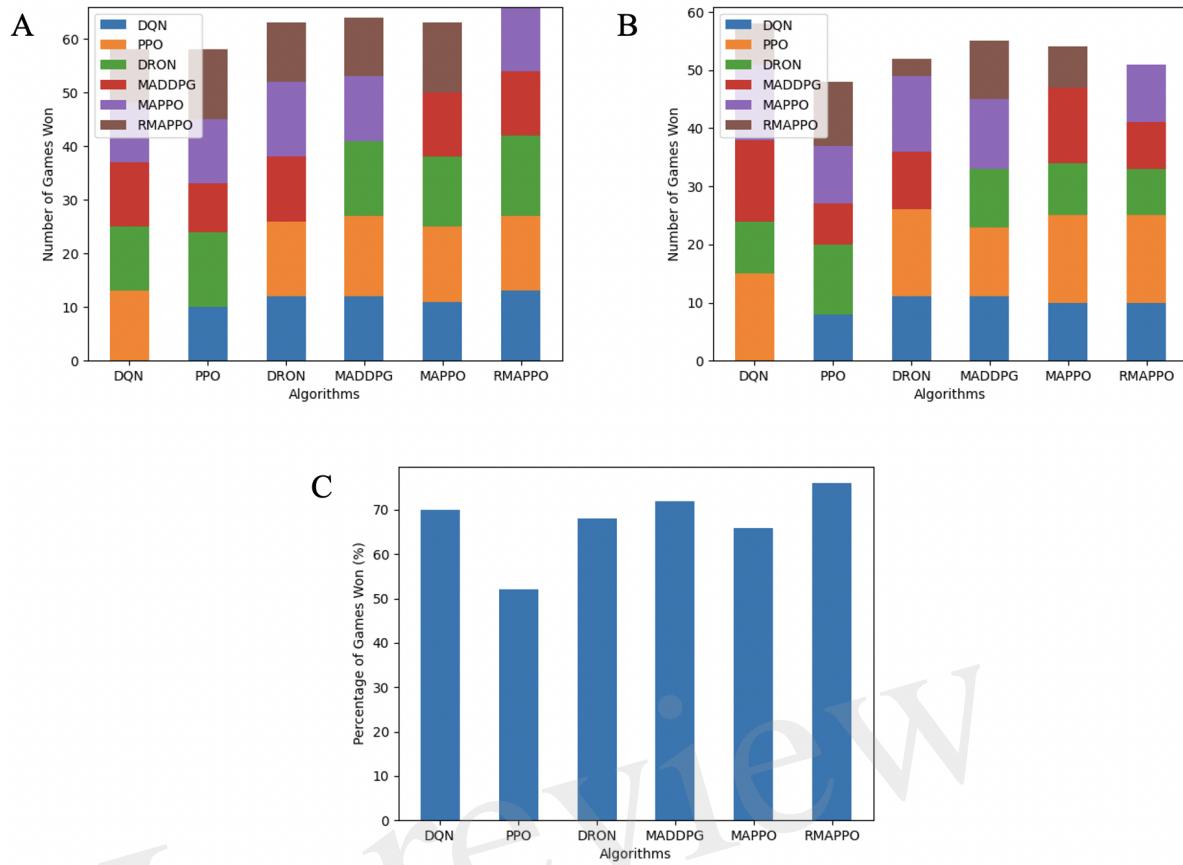
Algorithm	Simple Tag (Predator)	Algorithm	Simple Tag (Prey)
DQN	$3.24 \pm 7.99$	DQN	$-4.93 \pm 9.36$
PPO	$2.46 \pm 7.60$	PPO	$-47.0 \pm 57.9$
MADDPG	$4.40 \pm 9.24$	MADDPG	$-8.10 \pm 12.2$
RMAPPO	$14.0 \pm 20.5$	RMAPPO	$-16.4 \pm 21.2$
MAPPO	$13.4 \pm 19.1$	MAPPO	$-20.8 \pm 22.5$

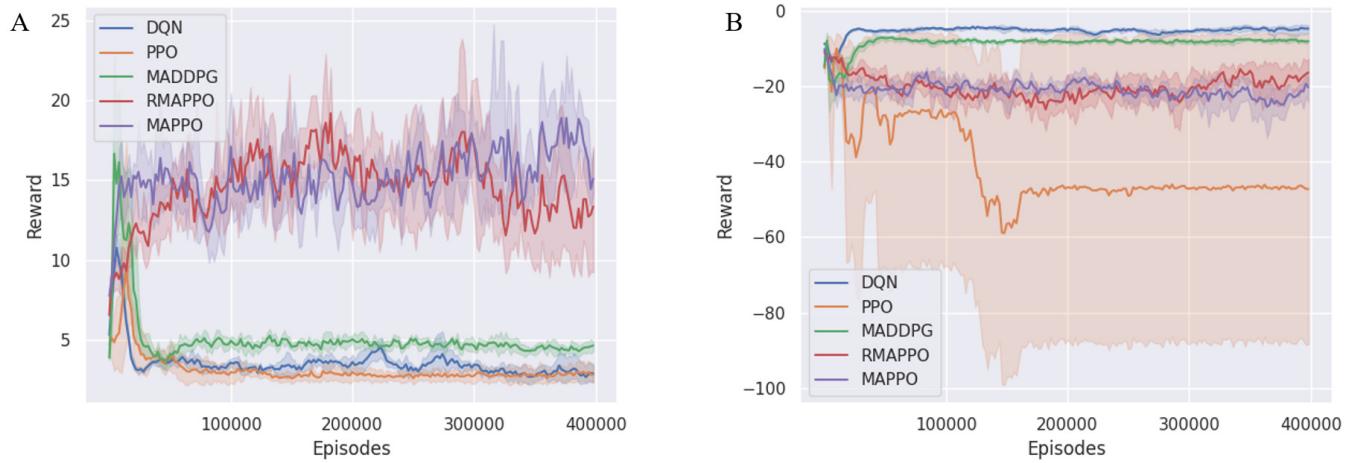
Algorithm	Simple Adversary (Adversary)	Algorithm	Simple Adversary (Cooperative Agent)
DQN	$-18.2 \pm 9.27$	DQN	$7.56 \pm 7.85$
PPO	$-52.7 \pm 21.7$	PPO	$34.2 \pm 20.0$
MADDPG	$-15.6 \pm 7.65$	MADDPG	$7.24 \pm 6.35$
RMAPPO	$-29.4 \pm 15.0$	RMAPPO	$10.8 \pm 16.4$
MAPPO	$-24.9 \pm 12.2$	MAPPO	$9.11 \pm 13.8$



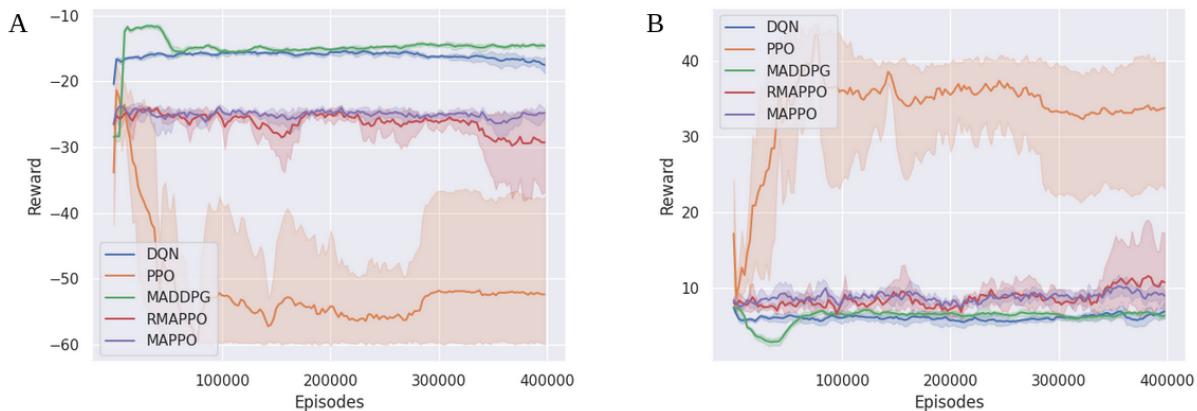
**Figure 4.** Performance of various algorithms when playing against other algorithms in Pong. **(A)** shows the number of games won as the first player, **(B)** shows the number of games won as the second player, **(C)** shows the overall win rate percentage.



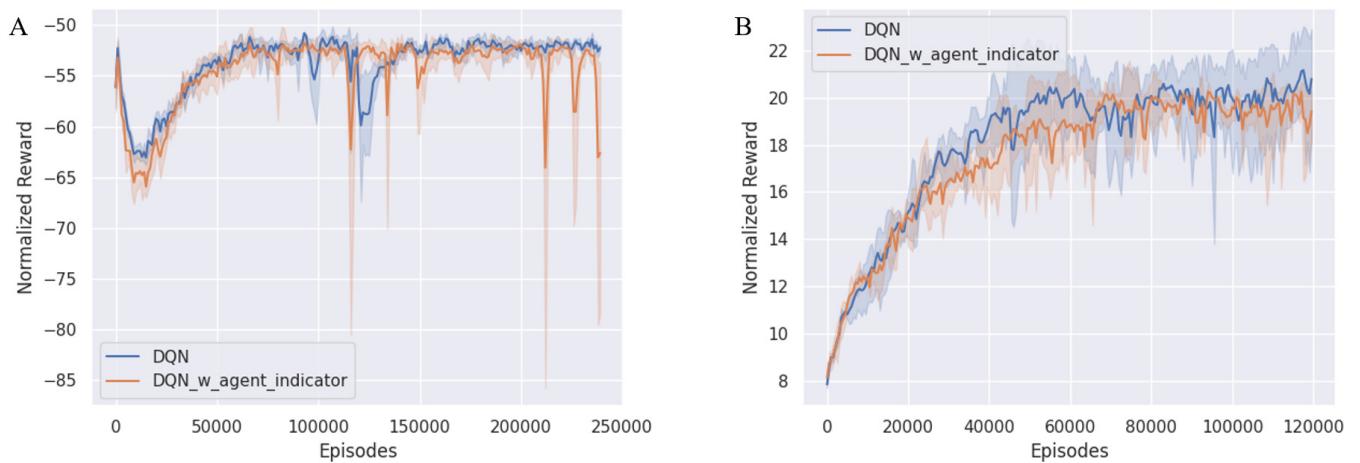
**Figure 5.** Performance of various algorithms when playing against other algorithms in the Space War environment. **(A)** shows the number of games won as the first player, **(B)** shows the number of games won as the second player, **(C)** shows the overall win rate percentage.



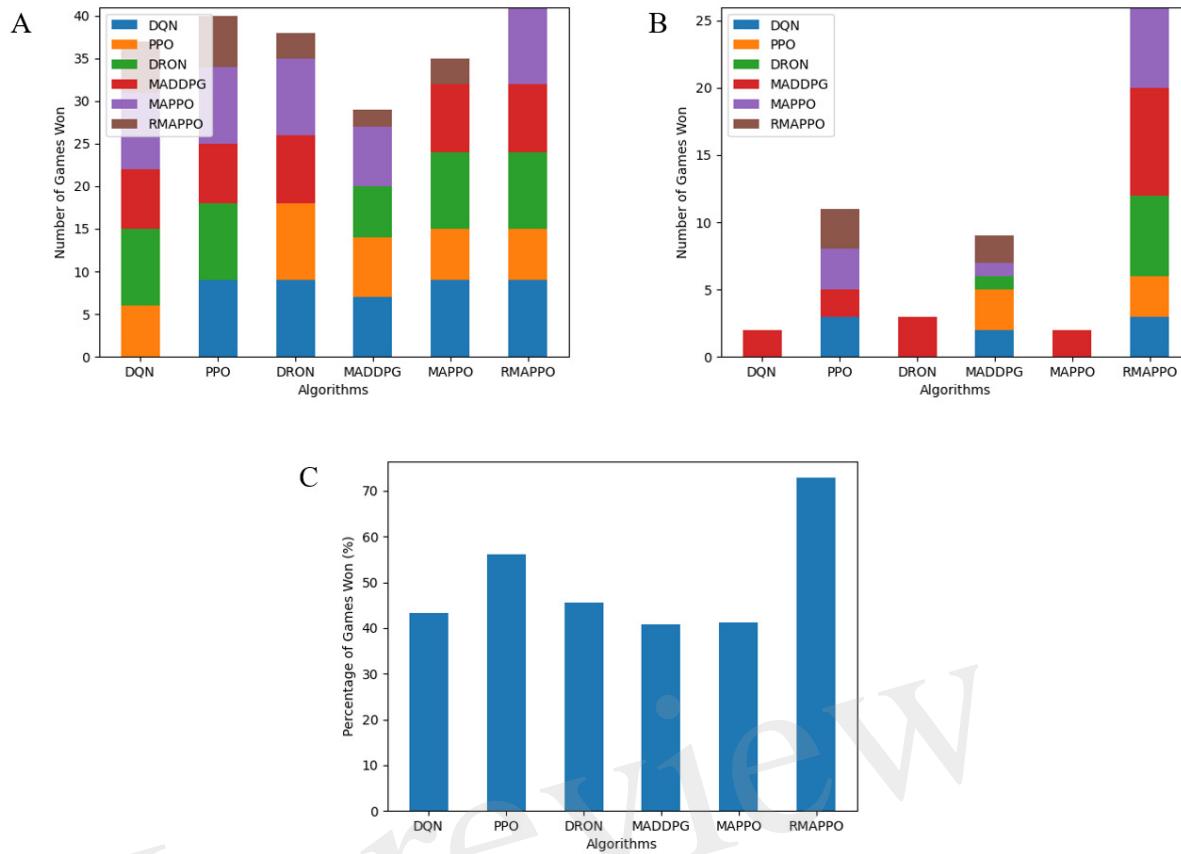
**Figure 6.** Training curves of various algorithms in the Simple Tag, a Predator-Prey environment. **(A)** shows the reward of a predator (all predators obtain the same reward), **(B)** shows the reward of the prey.



**Figure 7.** Training curves of various algorithms in the Simple Adversary environment. (A) shows the reward of the adversary, (B) shows the reward of a cooperative agent (both cooperative agents obtain the same reward).



**Figure 8.** Comparing DQN with (blue) and without (orange) agent indicators in (A) Simple Reference and (B) Space Invaders environment.



**Figure 9.** Performance of various algorithms when playing against other algorithms in Pong without agent indicators across 3 seeds. **(A)** shows the number of games won as the first player, **(B)** shows the number of games won as the second player, **(C)** shows the overall win rate percentage.

Figure 1.JPG

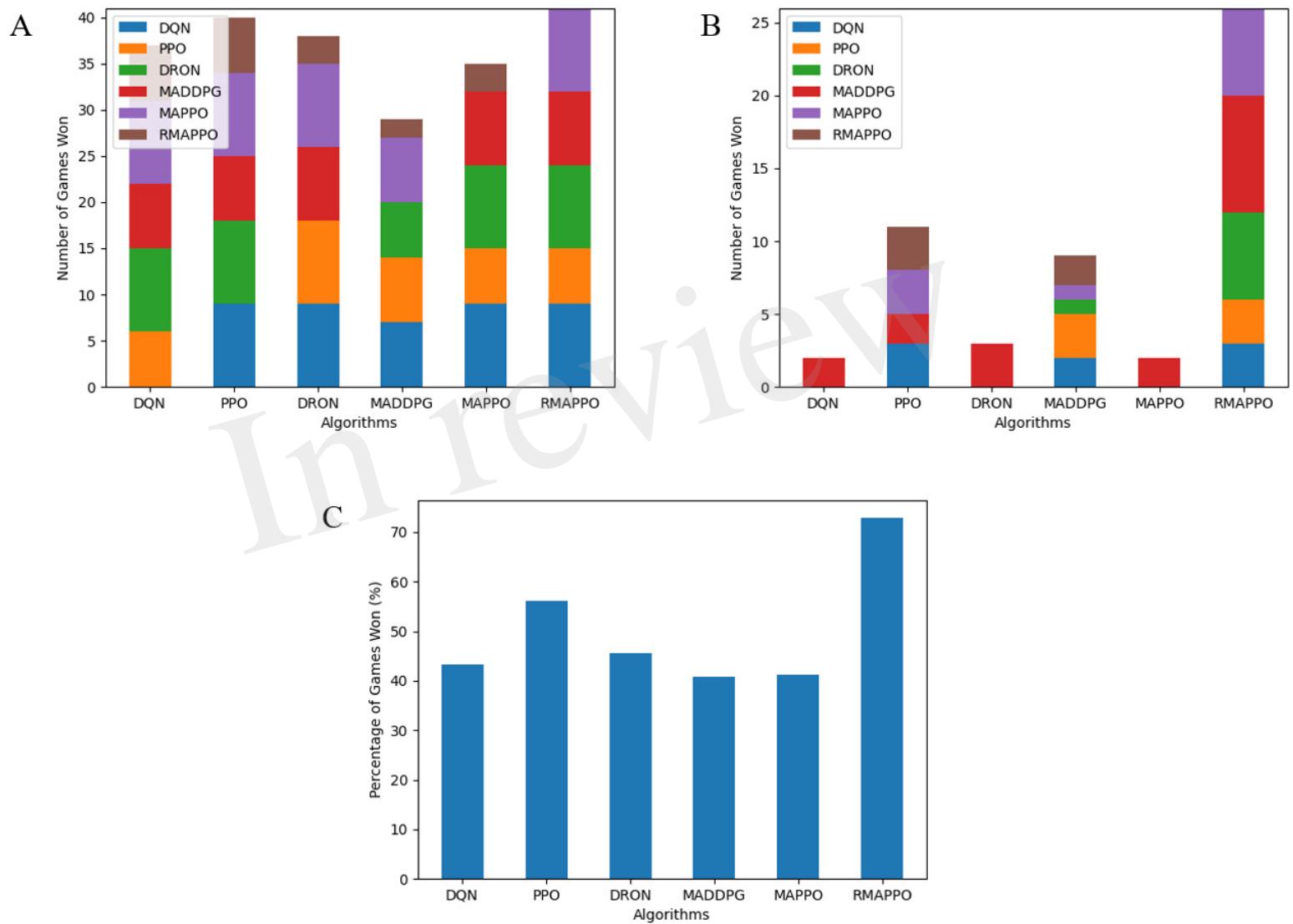


Figure 2.JPG

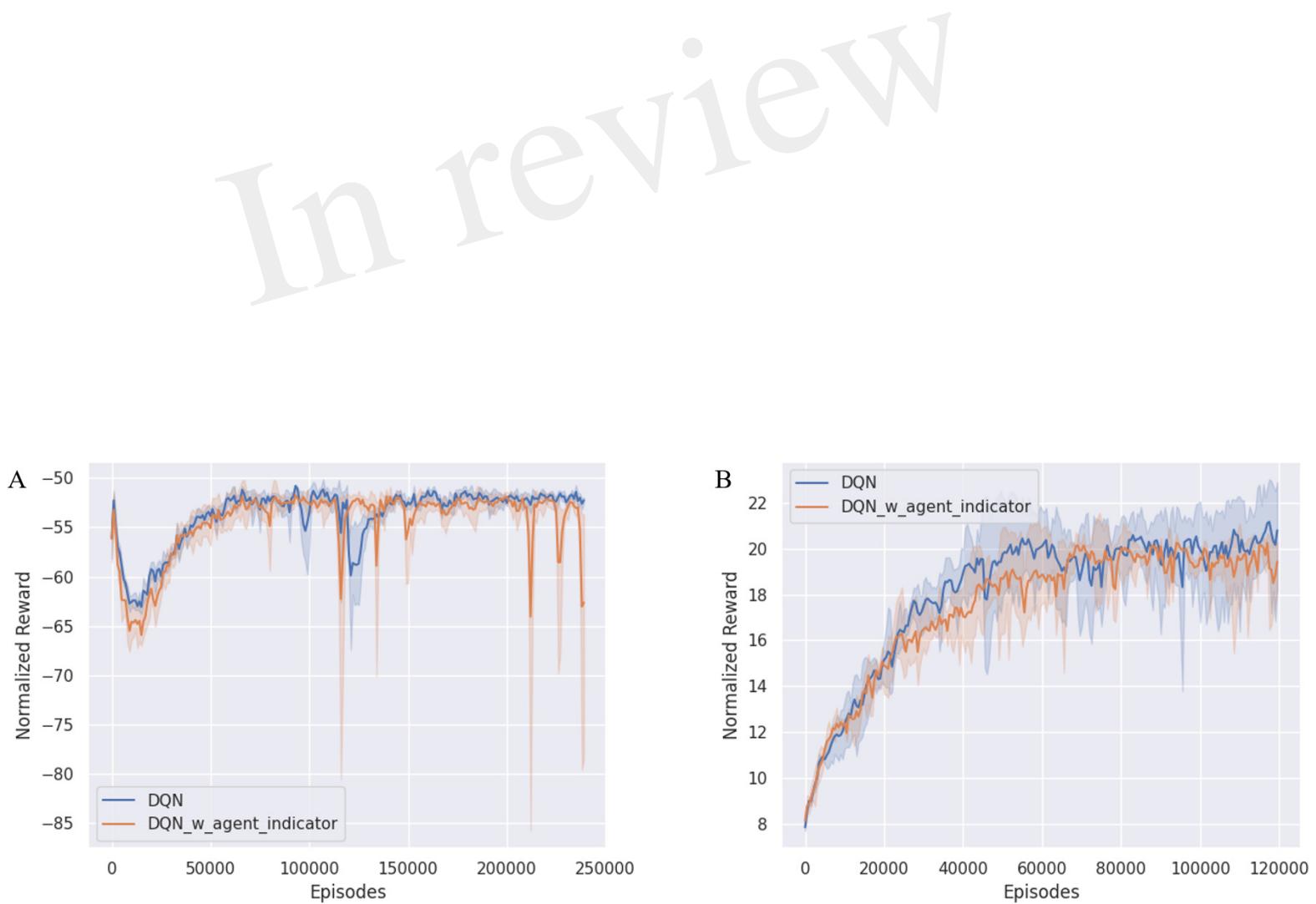


Figure 3.JPG

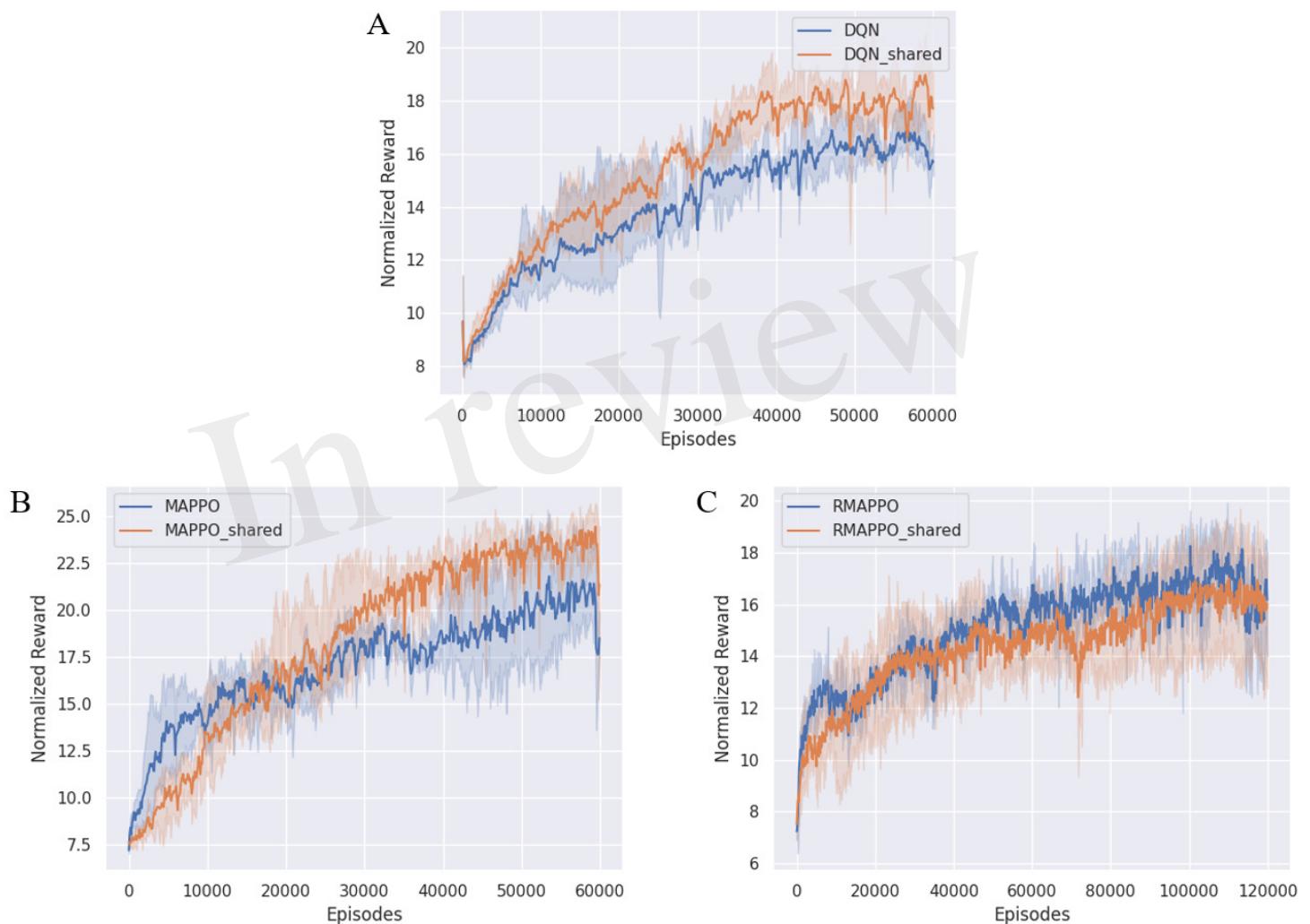


Figure 4.JPG

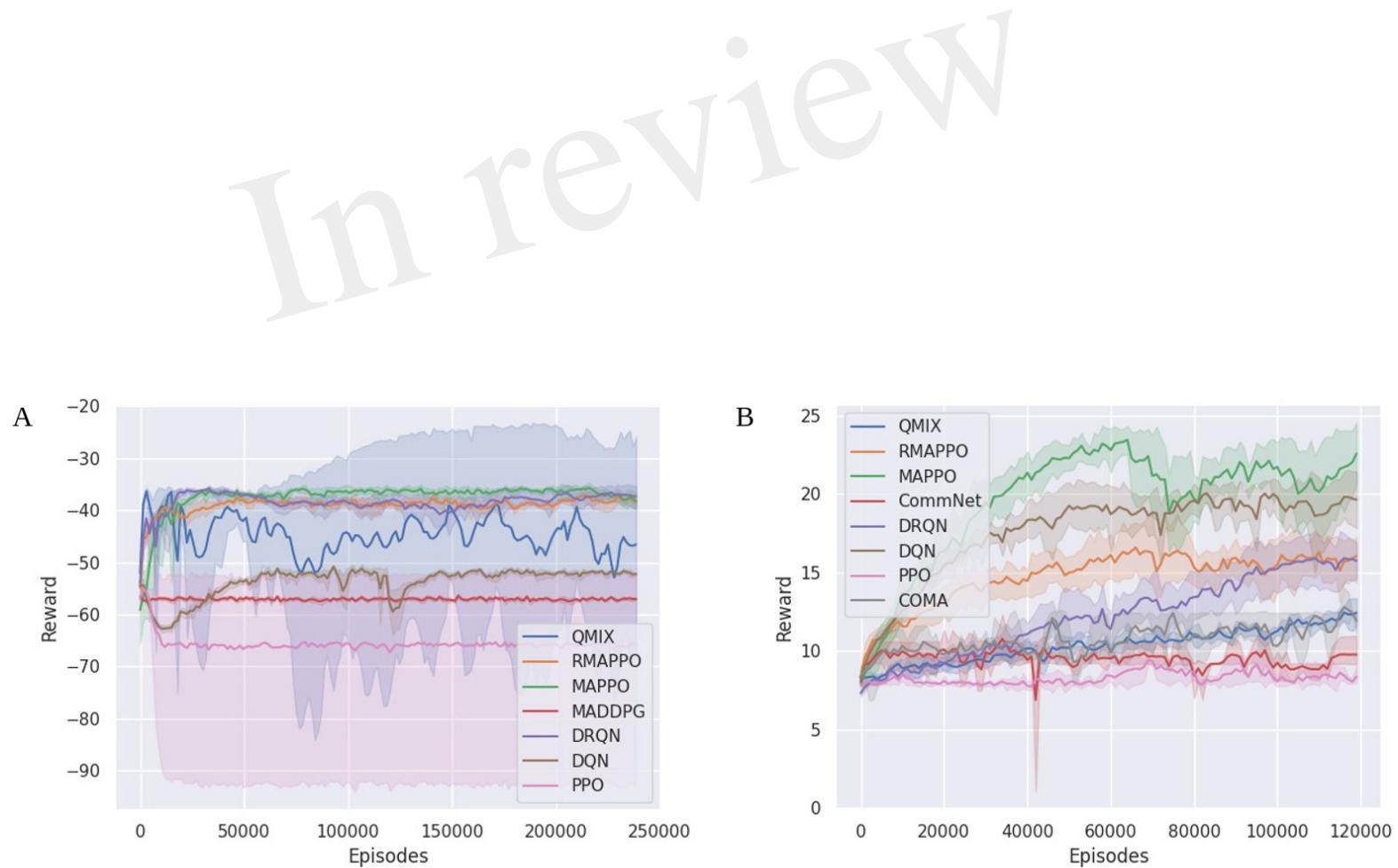


Figure 5.JPG

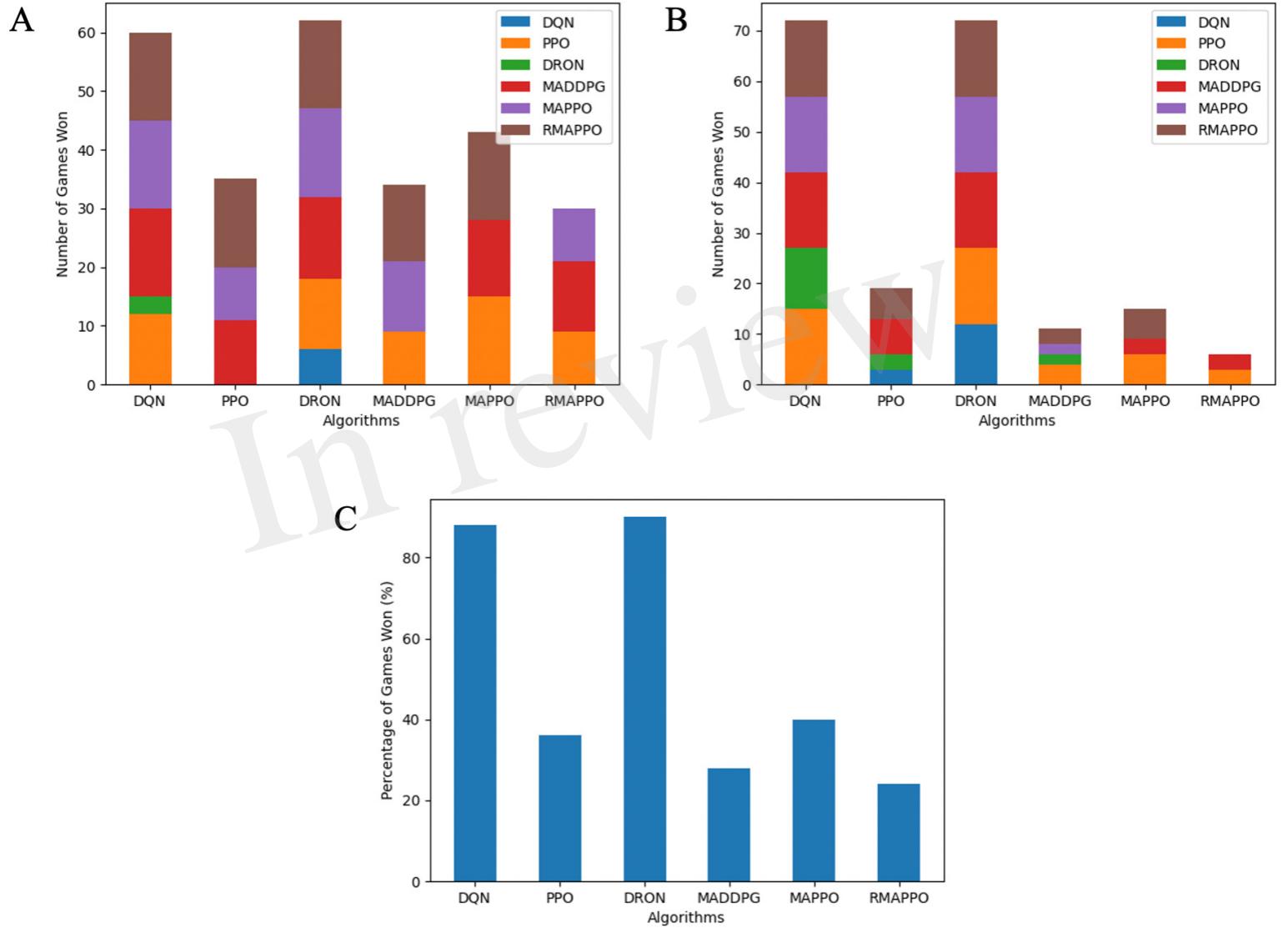


Figure 6.JPG

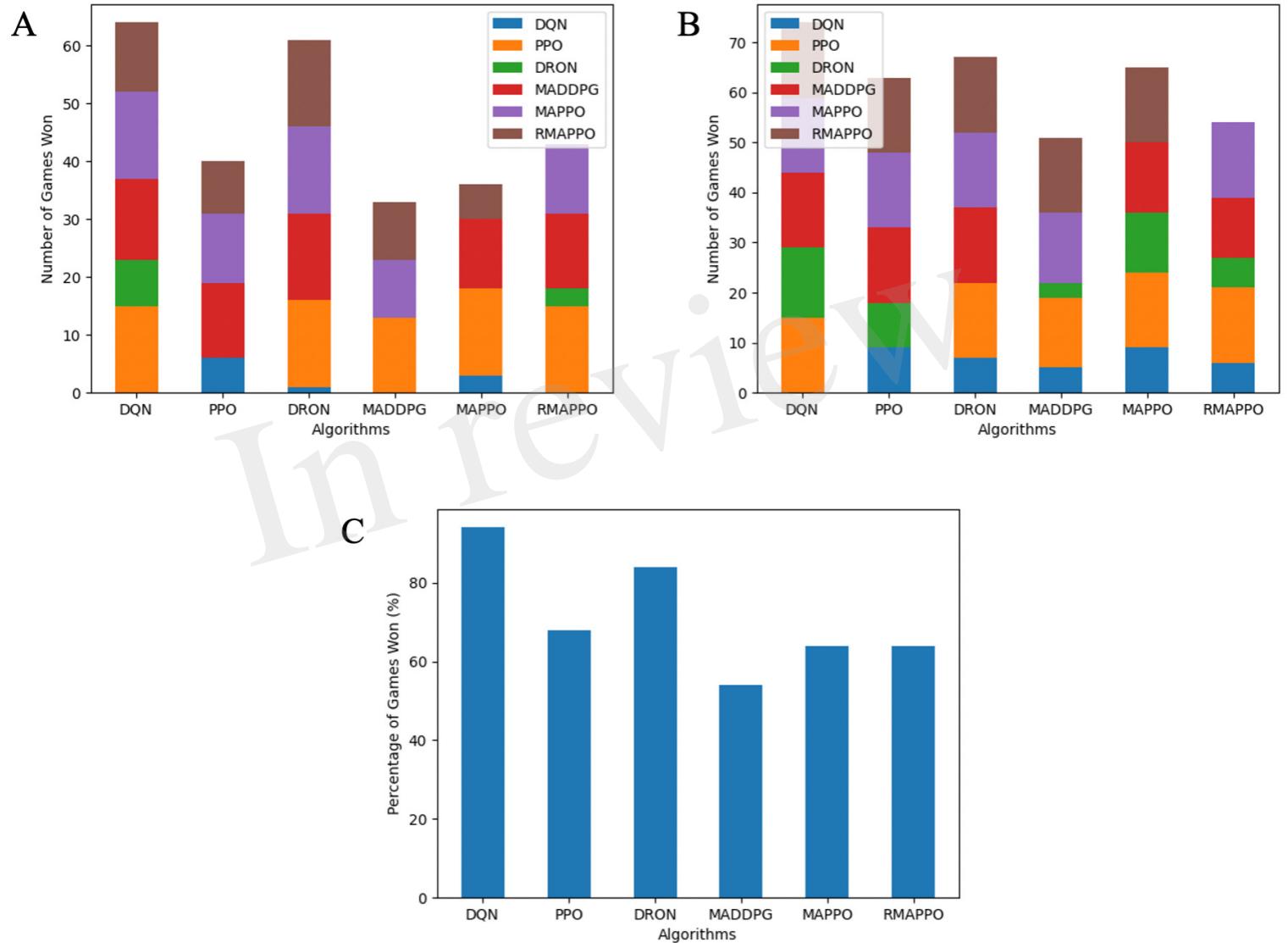


Figure 7.JPG

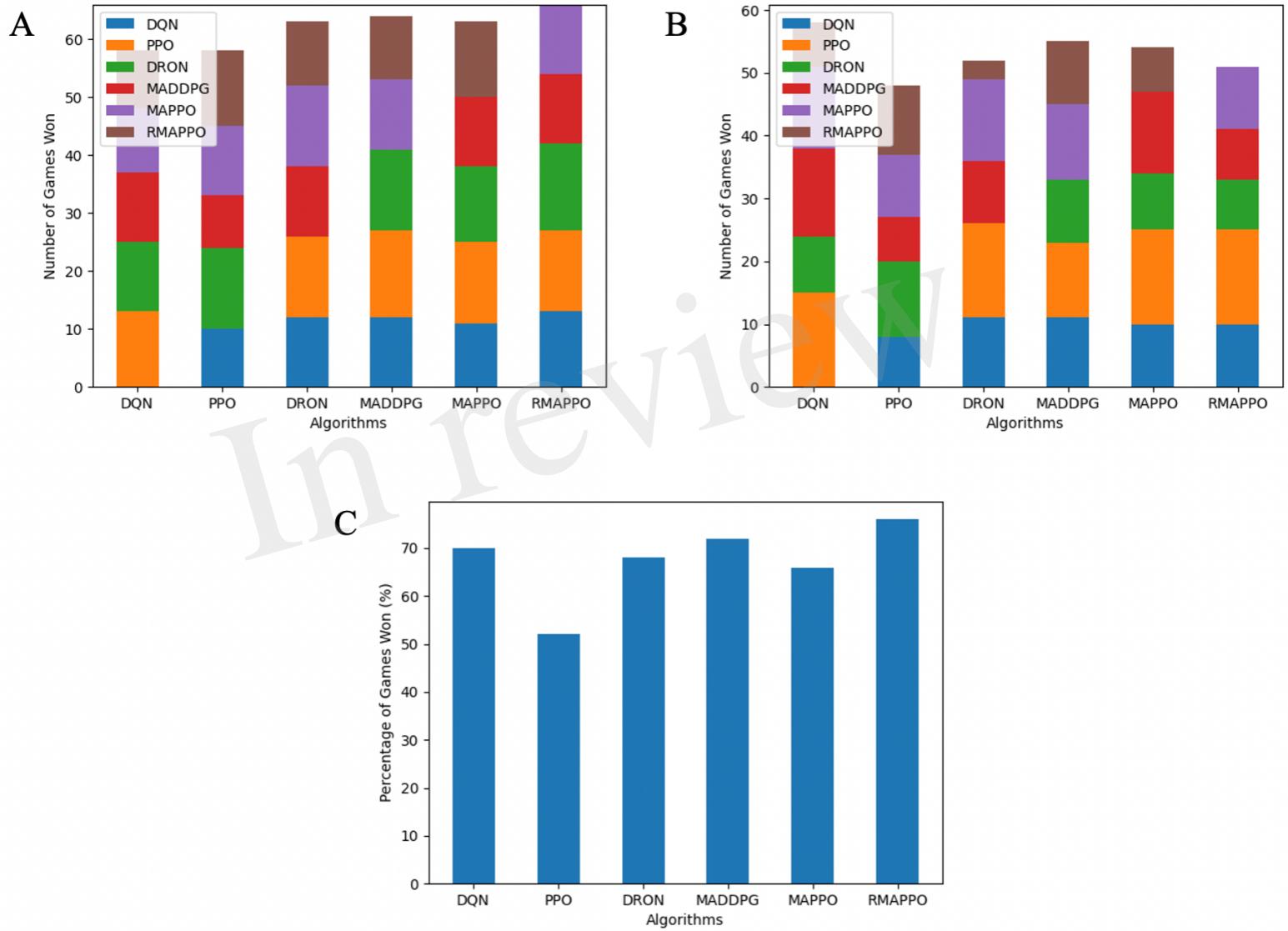


Figure 8.JPG

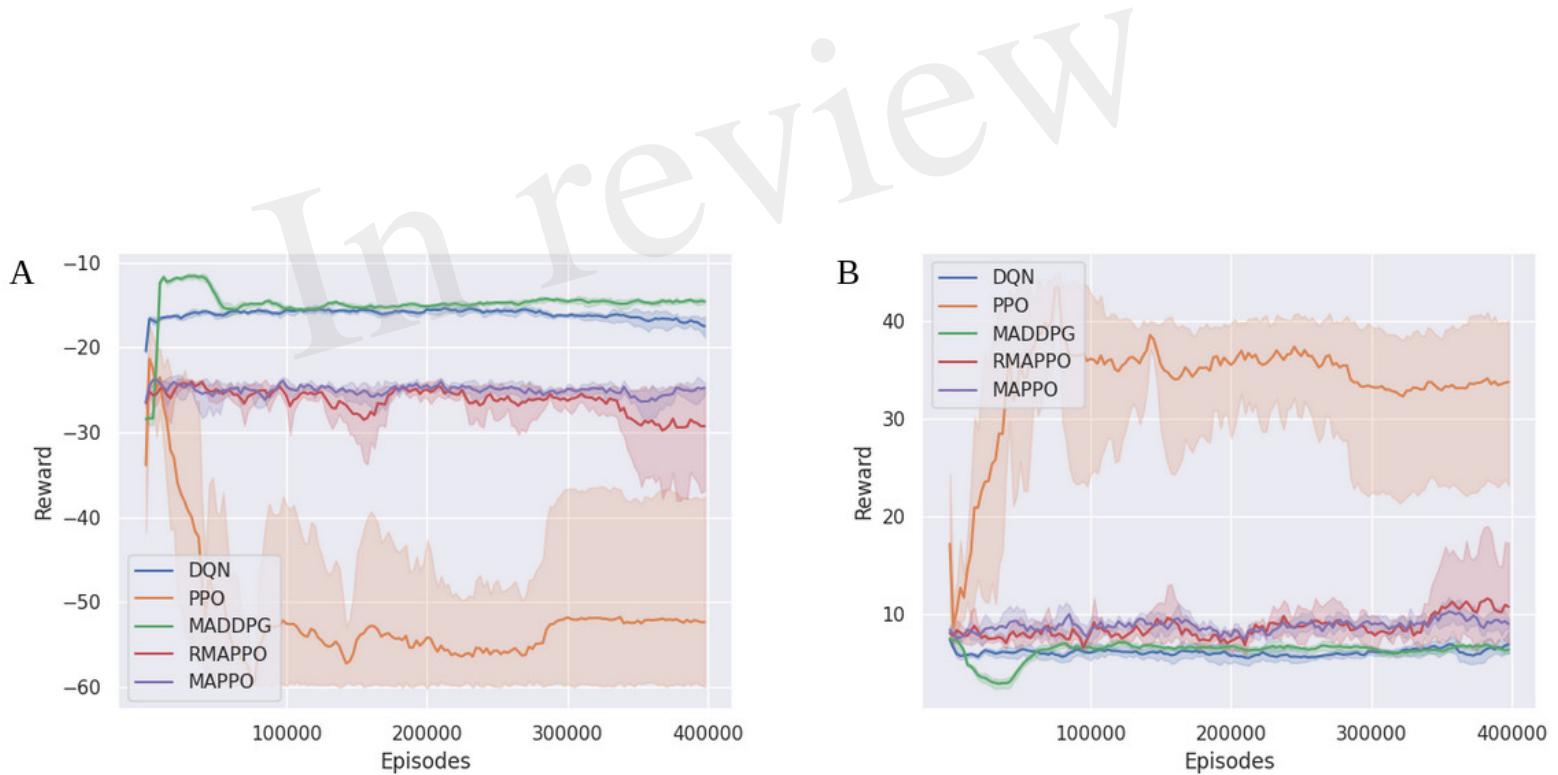


Figure 9.JPG

