
Investigation of Independent Reinforcement Learning Algorithms in Multi-Agent Environments

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Independent reinforcement learning algorithms have no theoretical guarantees for
2 finding the best policy in multi-agent settings. However, in practice, prior works
3 have reported good performance with independent algorithms in some domains
4 and bad performance in others. Moreover, a comprehensive study of the strengths
5 and weaknesses of independent algorithms is lacking in the literature. In this
6 paper, we carry out an empirical comparison of the performance of independent
7 algorithms on four PettingZoo environments that span the three main categories
8 of multi-agent environments, i.e., cooperative, competitive, and mixed. We show
9 that in fully-observable environments, independent algorithms can perform on par
10 with multi-agent algorithms in cooperative and competitive settings. For the mixed
11 environments, we show that agents trained via independent algorithms learn to
12 perform well individually, but fail to learn to cooperate with allies and compete
13 with enemies. We also show that adding recurrence improves the learning of
14 independent algorithms in cooperative partially observable environments.

15 1 Introduction

16 One of the simplest ways to apply reinforcement learning in multi-agent settings is to assume that
17 all agents are independent of each other. In other words, every other agent is seen as part of the
18 environment from any agent’s perspective. Independent algorithms (i.e., single-agent algorithms)
19 face the issue of non-stationarity in the multi-agent domain due to the violation of the Markovian
20 property in a Markov Decision Process [1]. As a result, independent algorithms lack convergence
21 guarantees, and are not theoretically sound in the multi-agent setting [2]. Despite these shortcomings,
22 independent algorithms have the advantage of requiring lower computational resources and being
23 easier to scale to large environments than traditional multi-agent algorithms which perform exact
24 opponent modelling of each agent. In practice, prior works have reported mixed performance for
25 independent algorithms in different multi-agent domains [3]–[9]. However, a study of the strengths
26 and weaknesses of independent algorithms across various categories within the multi-agent domain is
27 lacking in the literature.

28 In this paper, we investigate the empirical performance of independent algorithms in multi-agent
29 settings, and compare them to various multi-agent algorithms under the Centralized Training and De-
30 centralized Execution scheme [10], [11]. We evaluate these algorithms on 4 multi-agent environments
31 from the PettingZoo library [12], which span the 3 main categories of multi-agent environments (i.e.,
32 cooperative, competitive and mixed) [13]–[16]. We show that independent algorithms can perform on
33 par with multi-agent algorithms in the cooperative, fully-observable setting, and adding recurrence
34 allows them to perform well compared to multi-agent algorithms in partially observable environments.
35 In the competitive setting, we show that parameter sharing alongside the addition of agent indicators
36 allow independent algorithms to outperform some multi-agent algorithms, such as Multi-Agent
37 Proximal Policy Optimization [17], and Multi-Agent Deep Deterministic Policy Gradient [8], in

38 fully-observable environments. For the mixed setting, we show that agents of independent algorithms
39 learn to perform well individually, but fail in learning to cooperate with allies and compete against
40 enemies.

41 **2 Background Information**

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

44 **2.1 Reinforcement Learning**

45 In Reinforcement Learning (RL), an agent interacts with the environment by making sequential
46 decisions [18]. At every time step, denoted as t , the agent observes a state s_t from the environment,
47 and takes an action a_t . This action is executed in the environment, which returns a reward r_t and the
48 next state s_{t+1} that are determined by the reward function $R(s_t, a_t)$ and the transition probability,
49 $P(s_{t+1}|s_t, a_t)$, respectively. Critically, $R(s_t, a_t)$ and $P(s_{t+1}|s_t, a_t)$ are part of the environment, and
50 are usually unknown to the agent of a model-free RL algorithm. Since the transition probability
51 $P(s_{t+1}|s_t, a_t)$ conditions the next state s_{t+1} purely on the current state s_t and action a_t , it satisfies
52 the Markov property [19]. This interaction between the agent and the environment is called a Markov
53 Decision Process (MDP) [20]. The objective of an RL agent is to learn a policy $\pi(a_t|s_t)$, which maps
54 a state to an action that maximizes the expected cumulative reward it receives from the environment.
55 This is written as $\sum_t \gamma^t r_t$, where $\gamma \in [0, 1)$ represents a discount factor on future rewards.

56 **2.2 Multi-Agent Reinforcement Learning**

57 The single-agent MDP framework is extended to the Multi-Agent Reinforcement Learning (MARL)
58 setting in the form of stochastic games [21]. In an N -agent stochastic game, at every time step,
59 each of the n agents, identified by $j \in \{1, 2, \dots, n\}$ across all agents, takes an action u_t^j . The joint
60 action $\mathbf{u}_t \triangleq \{u_t^1, \dots, u_t^N\}$ determines the rewards obtained by each agent. State transitions of the
61 environment are determined by the transition probability $P(s_{t+1}|s_t, \mathbf{u}_t)$, which conditions on the
62 state and the joint action at timestep t .

63 **2.3 Centralized Training and Decentralized Execution**

64 While it is technically possible to learn a centralized controller that maps a state to a distribution
65 over the joint action space, the number of possible combinations of actions grows exponentially
66 with the number of agents. This makes centralized control intractable for environments with many
67 agents. Therefore, this paper is mainly focused on multi-agent algorithms which correspond to a
68 Centralized Training and Decentralized Execution (CTDE) scheme [10], [11]. A CTDE algorithm
69 has two phases. During the control phase, where policies are deployed in the environment, rather
70 than using a centralized controller to take actions for all agents, decentralized agents make decisions
71 based on their individual observations. During the prediction phase, centralized training is performed,
72 which allows for extra information (e.g. the state) to be utilized, as long as this is not required during
73 the control phase.

74 **2.4 Cooperative, Competitive and Mixed**

75 This paper follows the convention of classifying every multi-agent algorithm and environment studied
76 into one of three categories – cooperative, competitive, or mixed (cooperative-competitive) [13]–[16].

77 In the cooperative setting, agents collaborate with each other to achieve a common goal. As a result,
78 it is very common for all agents to share the same reward function (i.e., a team goal) [22]. Also
79 known as the multi-agent credit assignment problem, every agent has to deduce its own contributions
80 from the team reward [22]. Algorithms studied in this paper that explicitly address the multi-agent
81 credit-assignment problem include QMIX [9] and Counterfactual Multi-Agent Policy Gradients
82 (COMA) [7]. Additionally, the CommNet [23] extension on top of COMA is utilized for specific
83 cooperative environments. Other multi-agent algorithms that are considered for the cooperative
84 scenario include Multi-Agent Deep Deterministic Policy Gradient (MADDPG) [8] and Multi-Agent
85 Proximal Policy Optimization (MAPPO) [17].

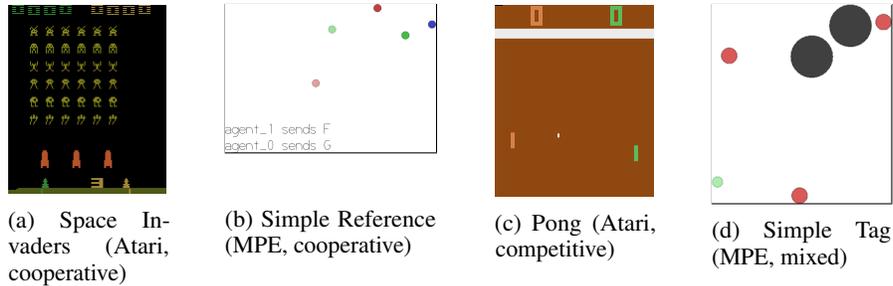


Figure 1: The four PettingZoo environments used in the experiments. All figures were obtained from <https://pettingzoo.ml/>

86 In the competitive setting, agents play a zero-sum game, where one agent’s gain is another agent’s
 87 loss. In other words, $\sum_a r(s, u, a) = 0 \forall s, u$. Algorithms that are studied specifically in this paper
 88 include Deep Reinforcement Opponent Network (DRON) [24], MADDPG and MAPPO. MADDPG
 89 and MAPPO learn a separate critic for every agent, which gives the algorithms flexibility to learn
 90 different behaviours for agents with different reward functions.

91 In a mixed or cooperative-competitive setting, environments are neither zero-sum (competitive) nor
 92 cooperative, and they do not necessarily need to be general-sum either. A common setting would
 93 be environments that require every agent to cooperate with some agents, and compete with others
 94 [13]–[15]. MADDPG and MAPPO are used here for the same reason as the competitive setting.

95 2.5 Independent Algorithms and Non-Stationarity

96 One naive approach for applying single-agent RL to the multi-agent setting would be the use of
 97 independent learners, where each agent treats every other agent as part of the environment, and learns
 98 purely based on individual observations. In a multi-agent setting, the transition probability P and
 99 reward function R are conditioned on the joint action u . Since all agents in the environment are
 100 learning, their policies constantly change. Therefore, from each independent learner’s perspective, the
 101 transition probability and reward function appear non-stationary, due to the lack of awareness of other
 102 agents’ actions. This violates the Markovian property of an MDP, which then causes independent
 103 algorithms to be slow to adapt to other agents’ changing policies, and as a result, face difficulties in
 104 converging to a good policy [24]–[26].

105 In this paper, we chose to use a popular off-policy algorithm, Deep Q-Network (DQN) [27], and
 106 an on-policy algorithm, Proximal Policy Optimization (PPO) [28]. In specific partially observable
 107 environments, Deep Recurrent Q-Network (DRQN) [29] is also utilized. Following the work of
 108 Gupta et al. [30], parameter sharing is utilized for all independent algorithms, such that experiences
 109 from all agents are trained simultaneously using a single network. This allows the training to be more
 110 efficient [30]. The aforementioned independent algorithms are tested in all 3 categories of multi-agent
 111 environments.

112 3 Experimental Setup

113 In this section, we introduce the environments used for the experiments, specify the various prepro-
 114 cessing that were applied, and other relevant implementation details.

115 3.1 Environments

116 The experiments were performed on multiple multi-agent environments from the PettingZoo library
 117 [12], which contains the Multi-Agent Particle Environments (MPE) [8], [31] and multi-agent variants
 118 of the Atari 2600 Arcade Learning Environment (ALE) [32], [33].

119 For the cooperative setting, experiments were performed on a modified version of the 2-player
 120 Space Invaders [32], [33], and the Simple Reference MPE environment [8], [31]. In Space Invaders
 121 (Fig. 1a), both agents share the common goal of shooting down all aliens. To make Space Invaders
 122 cooperative, we removed the (positive) reward that is given to a player whenever the other player

123 gets hit. Additionally, the environment rewards every agent individually by default. Since a number
124 of cooperative multi-agent algorithms (e.g., QMIX and COMA) assume that only a team reward is
125 given, we modified the reward function such that a team reward is given instead (i.e., both agents
126 receive the sum of their individual rewards). This setup exemplifies the multi-agent credit assignment
127 problem, the effect of which is studied more closely in the Section 4.1.1. On the other hand, in the
128 Simple Reference environment (Fig. 1b), two agents are rewarded by how close they are to their
129 target landmark. However, the target landmark of an agent is only known by the other agent, as a
130 result communication is required for both agents to navigate successfully to their target landmarks.

131 For the competitive setting, we performed experiments on the 2-player variant of the original Atari
132 Pong environment (Fig. 1c). For the mixed setting, we opted for the Simple Tag MPE environment
133 (Fig. 1d), which is a Predator-Prey environment [31]. This environment consists of 4 agents – 3
134 predators and a prey. The prey travels faster and has to avoid colliding with the predators, while the 3
135 predators travel slower and have to work together to capture the prey. The rewards received by the
136 prey and a predator sum to 0 (i.e., the prey gets a negative reward for collision, while the predators get
137 rewarded positively), and all predators receive the same reward. The prey is also negatively rewarded
138 if it strays away from the predefined area (a 1×1 unit square). This environment is general-sum
139 because it contains 3 predators and a single prey.

140 3.2 Preprocessing

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

143 For the Atari environments, following the recommendations of Marlos et al. [34], we performed the
144 following preprocessing - reward clipping, sticky actions, frame skipping, and no-op resets. The
145 number of steps per episode was also set to a limit of 200 for both Atari environments, as that
146 yielded the best results in general. Furthermore, the action spaces for both Atari environments were
147 shrunk to their effective action spaces in order to improve learning efficiency. For Pong specifically,
148 we also concatenated a one-hot vector of the agent’s index to the observations so that independent
149 algorithms can differentiate one from the other when parameter sharing is utilized. Further details of
150 the preprocessing performed can be found in Appendix A.

151 3.3 Implementation

152 Implementations of all algorithms were based on open-sourced libraries/reference implementations.
153 Default hyperparameters were used for all algorithms, and no hyperparameter tuning was performed.
154 Details of implementations, alongside their hyperparameters, can be found in Appendix C.

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

162 4 Experiment Results

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

166 4.1 Cooperative

167 **Simple Reference** We ran the various algorithms on the Simple Reference environment for 240k
168 episodes (6×10^6 steps). From Fig. 2a, it could be observed that all independent algorithms converged
169 to a lower score, except for DRQN, whose recurrence allowed it to vastly outperform DQN and
170 converge to a score on par with multi-agent algorithms. However, this trend was not observed
171 when comparing MAPPO to its recurrent variant (i.e., RMAPPO), as MAPPO performs equally

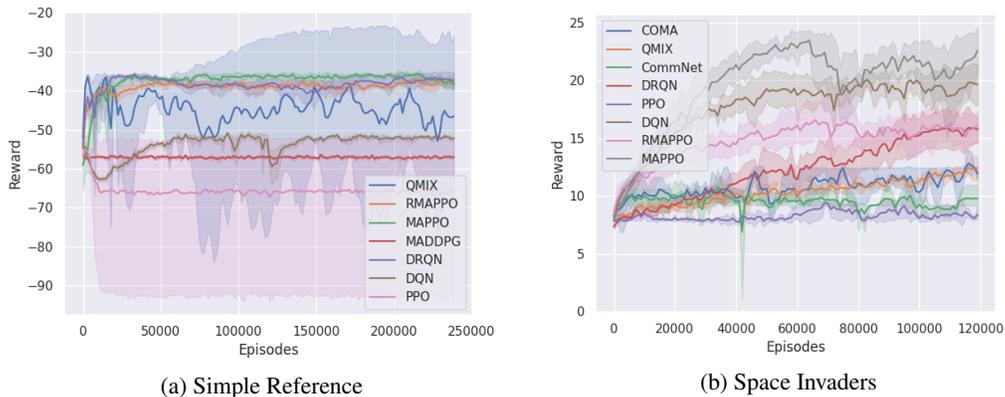


Figure 2: Training curves of various algorithms in two cooperative environments. For every algorithm, the solid line represents the mean reward per episode, while the shaded region represents the 95% confidence interval around the mean.

172 well as RMAPPO. We hypothesize that since MAPPO’s centralized critic learns based on the joint
 173 observation and action of both agents, this minimizes the amount of partial observability of every
 174 agent, and allows each agent to learn to communicate with other agents effectively without recurrence.
 175 In contrast, for independent algorithms, such as DQN, where the interactions between the agents are
 176 not explicitly learned (since all other agents are treated as part of the environment), adding recurrence
 177 could help mitigate some resulting partial observability, hence improving their performance, as
 178 described above.

179 **Space Invaders** Unlike the Simple Reference environment, the Space Invaders environment seemed
 180 to favour non-recurrent variants of algorithms (Fig. 2b). MAPPO vastly outperformed RMAPPO, and
 181 similarly DQN outperformed DRQN. This is also likely the underlying reasoning behind the compar-
 182 atively poorer performance of the multi-agent algorithms, such as QMIX, COMA and CommNet, all
 183 of which were implemented with recurrent neural networks under the CTDE scheme.

184 Additionally, since there is no unit collision in the Space Invaders environment (i.e., agents can move
 185 past each other without being blocked), they do not have to coordinate between themselves to achieve
 186 a high score in the environment; a good policy can be learned solely by having agents maximize their
 187 individual rewards. This explains the strong performance that was achieved by DQN. Also, since this
 188 is a cooperative task with both agents having identical goals, learning separate representations for
 189 individual agents is not very important; the learning of both agents assist each other. This is shown
 190 in Fig. 6b in Appendix B, where the addition of an agent indicator did not yield any performance
 191 improvement for DQN on Space Invaders.

192 Given such circumstances, it is interesting to observe the stronger performance of MAPPO com-
 193 pared to the independent algorithms. By conditioning on the joint action, MAPPO’s critic has full
 194 observability into the joint action that resulted in the team reward. Therefore, the observed reward is
 195 unbiased, which allows the learning process to be more efficient. In contrast, independent algorithms
 196 have to learn from a noisy team reward signal, where an agent could receive a large positive team
 197 reward even when it did nothing. This relates to the problem of credit assignment in MARL, noted in
 198 prior works [35].

199 4.1.1 Multi-Agent Credit Assignment Problem in Fully Observable Settings

200 In this section, we attempt to study the effect of using a team reward signal, rather than individual
 201 reward signals on various independent and multi-agent algorithms in a fully observable environment.
 202 When team rewards are the only rewards given, these reward signals are noisy for independent
 203 algorithms because the agent, which treats every other agent as part of the environment, does not
 204 know the actions taken by other agents. This makes it difficult for independent algorithms’ agents to
 205 learn how their individual actions contribute to the team reward signal. We performed the experiments
 206 on Space Invaders, in which the default agents receive individual rewards from the environment. To

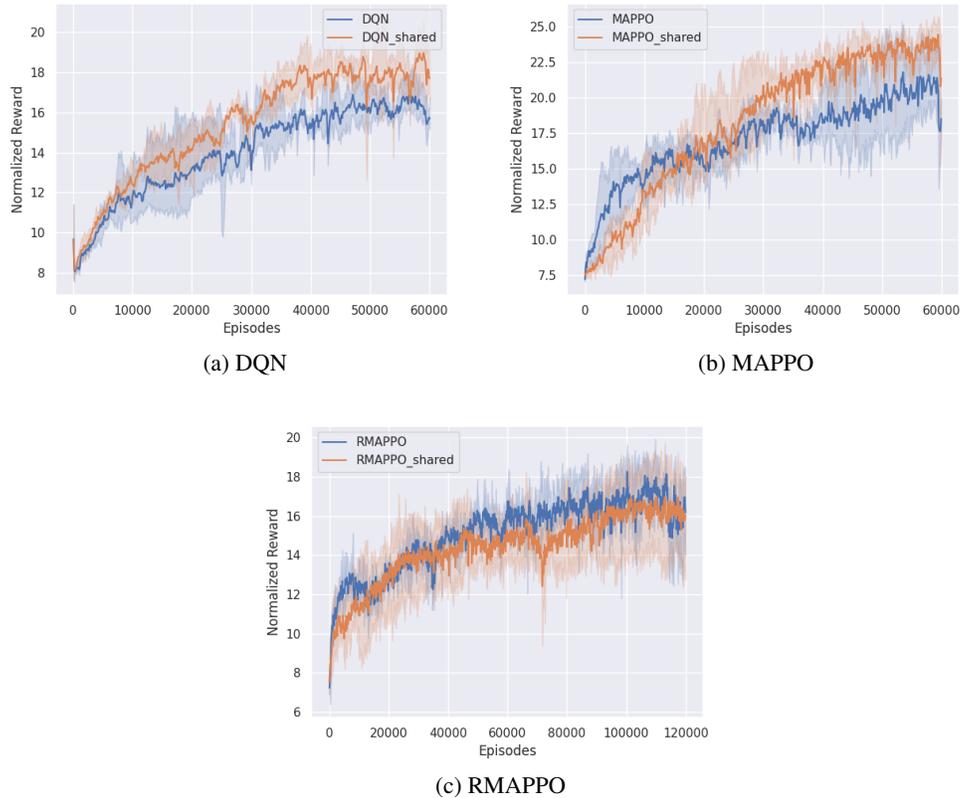


Figure 3: Training curves of various algorithms in Space Invaders, comparing when individual rewards are given (blue) to when team rewards are given (orange).

207 study the effect of the multi-agent credit assignment problem, we performed two runs per algorithm,
 208 one with team rewards only, and the other with individual rewards only (i.e., agents are rewarded
 209 independently by the environment).

210 For multi-agent algorithms, such as MAPPO (Fig. 3b) and RMAPPO (Fig. 3c), having a team reward
 211 does not have a large effect on the performance of the algorithms. This is expected because these
 212 algorithms have critics that learn from the joint action, which allow them to implicitly learn the
 213 estimated contribution of every agent without factorization.

214 On similar lines, regarding independent algorithms, we observe that having team rewards instead of
 215 individual ones do not impact their performance adversely (Fig. 3a). A plausible explanation could
 216 be that since all agents receive the same reward for a given joint action, this allows the independent
 217 algorithms to correlate actions from different agents that produced similar (high) rewards.

218 4.2 Competitive

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

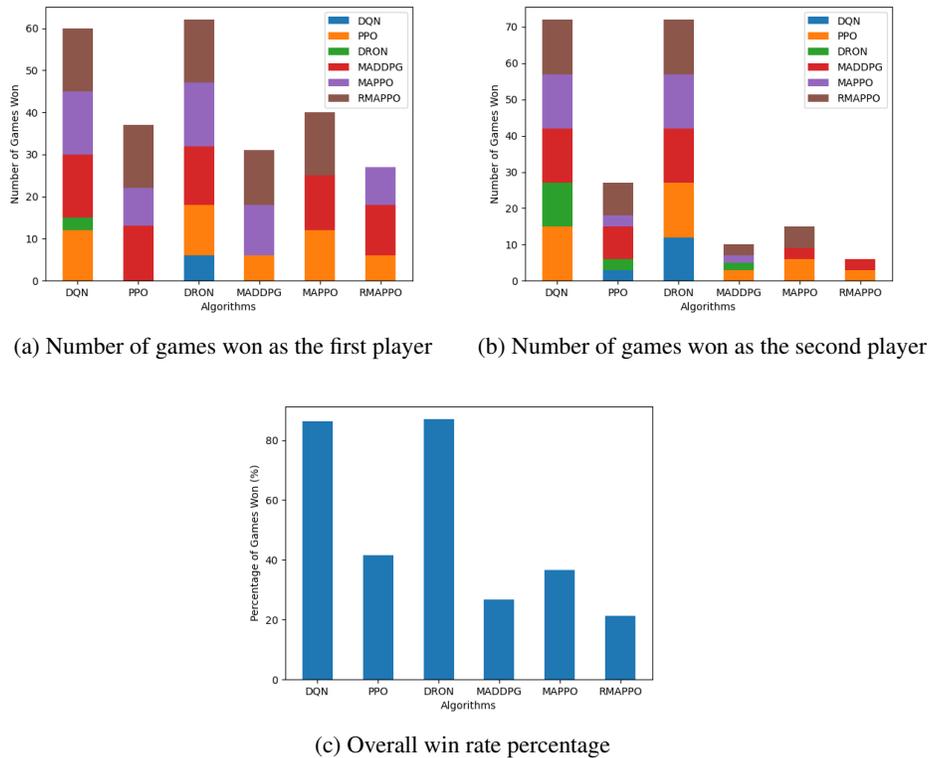


Figure 4: Performance of various algorithms when playing against other algorithms in Pong.

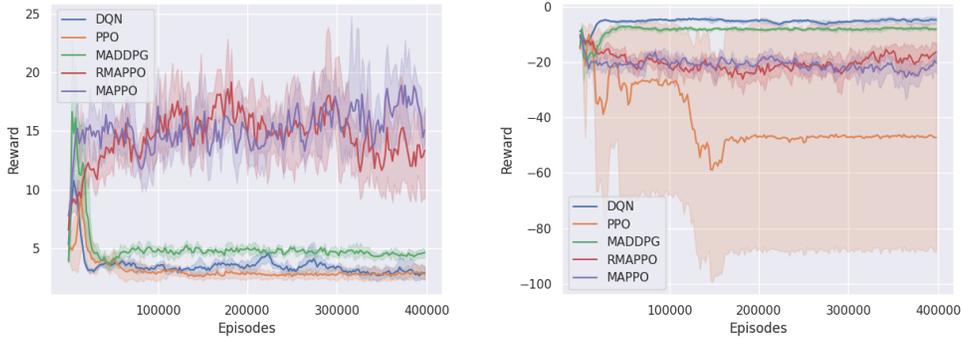
229 From the stacked bar charts shown in Fig. 4, a similar trend across the number of games won as
 230 the first and second player can be observed (Fig. 4a and 4b). DRON is consistently the best player,
 231 closely followed by DQN. Both of these algorithms were also the only algorithms to have a win rate
 232 of greater than 50% for the games they have played (Fig. 4c).

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

241 In addition to the above factors, it is possible that parameter sharing benefited agents of independent
 242 algorithms by allowing them to learn better representations of both players, since they were trained
 243 to play as both players simultaneously. Had these algorithms trained without parameter sharing,
 244 there would likely be a larger performance difference between independent algorithms and opponent
 245 modelling algorithms such as DRON. Instead of treating other agents as part of the environment,
 246 opponent modelling allows agents to adapt more quickly to the opponent’s changing strategies [24].
 247 However, the minimal improvement DRON has over DQN suggests that in the Pong environment,
 248 an agent’s policy may not be significantly affected by changes in the opponent agent’s policy (i.e.,
 249 individual agents can play the same way regardless of how their opponent played).

250 4.3 Mixed

251 In the Simple Tag (i.e., Predator-Prey) environment, the predators are incentivized to cooperate
 252 together to trap the prey, while the prey is incentivized to dodge the predators while staying within
 253 a predefined area. For our method of evaluation, we plot the training curves of the prey (Fig. 5b),



(a) Reward of a predator; all predators obtain the same reward

(b) Reward of the prey

Figure 5: Training curves of various algorithms in the Simple Tag, a Predator-Prey environment

254 and one of the predators (Fig. 5a), since all predators receive the same reward. Since the observation
 255 and action spaces differ between the predators and the prey, none of the agents have their parameters
 256 shared. We chose not to share the parameters of the predators to ensure that bias towards the predators
 257 was not introduced (since they would have 3 times the amount of data to work with compared to the
 258 prey).

259 In the case of DQN, the prey successfully learned to minimize the number of collisions with the
 260 predators, which can be observed by the strong performance achieved by the prey (Fig. 5b). However,
 261 similar to PPO, since the predators were trained completely independently (i.e., their parameters
 262 were not shared), they did not manage to learn how to cooperate with one another to capture the prey
 263 (Fig. 5a). It is interesting to observe that MADDPG converged to a policy similar to DQN, with
 264 the difference being that its predators have learned to cooperate better, thus getting slightly higher
 265 rewards compared to DQN’s predators (Fig. 5a). Subsequently, as a result of the higher rewards
 266 obtained by the predators, MADDPG achieves a slightly lower score for its prey (Fig. 5b).

267 MAPPO and RMAPPO, on the other hand, learned a different strategy. As we can observe from
 268 the comparatively noisier curves obtained from their predators and preys (Fig. 5a and 5b), there is
 269 a constant tug-of-war between the prey and the predators - as the predators learn how to cooperate
 270 better, their scores increase, which subsequently causes the prey to learn how to dodge, decreasing
 271 the predators’ scores, and vice versa. Since the predators of MAPPO and RMAPPO achieves a much
 272 higher score compared to all other algorithms, this is indicative that the predators have successfully
 273 learned to cooperate to trap the prey.

274 5 Conclusion

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

286 **Competitive** In the Pong environment, we saw that DRON and DQN were able to outperform all
 287 other algorithms. We argued that this is due to the fully observable nature of the Pong environment,
 288 in addition to the diminishing returns that learning from joint actions could yield. Furthermore, we

289 showed that with the use of agent indicators, independent algorithms were able to learn robust policies
290 for both competing agents using parameter sharing.

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

302 6 Future Work

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

314 References

- 315 [1] S. Choi, D.-Y. Yeung, and N. Zhang, “An environment model for nonstationary reinforcement learning,”
316 in *Advances in Neural Information Processing Systems*, S. Solla, T. Leen, and K. Müller, Eds., vol. 12,
317 MIT Press, 2000. [Online]. Available: [https://proceedings.neurips.cc/paper/1999/file/
318 e8d92f99edd25e2cef48eca48320a1a5-Paper.pdf](https://proceedings.neurips.cc/paper/1999/file/e8d92f99edd25e2cef48eca48320a1a5-Paper.pdf).
- 319 [2] M. Tan, “Multi-agent reinforcement learning: Independent vs. cooperative agents,” in *ICML*, 1993.
- 320 [3] E. Zawadzki, A. Lipson, and K. Leyton-Brown, *Empirically evaluating multiagent learning algorithms*,
321 2014. arXiv: 1401.8074 [cs.GT].
- 322 [4] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, “Multiagent
323 cooperation and competition with deep reinforcement learning,” *PLOS ONE*, vol. 12, no. 4, pp. 1–15,
324 Apr. 2017. DOI: 10.1371/journal.pone.0172395. [Online]. Available: [https://doi.org/10.
325 1371/journal.pone.0172395](https://doi.org/10.1371/journal.pone.0172395).
- 326 [5] Y. Shoham and K. Leyton-Brown, *Multiagent systems: Algorithmic, game-theoretic, and logical founda-
327 tions*. Cambridge University Press, 2008.
- 328 [6] OpenAI, : C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Dennison, D. Farhi, Q. Fischer,
329 S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P. d. O. Pinto,
330 J. Raiman, T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, and S. Zhang,
331 *Dota 2 with large scale deep reinforcement learning*, 2019. arXiv: 1912.06680 [cs.LG].
- 332 [7] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, “Counterfactual multi-agent policy
333 gradients,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, Apr. 2018.
334 [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/11794>.
- 335 [8] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, *Multi-agent actor-critic for mixed
336 cooperative-competitive environments*, 2020. arXiv: 1706.02275 [cs.LG].
- 337 [9] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, “QMIX: Mono-
338 tonic value function factorisation for deep multi-agent reinforcement learning,” in *Proceedings of the
339 35th International Conference on Machine Learning*, J. Dy and A. Krause, Eds., ser. Proceedings of
340 Machine Learning Research, vol. 80, PMLR, Jul. 2018, pp. 4295–4304. [Online]. Available: [http:
341 //proceedings.mlr.press/v80/rashid18a.html](http://proceedings.mlr.press/v80/rashid18a.html).

- 342 [10] L. Kraemer and B. Banerjee, “Multi-agent reinforcement learning as a rehearsal for decentralized
343 planning,” *Neurocomputing*, vol. 190, pp. 82–94, 2016, ISSN: 0925-2312. DOI: [https://doi.org/
344 10.1016/j.neucom.2016.01.031](https://doi.org/10.1016/j.neucom.2016.01.031). [Online]. Available: [https://www.sciencedirect.com/
345 science/article/pii/S0925231216000783](https://www.sciencedirect.com/science/article/pii/S0925231216000783).
- 346 [11] F. A. Oliehoek, M. T. J. Spaan, and N. A. Vlassis, “Optimal and approximate q-value functions for
347 decentralized pomdps,” *CoRR*, vol. abs/1111.0062, 2011. arXiv: 1111.0062. [Online]. Available:
348 <http://arxiv.org/abs/1111.0062>.
- 349 [12] J. K. Terry, B. Black, M. Jayakumar, A. Hari, R. Sullivan, L. Santos, C. Dieffendahl, N. L. Williams,
350 Y. Lokesh, C. Horsch, *et al.*, “Pettingzoo: Gym for multi-agent reinforcement learning,” *arXiv preprint
351 arXiv:2009.14471*, 2020.
- 352 [13] L. Busoniu, R. Babuska, and B. De Schutter, “A comprehensive survey of multiagent reinforcement
353 learning,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*,
354 vol. 38, no. 2, pp. 156–172, 2008. DOI: 10.1109/TSMCC.2007.913919.
- 355 [14] L. Canese, G. C. Cardarilli, L. Di Nunzio, R. Fazzolari, D. Giardino, M. Re, and S. Spanò, “Multi-agent
356 reinforcement learning: A review of challenges and applications,” *Applied Sciences*, vol. 11, no. 11, 2021,
357 ISSN: 2076-3417. DOI: 10.3390/app11114948. [Online]. Available: [https://www.mdpi.com/2076-
358 3417/11/11/4948](https://www.mdpi.com/2076-3417/11/11/4948).
- 359 [15] K. Zhang, Z. Yang, and T. Başar, *Multi-agent reinforcement learning: A selective overview of theories
360 and algorithms*, 2021. arXiv: 1911.10635 [cs.LG].
- 361 [16] S. Gronauer and K. Diepold, “Multi-agent deep reinforcement learning: A survey,” *Artificial Intelligence
362 Review*, pp. 1–49, 2021.
- 363 [17] C. Yu, A. Velu, E. Vinitzky, Y. Wang, A. Bayen, and Y. Wu, *The surprising effectiveness of mapo in
364 cooperative, multi-agent games*, 2021. arXiv: 2103.01955 [cs.LG].
- 365 [18] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- 366 [19] A. A. Markov, “The theory of algorithms,” *Trudy Matematicheskogo Instituta Imeni VA Steklova*, vol. 42,
367 pp. 3–375, 1954.
- 368 [20] R. BELLMAN, “A markovian decision process,” *Journal of Mathematics and Mechanics*, vol. 6, no. 5,
369 pp. 679–684, 1957, ISSN: 00959057, 19435274. [Online]. Available: [http://www.jstor.org/
370 stable/24900506](http://www.jstor.org/stable/24900506).
- 371 [21] L. S. Shapley, “Stochastic games,” *Proceedings of the National Academy of Sciences*, vol. 39, no. 10,
372 pp. 1095–1100, 1953, ISSN: 0027-8424. DOI: 10.1073/pnas.39.10.1095. eprint: [https://www.
373 pnas.org/content/39/10/1095.full.pdf](https://www.pnas.org/content/39/10/1095.full.pdf). [Online]. Available: [https://www.pnas.org/
374 content/39/10/1095](https://www.pnas.org/content/39/10/1095).
- 375 [22] Y.-h. Chang, T. Ho, and L. Kaelbling, “All learning is local: Multi-agent learning in global reward games,”
376 in *Advances in Neural Information Processing Systems*, S. Thrun, L. Saul, and B. Schölkopf, Eds., vol. 16,
377 MIT Press, 2004. [Online]. Available: [https://proceedings.neurips.cc/paper/2003/file/
378 c8067ad1937f728f51288b3eb986afaa-Paper.pdf](https://proceedings.neurips.cc/paper/2003/file/c8067ad1937f728f51288b3eb986afaa-Paper.pdf).
- 379 [23] S. Sukhbaatar, A. Szlam, and R. Fergus, “Learning multiagent communication with backpropaga-
380 tion,” in *Proceedings of the 30th International Conference on Neural Information Processing Systems*,
381 ser. NIPS’16, Barcelona, Spain: Curran Associates Inc., 2016, pp. 2252–2260, ISBN: 9781510838819.
- 382 [24] H. He, J. Boyd-Graber, K. Kwok, and H. Daumé, “Opponent modeling in deep reinforcement learning,”
383 in *Proceedings of the 33rd International Conference on International Conference on Machine Learning -
384 Volume 48*, ser. ICML’16, New York, NY, USA: JMLR.org, 2016, pp. 1804–1813.
- 385 [25] G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, “Dealing with non-stationarity in multi-
386 agent deep reinforcement learning,” *CoRR*, vol. abs/1906.04737, 2019. arXiv: 1906.04737. [Online].
387 Available: <http://arxiv.org/abs/1906.04737>.
- 388 [26] P. Hernandez-Leal, M. Kaisers, T. Baarslag, and E. M. de Cote, *A survey of learning in multiagent
389 environments: Dealing with non-stationarity*, 2019. arXiv: 1707.09183 [cs.MA].
- 390 [27] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller,
391 A. K. Fidjeland, G. Ostrovski, *et al.*, “Human-level control through deep reinforcement learning,” *nature*,
392 vol. 518, no. 7540, pp. 529–533, 2015.
- 393 [28] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algo-
394 rithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- 395 [29] M. Hausknecht and P. Stone, *Deep recurrent q-learning for partially observable mdps*, 2017. arXiv:
396 1507.06527 [cs.LG].
- 397 [30] J. K. Gupta, M. Egorov, and M. Kochenderfer, “Cooperative multi-agent control using deep reinforcement
398 learning,” in *Autonomous Agents and Multiagent Systems*, G. Sukthankar and J. A. Rodriguez-Aguilar,
399 Eds., Cham: Springer International Publishing, 2017, pp. 66–83, ISBN: 978-3-319-71682-4.
- 400 [31] I. Mordatch and P. Abbeel, “Emergence of grounded compositional language in multi-agent populations,”
401 *arXiv preprint arXiv:1703.04908*, 2017.

- 402 [32] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling, "The arcade learning environment: An evaluation
403 platform for general agents," *Journal of Artificial Intelligence Research*, vol. 47, pp. 253–279, Jun. 2013.
- 404 [33] J. K. Terry and B. Black, "Multiplayer support for the arcade learning environment," *arXiv preprint*
405 *arXiv:2009.09341*, 2020.
- 406 [34] M. C. Machado, M. G. Bellemare, E. Talvitie, J. Veness, M. J. Hausknecht, and M. Bowling, "Revisiting
407 the arcade learning environment: Evaluation protocols and open problems for general agents," *CoRR*,
408 vol. abs/1709.06009, 2017. arXiv: 1709.06009. [Online]. Available: <http://arxiv.org/abs/1709.06009>.
- 409
- 410 [35] P. Hernandez-Leal, B. Kartal, and M. E. Taylor, "A survey and critique of multiagent deep reinforcement
411 learning," *Journal of Autonomous Agents and Multi-Agent Systems (JAAMAS)*, vol. 33, no. 6, pp. 750–797,
412 2019.
- 413 [36] J. K. Terry, B. Black, and A. Hari, "Supersuit: Simple microwrappers for reinforcement learning
414 environments," *arXiv preprint arXiv:2008.08932*, 2020.
- 415 [37] M. Li, *Machin*, <https://github.com/iffix/machin>, 2020.
- 416 [38] A. Raffin, A. Hill, M. Ernestus, A. Gleave, A. Kanervisto, and N. Dormann, *Stable baselines3*, <https://github.com/DLR-RM/stable-baselines3>, 2019.
- 417
- 418 [39] starry-sky6688, *Starcraft*, <https://github.com/starry-sky6688/StarCraft>, 2019.
- 419 [40] H. van Hasselt, A. Guez, and D. Silver, *Deep reinforcement learning with double q-learning*, 2015. arXiv:
420 1509.06461 [cs.LG].
- 421 [41] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas, *Dueling network architec-*
422 *tures for deep reinforcement learning*, 2016. arXiv: 1511.06581 [cs.LG].
- 423 [42] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, *Prioritized experience replay*, 2016. arXiv: 1511.05952
424 [cs.LG].