

Investigation of Independent Reinforcement Learning Algorithms in Multi-Agent Environments

Ken Ming Lee^{1*}, Sriram Ganapathi Subramanian¹, Mark Crowley¹

¹University 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:
28 Dec 2021

Journal website link:
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: 152

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 PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. We show that in fully-observable environments, independent algorithms can perform on par with multi-agent algorithms in cooperative and competitive settings. For the mixed environments, we 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 environments.

Contribution to the field

Independent (i.e., single-agent) 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 PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. We show that in fully-observable environments, independent algorithms can perform on par with multi-agent algorithms in cooperative and competitive settings. For the mixed environments, we 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 environments.

Funding statement

This project was made possible through a grants from the Canada Wildfire Strategic Network and Discovery Grant Program of from 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.

In review

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^{1,*}, Sriram Ganapathi Subramanian¹ and Mark Crowley¹

¹*Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON, Canada*

Correspondence*:

Ken Ming Lee
km23lee@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 PettingZoo environments that span the three main categories of multi-agent environments, i.e., cooperative, competitive, and mixed. We show that in fully-observable environments, independent algorithms can perform on par with multi-agent algorithms in cooperative and competitive settings. For the mixed environments, we 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 environments.

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

16 [Number of words: 44844494, Number of figures and tables: 68]

1 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 lower computational resources and being easier to scale to large environments than traditional multi-agent algorithms which perform exact opponent modelling of each agent. In practice, prior works have reported mixed performance for independent algorithms in different multi-agent domains (Zawadzki et al., 2014; Tampuu et al., 2017; Shoham and Leyton-Brown, 2008; Berner et al., 2019; Foerster et al.,

(Berner et al., 2019; Foerster et al., 2018; Lowe et al., 2017; Rashid et al., 2018; Shoham and Leyton-Brown, 2008). However, a study of the strengths and weaknesses of independent algorithms across various categories within the multi-agent domain is lacking in the literature.

In this paper, we investigate the empirical performance of independent algorithms in multi-agent settings, and compare them to various multi-agent algorithms under the Centralized Training and Decentralized Execution scheme (Kraemer and Banerjee, 2016; Oliehoek et al., 2008). We evaluate these algorithms on 4 multi-agent environments from the PettingZoo library (Terry et al., 2020b), which span the 3 main categories of multi-agent environments (i.e., cooperative, competitive, and mixed) (Busoniu et al., 2008; Canese et al., 2021; Zhang et al., 2021; Gronauer and Diepold, 2021) (Busoniu et al., 2008; Canese et al., 2021; Gronauer and Diepold, 2021; Zhang et al., 2021). We show that independent algorithms can perform on par with multi-agent algorithms in the cooperative, fully-observable setting, and adding recurrence allows them to perform well compared to multi-agent algorithms in partially observable environments. In the competitive setting, we show that parameter sharing alongside the addition of agent indicators allow independent algorithms to outperform some multi-agent algorithms, such as Multi-Agent Proximal Policy Optimization (Yu et al., 2021), and Multi-Agent Deep Deterministic Policy Gradient (Lowe et al., 2017), in fully-observable environments. For the mixed setting, we show that agents of independent algorithms learn to perform well individually, but fail in learning to cooperate with allies and compete against enemies.

2 BACKGROUND INFORMATION

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

2.1 Reinforcement Learning

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

2.2 Multi-Agent Reinforcement Learning

The single-agent MDP framework is extended to the Multi-Agent Reinforcement Learning (MARL) setting in the form of stochastic games (Shapley, 1953). In an N -agent stochastic game, at every time step, 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 $\mathbf{u}_t \triangleq \{u_t^1, \dots, u_t^N\}$ determines the rewards obtained by each agent is written as $\mathbf{u}_t \triangleq \langle u_t^1, \dots, u_t^N \rangle$. Every agent receives an agent specific reward through the reward function $R(s_t, \mathbf{u}_t, j)$. State transitions of the

environment are determined by the transition probability $P(s_{t+1}|s_t, u_t)P(s_{t+1}|s_t, u_t)$, which conditions on the state and the joint action at timestep t .

2.3 Centralized Training and Decentralized Execution

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

2.4 Cooperative, Competitive and Mixed

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

In the cooperative setting, agents collaborate with each other to achieve a common goal. As a result, it is very common for all agents to share the same reward function (i.e., a team goal) (Chang et al., 2004). Also known as the multi-agent credit assignment problem, every agent has to deduce its own contributions from the team reward (Chang et al., 2004). Algorithms studied in this paper that explicitly address the multi-agent credit-assignment problem include QMIX (Rashid et al., 2018) and Counterfactual Multi-Agent Policy Gradients (COMA) (Foerster et al., 2018). Additionally, the CommNet (Sukhbaatar et al., 2016) extension on top of COMA is utilized for specific cooperative environments. Other multi-agent algorithms that are considered for the cooperative scenario include Multi-Agent Deep Deterministic Policy Gradient (MADDPG) (Lowe et al., 2017) and Multi-Agent Proximal Policy Optimization (MAPPO) (Yu et al., 2021).

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

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

2.5 Independent Algorithms and Non-Stationarity

One naive approach for applying single-agent RL to the multi-agent setting would be the use of independent learners, where each agent treats every other agent as part of the environment, and learns purely based on individual observations. In a multi-agent setting, the transition probability P and

reward function R are conditioned on the joint action u . Since all agents in the environment are learning, their policies constantly change. Therefore, from each independent learner’s perspective, the transition probability and reward function appear non-stationary, due to the lack of awareness of other agents’ actions. This violates the Markovian property of an MDP, which then causes independent algorithms to be slow to adapt to other agents’ changing policies, and as a result, face difficulties in converging to a good policy (Papoudakis et al., 2019; Hernandez-Leal et al., 2017; He et al., 2016) (He et al., 2016; Hernandez-Leal et al., 2017; Papoudakis et al., 2019).

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

3 EXPERIMENTAL SETUP

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

3.1 Environments

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

For the cooperative setting, experiments were performed on a modified version of the 2-player Space Invaders (Bellemare et al., 2013; Terry and Black, 2020), and the Simple Reference MPE environment (Lowe et al., 2017; Mordatch and Abbeel, 2017). In Space Invaders, both agents share the common goal of shooting down all aliens. To make Space Invaders cooperative, we removed the (positive) reward that is given to a player whenever the other player gets hit. Additionally, the environment rewards every agent individually by default. Since a number of cooperative multi-agent algorithms (e.g., QMIX and COMA) assume that only a team reward is given, we modified the reward function such that a team reward is given instead (i.e., both agents receive the sum of their individual rewards). This setup exemplifies the multi-agent credit assignment problem, the effect of which is studied more closely in the Section 4.1.1. On the other hand, in the Simple Reference environment, two agents are rewarded by how close they are to their target landmark. However, the target landmark of an agent is only known by the other agent, as a result communication is required for both agents to navigate successfully to their target landmarks.

For the competitive setting, we performed experiments on the 2-player variant of the original Atari Pong environment. For the mixed setting, we opted for the Simple Tag MPE environment, which is a Predator-Prey environment (Mordatch and Abbeel, 2017). This environment consists of 4 agents – 3 predators and a prey. The prey travels faster and has to avoid colliding with the predators, while the 3 predators travel slower and have to work together to capture the prey. The rewards received by the prey and a predator sum to 0 (i.e., the prey gets a negative reward for collision, while the predators get rewarded positively), and all predators receive the same reward. The prey is also negatively rewarded if it strays

away from the predefined area (a 1×1 unit square). This environment is general-sum because it contains 3 predators and a single prey.

3.2 Preprocessing

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

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

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

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

Furthermore, the action spaces for both Atari environments were shrunk to their effective action spaces in order to improve learning efficiency. For Pong specifically, we also concatenated a one-hot vector of the agent's index to the observations so that independent algorithms can differentiate one from the other when parameter sharing is utilized. The effect of this addition is studied more closely in Section 4.4.

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

3.3 Implementation

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

- Implementation of DQN and DRON were based on the Machin library (Li, 2020).
- Implementation of independent PPO was based on Stable Baselines3 (Raffin et al., 2019).
- Implementation of DRQN, QMIX, COMA and CommNet came from a popular public repository by the name of StarCraft (starry sky6688, 2019).
- Implementation of MADDPG came from the original code implementation (Lowe et al., 2017).
- Implementation of MAPPO came from the original code implementation (Yu et al., 2021).

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

All experiments were performed across 5 different seeds. Parameter sharing was utilized for all algorithms throughout the experiments for all environments with homogeneous state and action spaces. For multi-agent algorithms that perform centralized training (e.g., QMIX, COMA, MADDPG etc.), the global states were represented by the concatenation of all agents' local observations. We also used the 128-byte Atari RAM

as state inputs, rather than visual observations. This allows the algorithms to focus their learning on control rather than on both control and perception, improving learning efficiency.

4 EXPERIMENTAL RESULTS AND DISCUSSION

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

4.1 Cooperative

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

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

Additionally, since there is no unit collision in the Space Invaders environment (i.e., agents can move past each other without being blocked), they do not have to coordinate between themselves to achieve a high score in the environment; a good policy can be learned solely by having agents maximize their individual rewards. This explains the strong performance that was achieved by DQN. Also, since this is a cooperative task with both agents having identical goals, learning separate representations for individual agents is not very important; the learning of both agents assist each other. This is shown in figure 5B in Section 4.4, where the addition of an agent indicator did not yield any performance improvement for DQN on Space Invaders.

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

4.1.1 Multi-Agent Credit Assignment Problem in Fully Observable Settings

In this section, we attempt to study the effect of using a team reward signal, rather than individual reward signals on various independent and multi-agent algorithms in a fully observable environment. When team

rewards are the only rewards given, these reward signals are noisy for independent algorithms because the agent, which treats every other agent as part of the environment, does not know the actions taken by other agents. This makes it difficult for independent algorithms' agents to learn how their individual actions contribute to the team reward signal. We performed the experiments on Space Invaders, in which the default agents receive individual rewards from the environment. To study the effect of the multi-agent credit assignment problem, we performed two runs per algorithm, one with team rewards only, and the other with individual rewards only (i.e., agents are rewarded independently by the environment).

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

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

4.2 Competitive

The 2-player Pong environment was used for the competitive setting. All algorithms were first trained using parameter sharing with the addition of agent indicators (the effect of which is detailed in Section 4.4) for 60k episodes (1.2×10^6 steps), their network parameters were then saved. Since Pong is a zero-sum game, we evaluated them by putting them head-to-head against each other for 3 episodes for all possible combinations. After that, their positions were swapped, and the entire process was repeated. Swapping their positions is crucial for evaluation, because the first player (playing the right paddle) is always the serving player, therefore the first player always has an advantage over the second player (which plays the left paddle). This advantage is further exacerbated because the winning side always gets to serve subsequent openings. The entire evaluation process was repeated across all 5 seeds.

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

An interesting observation that can be made is the strong performance of independent algorithms, compared to other multi-agent algorithms. Since Pong is fully observable, critics that learn based on the joint observation of both agents do not necessarily provide any new information. Furthermore, since Pong is a highly reactive environment, an agent can learn a good policy solely by understanding how to position its paddle according to the trajectory of the ball (towards the agent). While learning on the joint action could allow agents to learn to better predict the incoming trajectory of the ball, it can be observed that the additional layer of complexity causes the sample efficiency to decrease and only yields diminishing returns.

In addition to the above factors, it is possible that parameter sharing benefited agents of independent algorithms by allowing them to learn better representations of both players, since they were trained to play as both players simultaneously. Had these algorithms trained without parameter sharing, there would likely be a larger performance difference between independent algorithms and opponent modelling algorithms such as DRON. Instead of treating other agents as part of the environment, opponent modelling allows agents to adapt more quickly to the opponent's changing strategies (He et al., 2016). However, the minimal

improvement DRON has over DQN suggests that in the Pong environment, an agent's policy may not be significantly affected by changes in the opponent agent's policy (i.e., individual agents can play the same way regardless of how their opponent played).

4.3 Mixed

In the Simple Tag (i.e., Predator-Prey) environment, the predators are incentivized to cooperate together to trap the prey, while the prey is incentivized to dodge the predators while staying within a predefined area. For our method of evaluation, we plot the training curves of the prey (figure 4B), and one of the predators (figure 4A), since all predators receive the same reward. Since the observation and action spaces differ between the predators and the prey, none of the agents have their parameters shared. We chose not to share the parameters of the predators to ensure that bias towards the predators was not introduced (since they would have 3 times the amount of data to learn from compared to the prey).

In the case of DQN, the prey successfully learned to minimize the number of collisions with the predators, which can be observed by the strong performance achieved by the prey (figure 4B). However, similar to PPO, since the predators were trained completely independently (i.e., their parameters were not shared), they did not manage to learn how to cooperate with one another to capture the prey (figure 4A). It is interesting to observe that MADDPG converged to a policy similar to DQN, with the difference being that its predators have learned to cooperate better, thus getting slightly higher rewards compared to DQN's predators (figure 4A). Subsequently, as a result of the higher rewards obtained by the predators, MADDPG achieves a slightly lower score for its prey (figure 4B).

MAPPO and RMAPPO, on the other hand, learned a different strategy. As we can observe from the comparatively noisier curves obtained from their predators and preys (figure 4A and 4B), there is a constant tug-of-war between the prey and the predators - as the predators learn how to cooperate better, their scores increase, which subsequently causes the prey to learn how to dodge, decreasing the predators' scores, and vice versa. Since the predators of MAPPO and RMAPPO achieves a much higher score compared to all other algorithms, this is indicative that the predators have successfully learned to cooperate to trap the prey.

4.4 Importance of Agent Indicator

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

Interestingly, in both cooperative environments, there was no noticeable improvement in the performance of DQN when an agent indicator was added (figure 5A and 5B). As was previously discussed, in the case of Space Invaders, since both agents have identical goals and similar representations, there is little need to distinguish between either agent. On the other hand, due to the partially observable nature of the Simple Reference environment, DQN performed similarly poorly, regardless of whether agent indicators were present. In this case, the addition of recurrence would have resulted in a much more significant difference instead, as was previously shown.

Conversely, for the Pong environment, even though it is also fully observable (akin to Space Invaders), the representation of both agents are not interchangeable. Utilizing parameter sharing without agent indicators, all algorithms struggled to learn due to the inability to tell which paddle were they controlling at every timestep. The only exception was RMAPPO (figure 6), which was able to condition on the sequence of previous observations and actions to infer which paddle was it controlling.

5 CONCLUSION

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

5.1 Cooperative

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

5.2 Competitive

In the Pong environment, we saw that DRQN and DQN were able to outperform all other 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. Furthermore, we showed that with the use of agent indicators, independent algorithms were able to learn robust policies for both competing agents using parameter sharing.

5.3 Mixed

In the Predator-Prey environment, we saw that since there were no parameter sharing to induce cooperation, predators from independent algorithms were unable to learn how to cooperate with each other to capture the prey. Conversely, in DQN we saw that its prey was able to achieve the highest score consistently, showing that the prey has learned to dodge the predators effectively while staying within the predefined area. Interestingly, we also saw how MADDPG's training curve for its predators and prey shows resemblance to that of DQN, suggesting that it also faced difficulties in learning strategies for the predators to coordinate and capture the prey. MAPPO and RMAPPO, on the other hand, were the only algorithms that managed to achieve high scores for their predators, suggesting that their predators have learned how to collaborate with each other to hunt the prey. The noisiness of their graphs suggest that there is a constant tug-of-war between the prey and the predators, as one tries to outsmart the other.

6 FUTURE WORK

In this section, we highlight some future work that could potentially bring more insights into having a broader understanding of dealing with non-stationarity and partial observability for independent algorithms, both of which are common in the multi-agent setting. In the Space Invaders environment, we observed that independent algorithms were able to learn well with just a team reward. Future work could be done to determine if this was only the case for fully observable environments, or under what conditions would independent algorithms still be able to cope with the multi-agent credit assignment problem. It would also be interesting to study the performance of non-recurrent variants of multi-agent algorithms such as QMIX

and COMA in fully observable environments. Since the experiments performed in this paper only included fully-observable competitive and mixed environments, future work can also include a more diverse set of environments, such as partially observable competitive and mixed environments.

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 CONTRIBUTIONS

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.

FUNDING

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

ACKNOWLEDGMENTS

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

REFERENCES

- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research* 47, 253–279
- Bellman, R. (1957). A markovian decision process. *Journal of Mathematics and Mechanics* 6, 679–684
- Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*
- Busoniu, L., Babuska, R., and De Schutter, B. (2008). A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 38, 156–172. doi:10.1109/TSMCC.2007.913919
- Canese, L., Cardarilli, G. C., Di Nunzio, L., Fazzolari, R., Giardino, D., Re, M., et al. (2021). Multi-agent reinforcement learning: A review of challenges and applications. *Applied Sciences* 11. doi:10.3390/app11114948
- Chang, Y.-h., Ho, T., and Kaelbling, L. (2004). All learning is local: Multi-agent learning in global reward games. In *Advances in Neural Information Processing Systems*, eds. S. Thrun, L. Saul, and B. Schölkopf (MIT Press), vol. 16
- Choi, S., Yeung, D.-Y., and Zhang, N. (2000). An environment model for nonstationary reinforcement learning. In *Advances in Neural Information Processing Systems*, eds. S. Solla, T. Leen, and K. Müller (MIT Press), vol. 12
- Foerster, J., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. (2018). Counterfactual multi-agent policy gradients. *Proceedings of the AAAI Conference on Artificial Intelligence* 32

- Gronauer, S. and Diepold, K. (2021). Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, 1–49
- Gupta, J. K., Egorov, M., and Kochenderfer, M. (2017). Cooperative multi-agent control using deep reinforcement learning. In *Autonomous Agents and Multiagent Systems*, eds. G. Sukthankar and J. A. Rodriguez-Aguilar (Cham: Springer International Publishing), 66–83
- Hausknecht, M. and Stone, P. (2015). Deep recurrent q-learning for partially observable mdps. In *2015 AAAI fall symposium series*
- He, H., Boyd-Graber, J., Kwok, K., and Daumé, H. (2016). Opponent modeling in deep reinforcement learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48 (JMLR.org)*, ICML'16, 1804–1813
- Hernandez-Leal, P., Kaisers, M., Baarslag, T., and de Cote, E. M. (2017). A survey of learning in multiagent environments: Dealing with non-stationarity. *arXiv preprint arXiv:1707.09183*
- Hernandez-Leal, P., Kartal, B., and Taylor, M. E. (2019). A survey and critique of multiagent deep reinforcement learning. *Journal of Autonomous Agents and Multi-Agent Systems (JAAMAS)* 33, 750–797
- Kraemer, L. and Banerjee, B. (2016). Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing* 190, 82–94. doi:<https://doi.org/10.1016/j.neucom.2016.01.031>
- [Dataset] Li, M. (2020). Machin. <https://github.com/iffix/machin>
- Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, P., and Mordatch, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*
- Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M. (2018). Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research* 61, 523–562
- Markov, A. A. (1954). The theory of algorithms. *Trudy Matematicheskogo Instituta Imeni VA Steklova* 42, 3–375
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., et al. (2015). Human-level control through deep reinforcement learning. *nature* 518, 529–533
- Mordatch, I. and Abbeel, P. (2017). Emergence of grounded compositional language in multi-agent populations. *arXiv preprint arXiv:1703.04908*
- Oliehoek, F. A., Spaan, M. T., and Vlassis, N. (2008). Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research* 32, 289–353
- Papoudakis, G., Christianos, F., Rahman, A., and Albrecht, S. V. (2019). Dealing with non-stationarity in multi-agent deep reinforcement learning. *arXiv preprint arXiv:1906.04737*
- [Dataset] Raffin, A., Hill, A., Ernestus, M., Gleave, A., Kanervisto, A., and Dormann, N. (2019). Stable baselines3. <https://github.com/DLR-RM/stable-baselines3>
- Rashid, T., Samvelyan, M., Schroeder, C., Farquhar, G., Foerster, J., and Whiteson, S. (2018). QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, eds. J. Dy and A. Krause (PMLR), vol. 80 of *Proceedings of Machine Learning Research*, 4295–4304
- Schaul, T., Quan, J., Antonoglou, I., and Silver, D. (2015). Prioritized experience replay. *arXiv preprint arXiv:1511.05952*
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*
- Shapley, L. S. (1953). Stochastic games. *Proceedings of the National Academy of Sciences* 39, 1095–1100. doi:[10.1073/pnas.39.10.1095](https://doi.org/10.1073/pnas.39.10.1095)

- 418 Shoham, Y. and Leyton-Brown, K. (2008). *Multiagent systems: Algorithmic, game-theoretic, and logical*
 419 *foundations* (Cambridge University Press)
- 420 [Dataset] starry sky6688 (2019). Starcraft. <https://github.com/starry-sky6688/>
 421 StarCraft
- 422 Sukhbaatar, S., Szlam, A., and Fergus, R. (2016). Learning multiagent communication with
 423 backpropagation. In *Proceedings of the 30th International Conference on Neural Information Processing*
 424 *Systems* (Red Hook, NY, USA: Curran Associates Inc.), NIPS'16, 2252–2260
- 425 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction* (MIT press)
- 426 Tampuu, A., Matiisen, T., Kodelja, D., Kuzovkin, I., Korjus, K., Aru, J., et al. (2017). Multiagent
 427 cooperation and competition with deep reinforcement learning. *PLOS ONE* 12, 1–15. doi:10.1371/
 428 journal.pone.0172395
- 429 Tan, M. (1993). Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings*
 430 *of the tenth international conference on machine learning*. 330–337
- 431 Terry, J. K. and Black, B. (2020). Multiplayer support for the arcade learning environment. *arXiv preprint*
 432 *arXiv:2009.09341*
- 433 Terry, J. K., Black, B., and Hari, A. (2020a). Supersuit: Simple microwrappers for reinforcement learning
 434 environments. *arXiv preprint arXiv:2008.08932*
- 435 Terry, J. K., Black, B., Jayakumar, M., Hari, A., Sullivan, R., Santos, L., et al. (2020b). Pettingzoo: Gym
 436 for multi-agent reinforcement learning. *arXiv preprint arXiv:2009.14471*
- 437 Van Hasselt, H., Guez, A., and Silver, D. (2016). Deep reinforcement learning with double q-learning. In
 438 *Proceedings of the AAAI conference on artificial intelligence*. vol. 30
- 439 Wang, Z., Schaul, T., Hessel, M., Hasselt, H., Lanctot, M., and Freitas, N. (2016). Dueling network
 440 architectures for deep reinforcement learning. In *International conference on machine learning* (PMLR),
 441 1995–2003
- 442 Yu, C., Velu, A., Vinitisky, E., Wang, Y., Bayen, A., and Wu, Y. (2021). The surprising effectiveness of
 443 mappo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*
- 444 Zawadzki, E., Lipson, A., and Leyton-Brown, K. (2014). Empirically evaluating multiagent learning
 445 algorithms. *arXiv preprint arXiv:1401.8074*
- 446 Zhang, K., Yang, Z., and Başar, T. (2021). Multi-agent reinforcement learning: A selective overview of
 447 theories and algorithms. *Handbook of Reinforcement Learning and Control*, 321–384

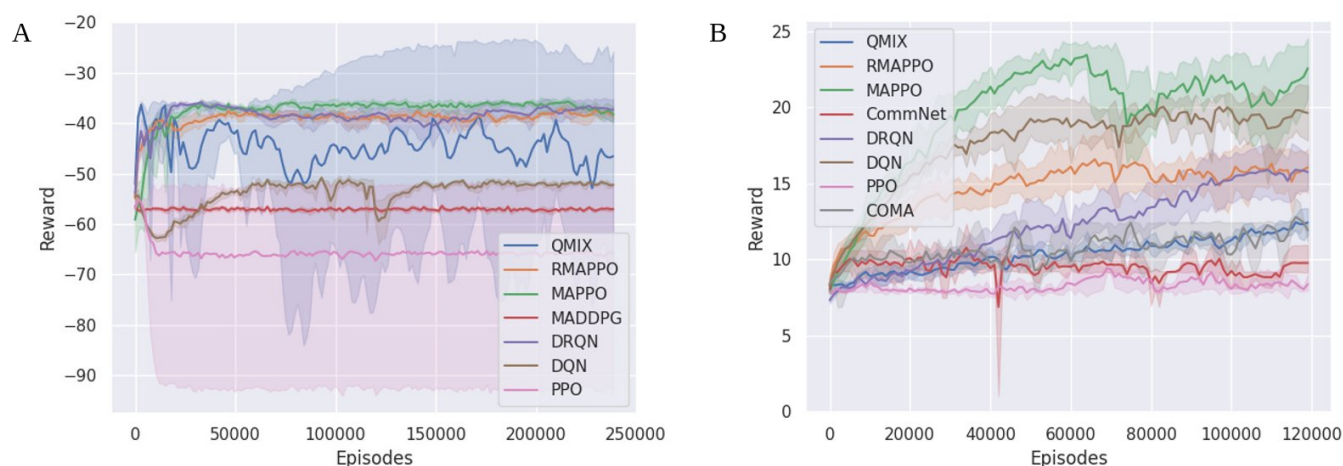


Figure 1. Training curves of various algorithms in two cooperative environments. For every algorithm, the solid line represents the mean reward per episode, while the shaded region represents the 95% confidence interval around the mean. (A) shows training curve for Simple Reference environment, (B) shows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

FIGURE CAPTIONS

Table 1. Final scores (mean and standard deviation) of algorithms obtained over the last 100 episodes across all 5 seeds in the Space Invaders environment.

Algorithms	Space Invaders
QMIX	12.5±3.28
RMAPPO	16.2±3.31
MAPPO	22.5±3.45
CommNet	9.78±0.98
COMA	12.2±1.61
DRQN	15.7±3.70
DQN	19.8±3.52
PPO	7.92±2.69

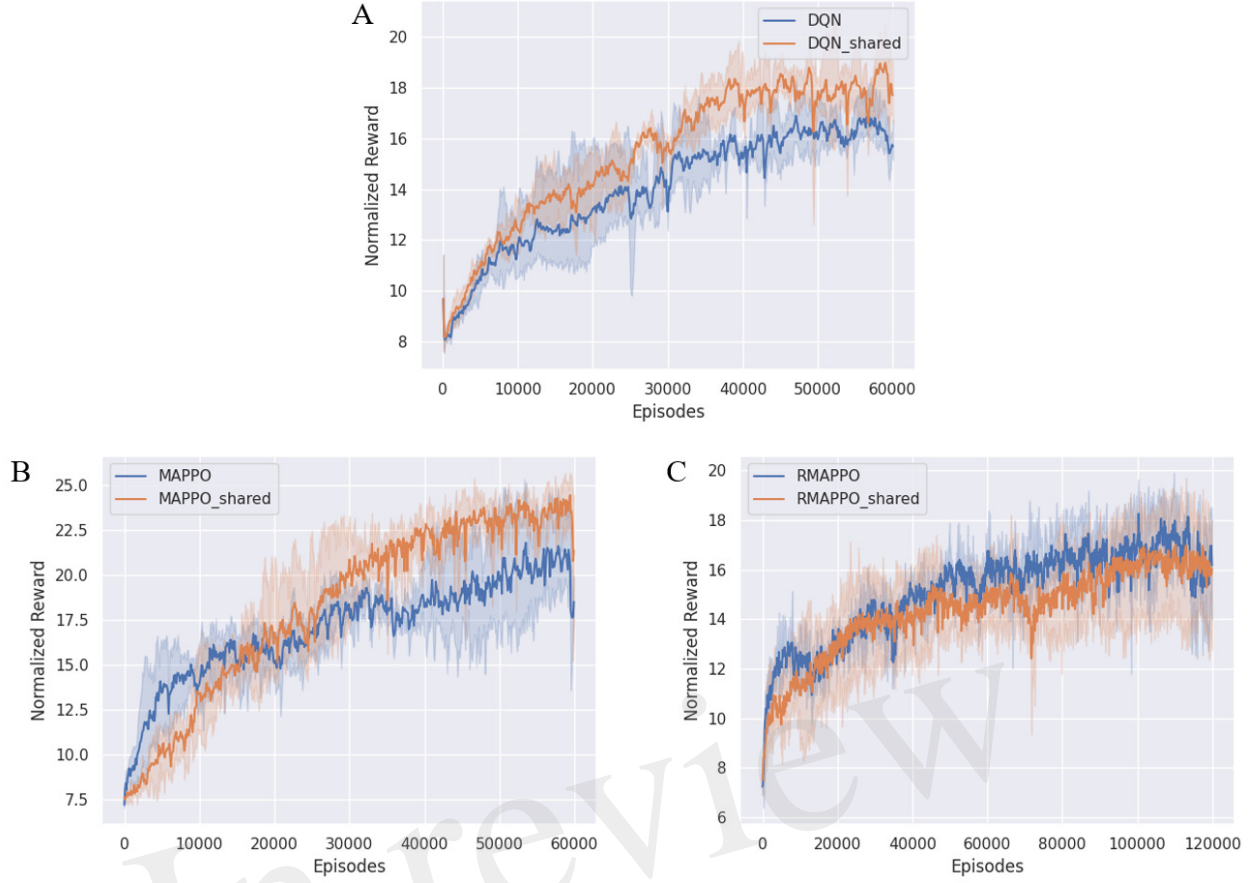


Figure 2. Training curves of various algorithms in Space Invaders, comparing when individual rewards are given (blue) to when team rewards are given (orange). (A) shows training curve of DQN, (B) shows training curve of MAPPO, (C) shows training curve of RMAPPO.

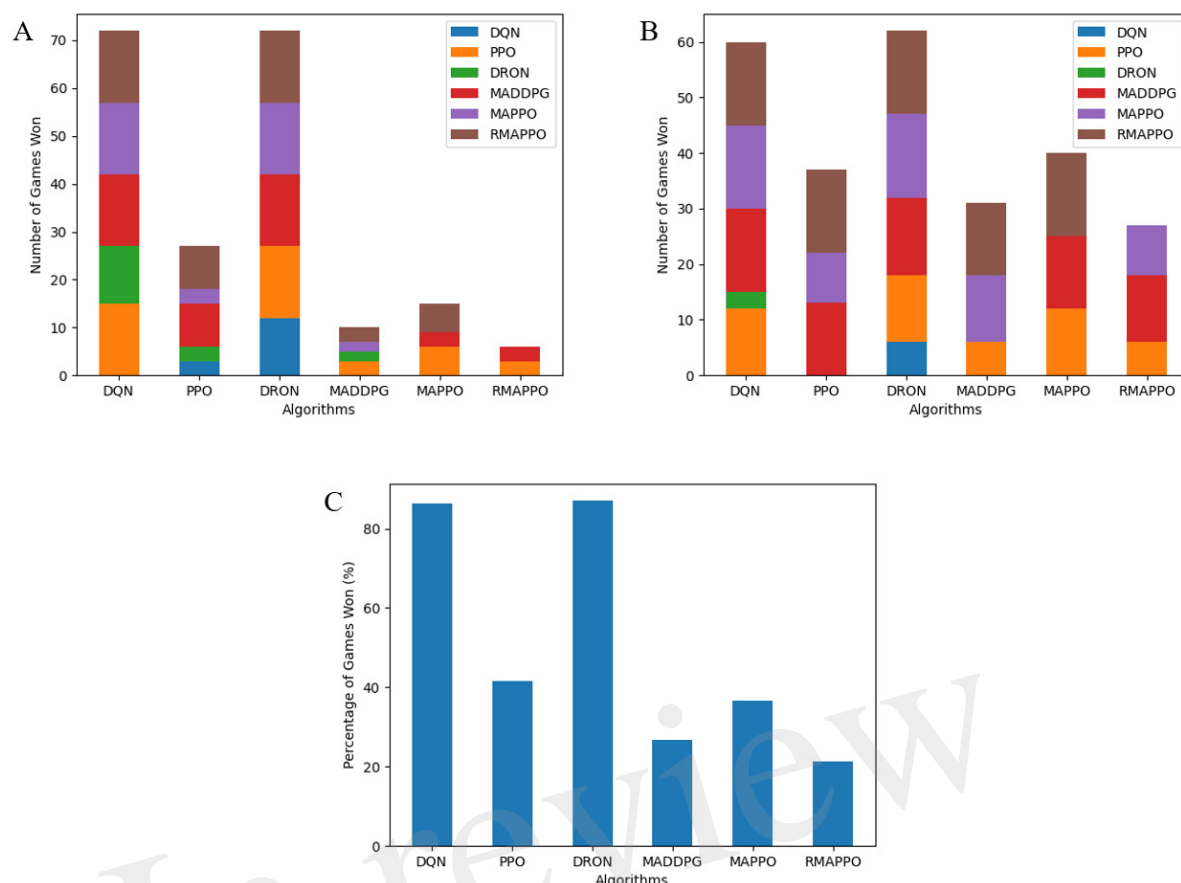


Figure 3. Performance of various algorithms when playing against other algorithms in Pong. (A) shows the number of games won as the first player, (B) shows the number of games won as the second player, (C) shows the overall win rate percentage.

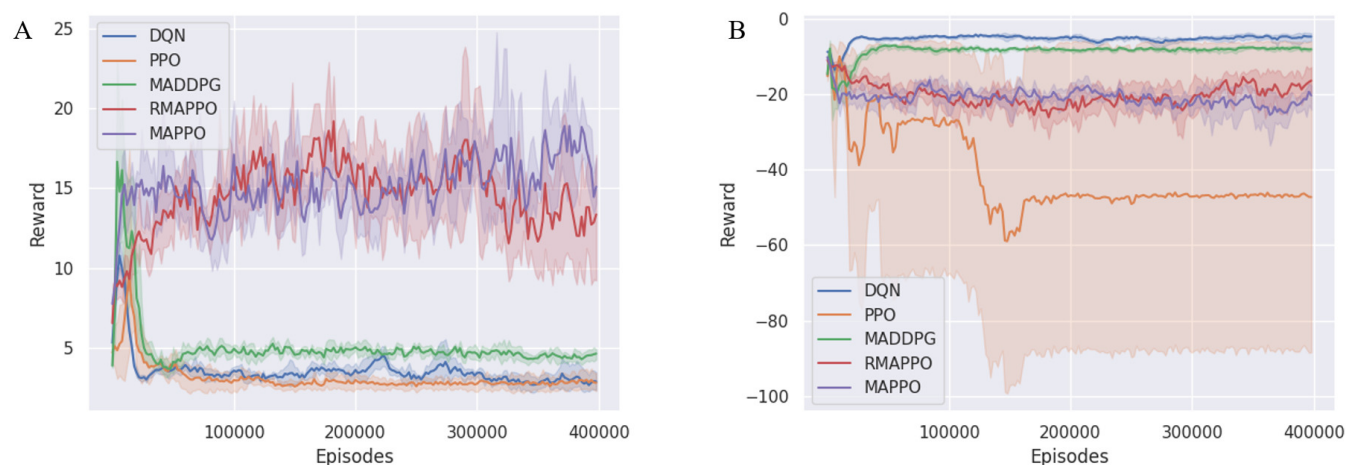


Figure 4. Training curves of various algorithms in the Simple Tag, a Predator-Prey environment. (A) shows the reward of a predator (all predators obtain the same reward), (B) shows the reward of the prey.

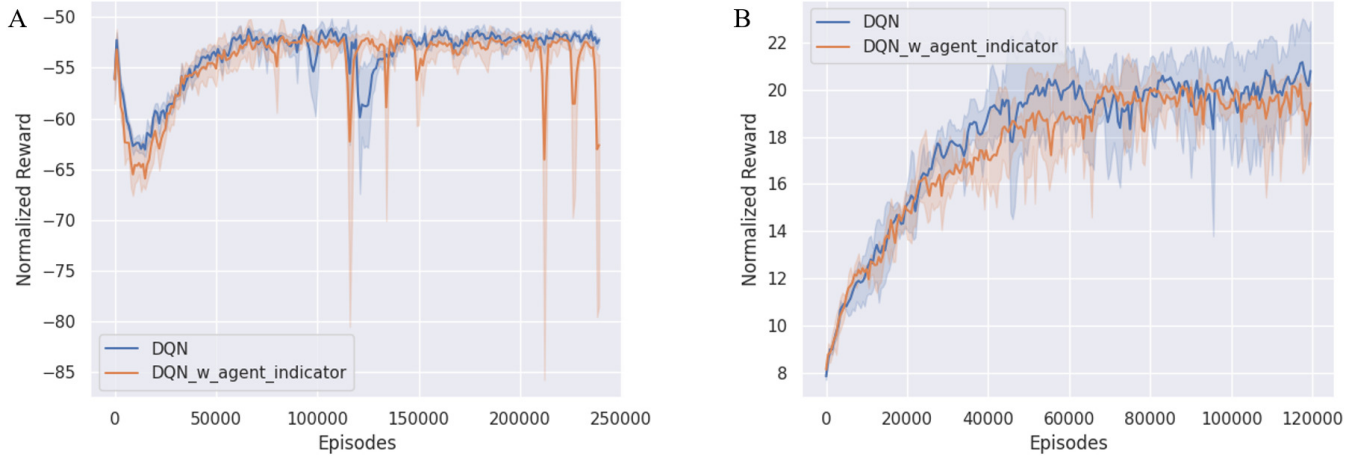


Figure 5. Comparing DQN with (blue) and without (orange) agent indicators in (A) Simple Reference and (B) Space Invaders environment.

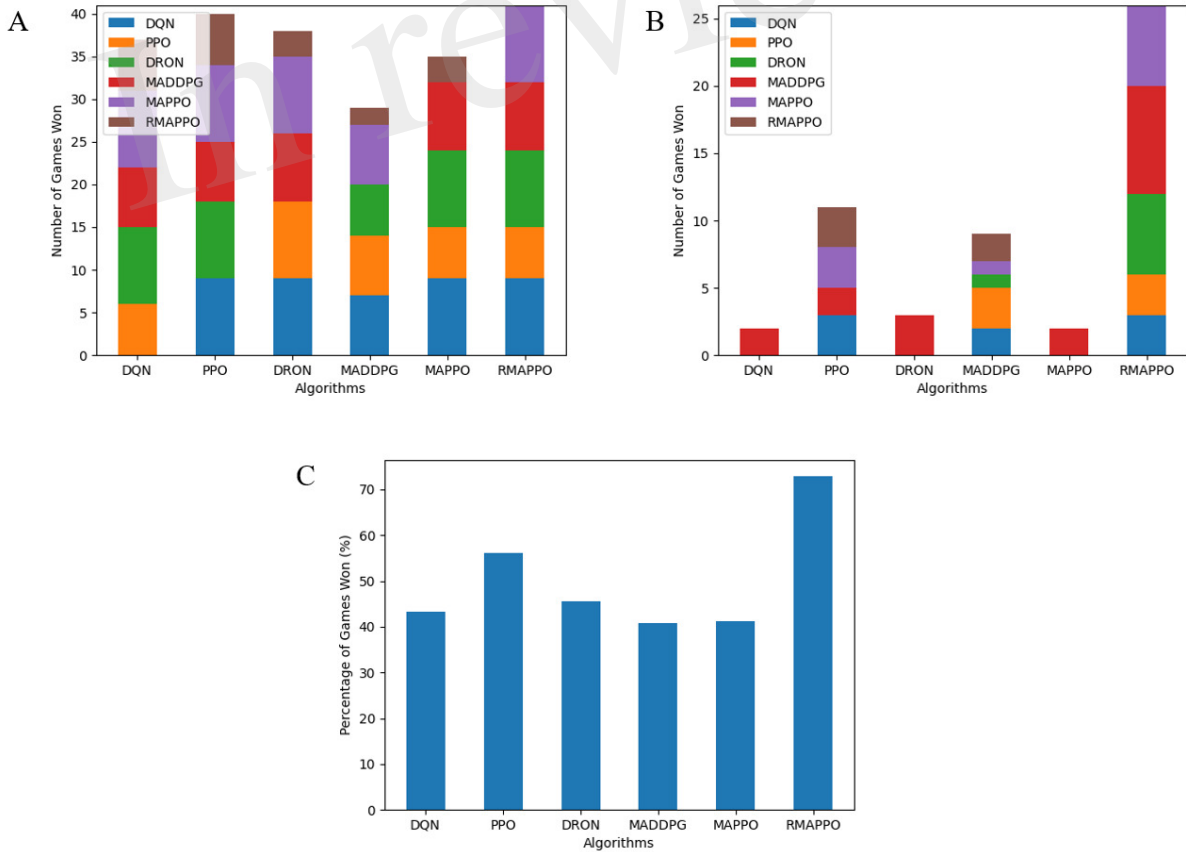


Figure 6. Performance of various algorithms when playing against other algorithms in Pong without agent indicators across 3 seeds. (A) shows the number of games won as the first player, (B) shows the number of games won as the second player, (C) shows the overall win rate percentage.

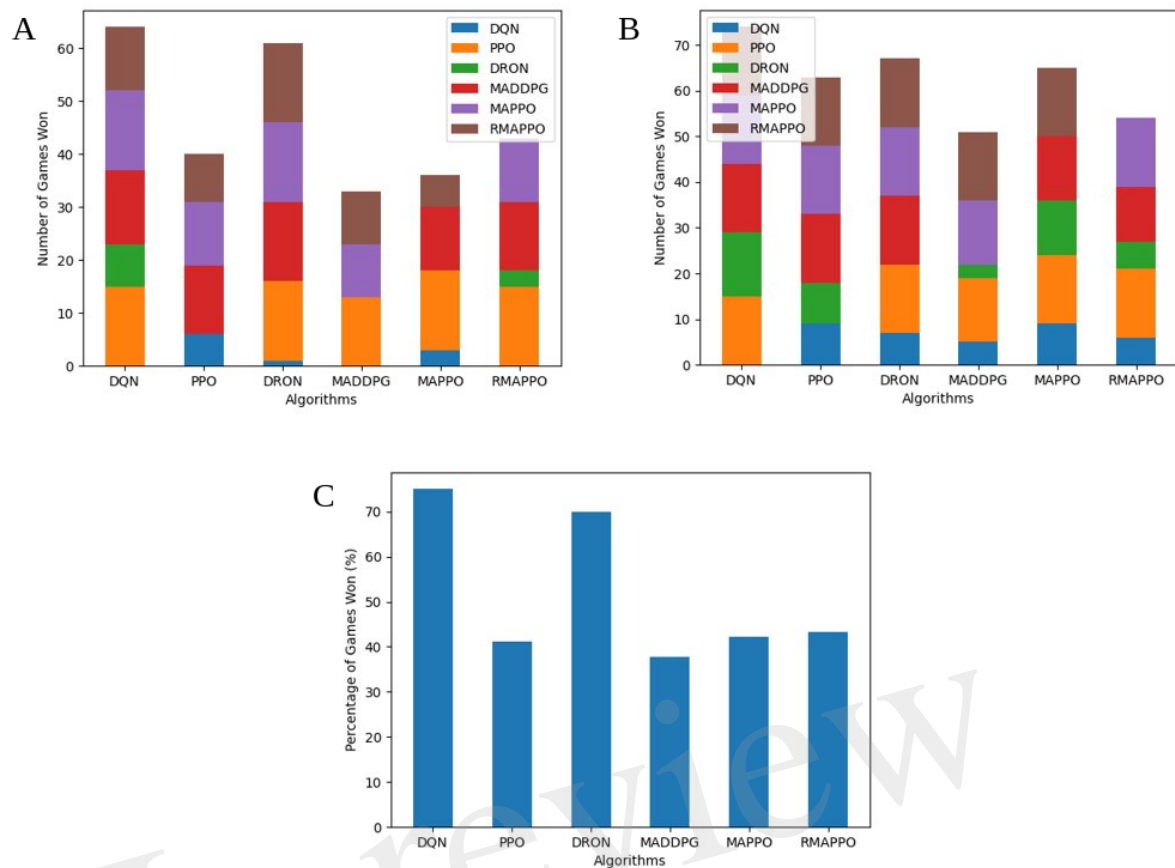


Figure 7. Performance of various algorithms when playing against other algorithms in the Boxing environment. (A) shows the number of games won as the first player, (B) shows the number of games won as the second player, (C) shows the overall win rate percentage.

Figure 1.JPEG

In review

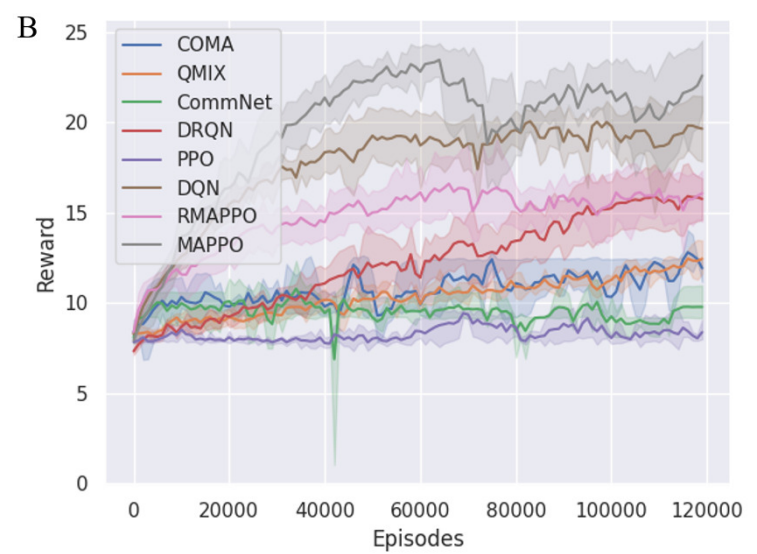
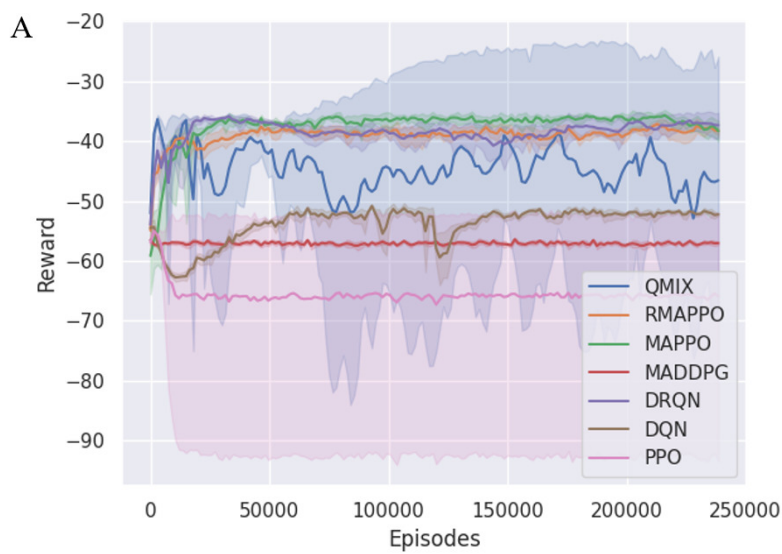


Figure 2.JPEG

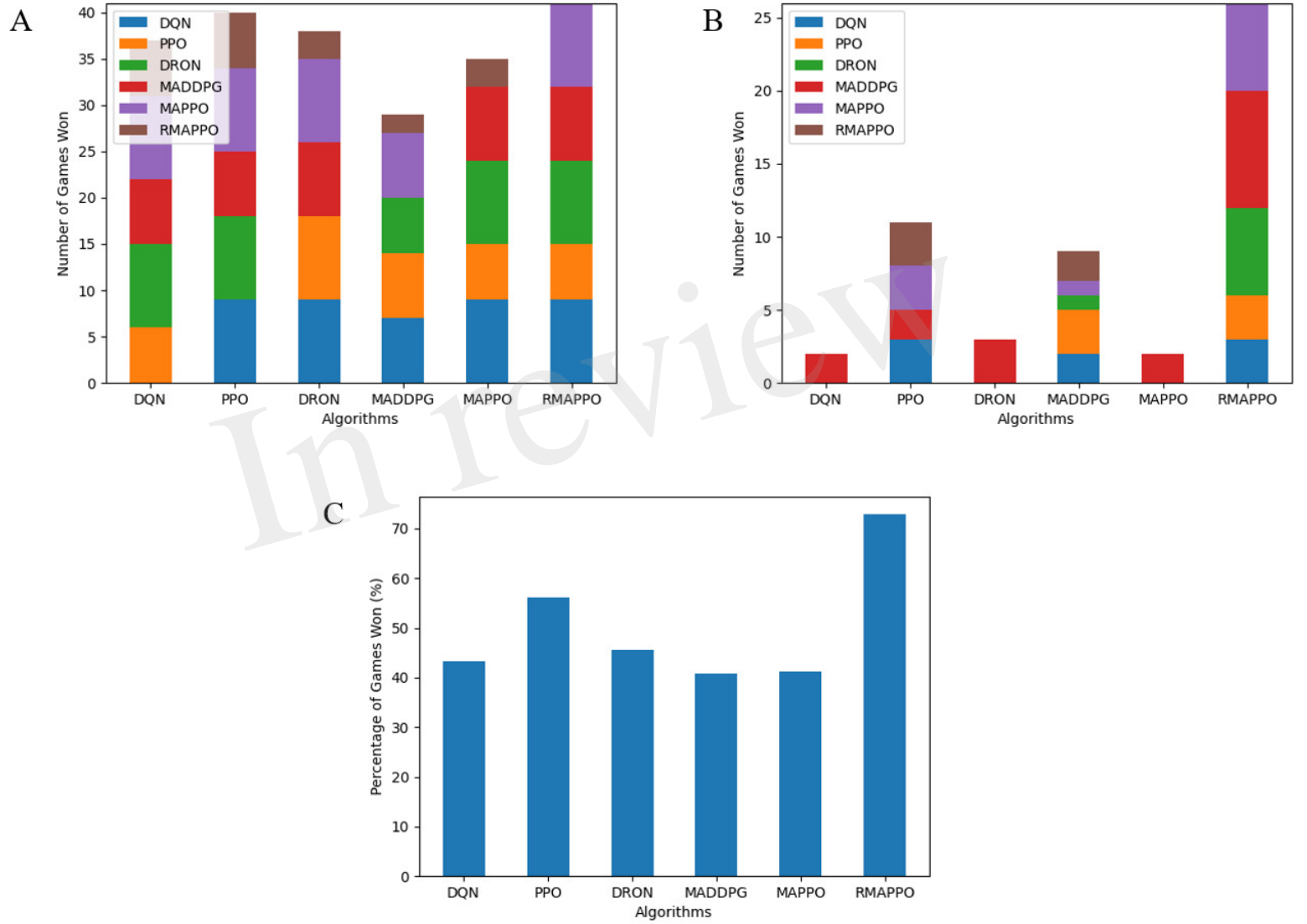


Figure 3.JPEG

In review

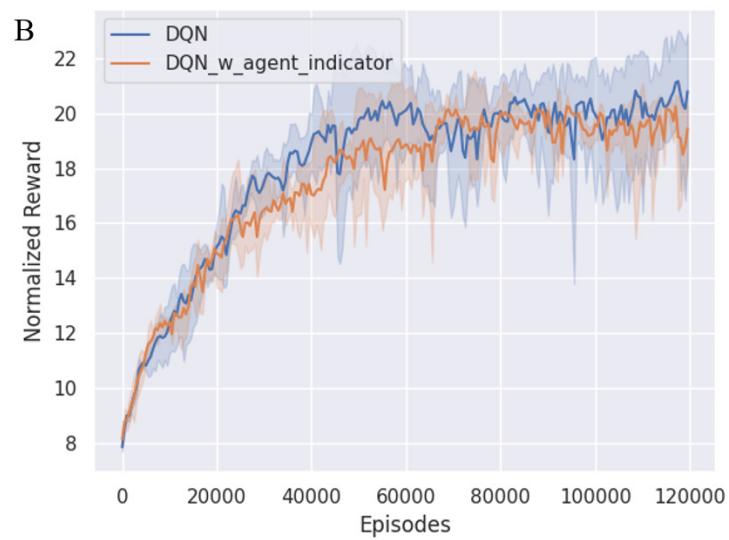
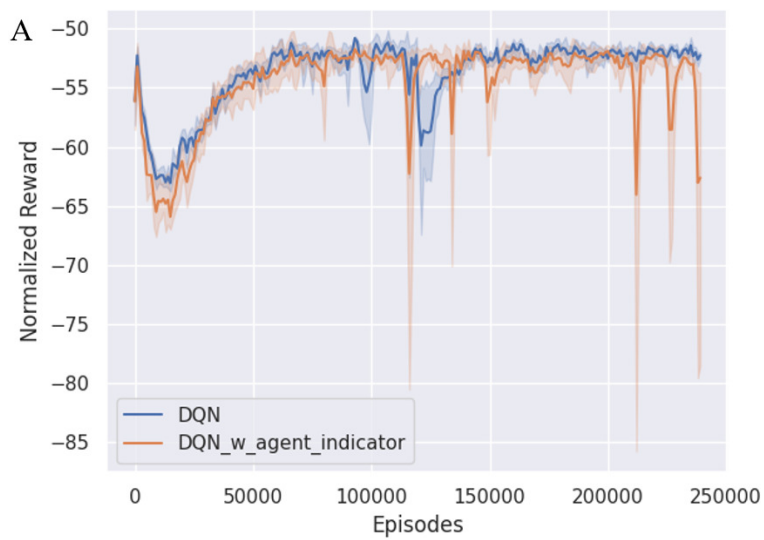


Figure 4.JPEG

In review

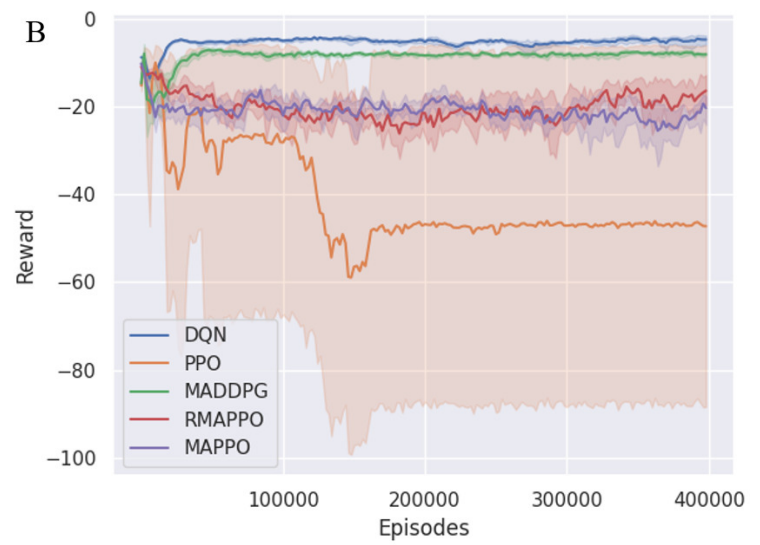
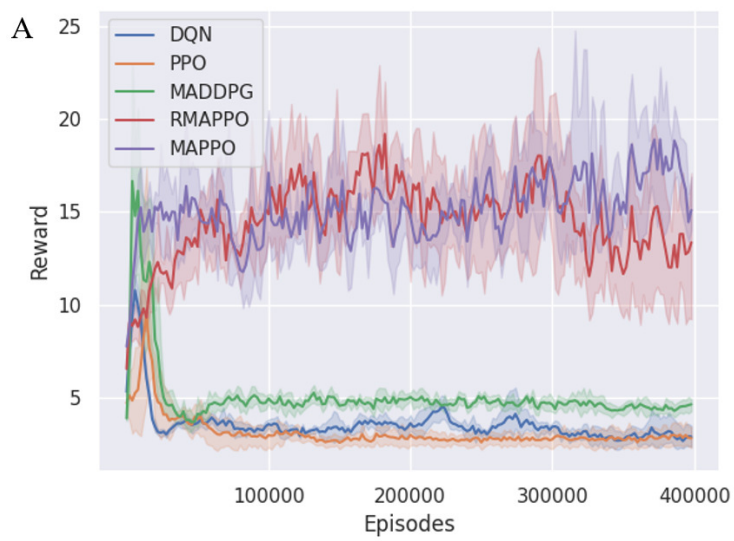


Figure 5.JPEG

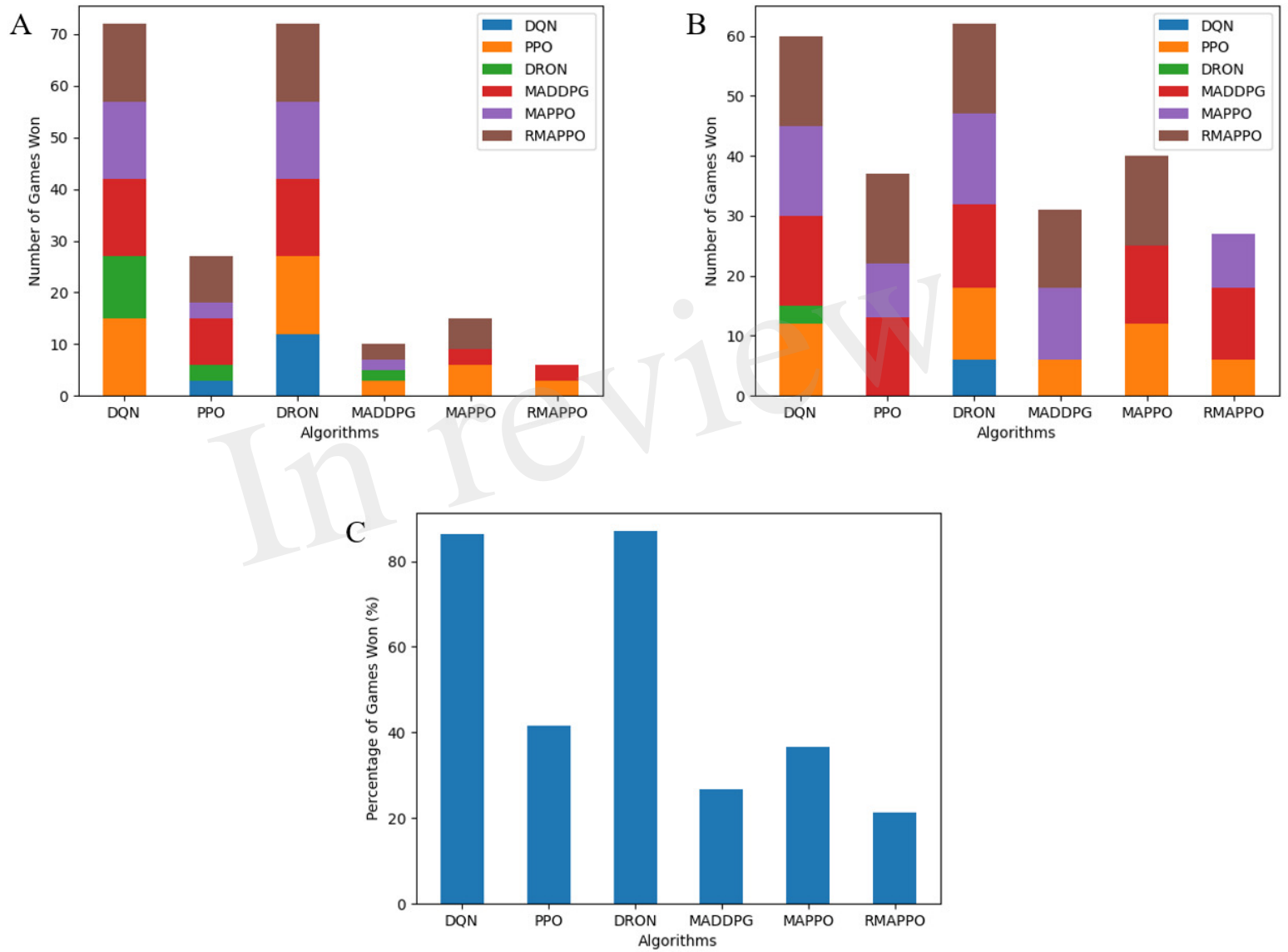


Figure 6.JPEG

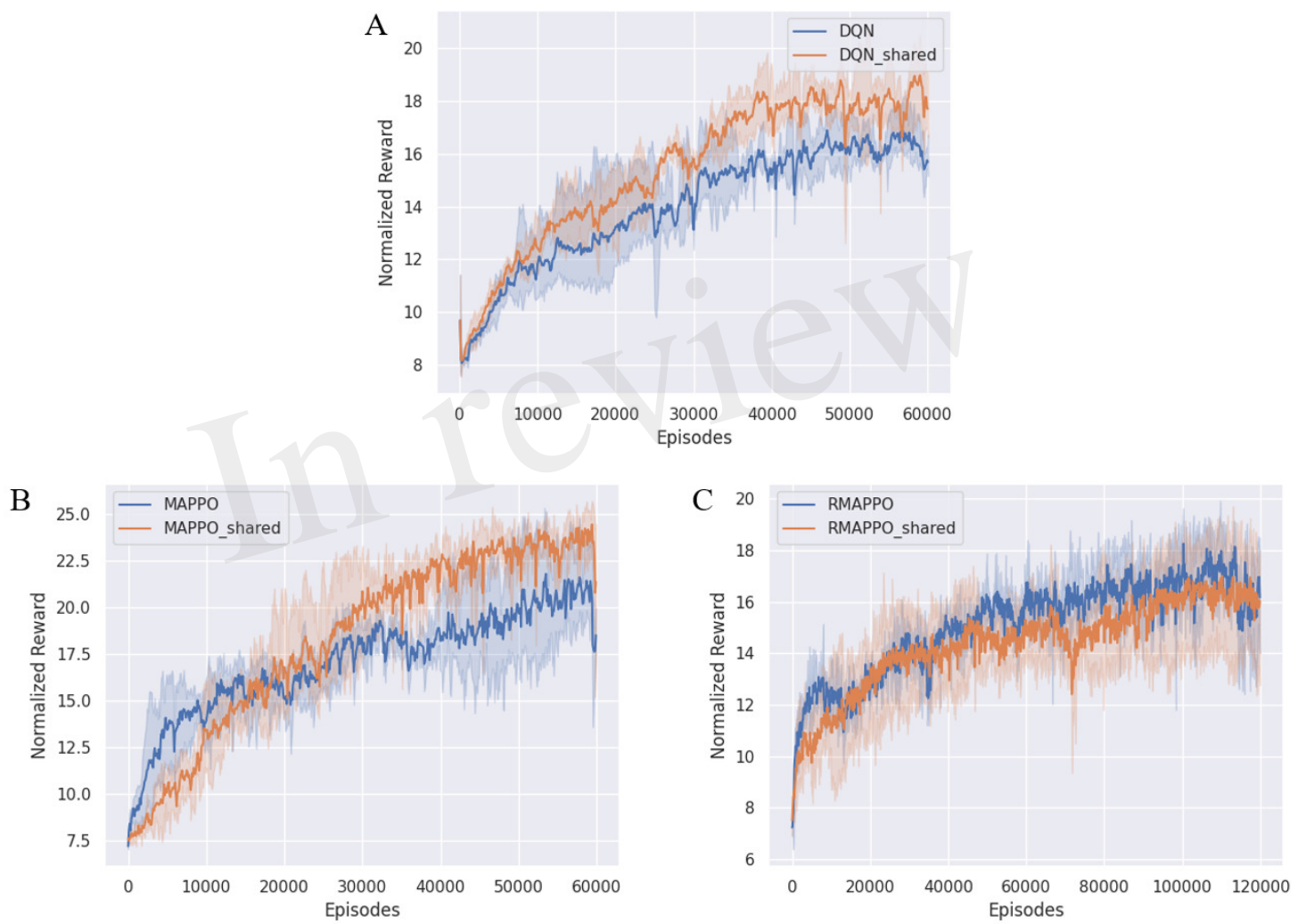


Figure 7.JPEG

