

An Extended Benchmarking of Multi-Agent Reinforcement Learning Algorithms in Complex Fully Cooperative Tasks

George Papadopoulos*
 University of Piraeus
 Piraeus, Greece
 georgepap@unipi.gr

Chaido Poulianou
 University of Piraeus
 Piraeus, Greece
 cpoulianou@uth.gr

Andreas Kontogiannis*
 NTUA & Archimedes AI
 Athens, Greece
 andreas.kontogiannis@mail.ntua.gr

Ioannis Koumentis
 University of Piraeus
 Piraeus, Greece
 iokoumen@unipi.gr

Foteini Papadopoulou
 Radboud University
 Nijmegen, Netherlands
 foteini.papadopoulou@ru.nl

George Vouros
 University of Piraeus
 Piraeus, Greece
 georgev@unipi.gr

ABSTRACT

Multi-Agent Reinforcement Learning (MARL) has recently emerged as a significant area of research. However, MARL evaluation often lacks systematic diversity, hindering a comprehensive understanding of algorithms' capabilities. In particular, cooperative MARL algorithms are predominantly evaluated on benchmarks such as SMAC and GRF, which primarily feature team game scenarios without assessing adequately various aspects of agents' capabilities required in fully cooperative real-world tasks such as multi-robot cooperation and warehouse, resource management, search and rescue, and human-AI cooperation. Moreover, MARL algorithms are mainly evaluated on low dimensional state spaces, and thus their performance on high-dimensional (e.g., image) observations is not well-studied. To fill this gap, this paper highlights the crucial need for expanding systematic evaluation across a wider array of existing benchmarks. To this end, we conduct extensive evaluation and comparisons of well-known MARL algorithms on complex *fully* cooperative benchmarks, including tasks with *images* as agents' observations. Interestingly, our analysis shows that many algorithms, hailed as state-of-the-art on SMAC and GRF, may underperform standard MARL baselines on fully cooperative benchmarks. Finally, towards more systematic and better evaluation of cooperative MARL algorithms, we have open-sourced PyMARLzoo+, an extension of the widely used (E)PyMARL libraries, which addresses an open challenge from [49], facilitating seamless integration and support with all benchmarks of PettingZoo, as well as Overcooked, PressurePlate, Capture Target and Box Pushing.

KEYWORDS

Fully Cooperative Multi-Agent Reinforcement Learning, Benchmarking, Image-based Observations, Open-Source Framework

ACM Reference Format:

George Papadopoulos, Andreas Kontogiannis, Foteini Papadopoulou, Chaido Poulianou, Ioannis Koumentis, and George Vouros. 2025. An Extended

*Equal Contribution.



This work is licensed under a Creative Commons Attribution International 4.0 License.

Benchmarking of Multi-Agent Reinforcement Learning Algorithms in Complex Fully Cooperative Tasks. In *Proc. of the 24th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2025)*, Detroit, Michigan, USA, May 19 – 23, 2025, IFAAMAS, 31 pages.

1 INTRODUCTION

In *fully* cooperative Multi-Agent Reinforcement Learning (MARL) problems, the goal is to train learnable agents in order to maximize their shared cumulative reward, through excessive coordination, sharing of tasks, collaborative exploration, with appropriate decisions for timing and action, and sharing of capabilities. Fully cooperative MARL is of remarkable interest, as it can naturally model many real-world applications, including multi-robot collaboration [6] and warehouse [36], search and rescue [39], human-AI coordination [7], air traffic management [25], logistics networks [29], and supply-chain optimization [23]. Recently, cooperative MARL algorithms are mainly evaluated in settings where adversarial non-learnable agents that interact with the learnable cooperative ones exist: This has received a surge of approaches and methodologies, e.g., see [21, 31, 48, 61], also drawing motivation from multiplayer video-games and team sports.

Despite recent efforts [4, 20, 36, 63] aiming to provide a comprehensive understanding of standard cooperative MARL algorithms' capabilities through benchmarking, MARL evaluation still lacks systematic diversity and reliability for the following reasons:

- Most state-of-the-art (SoTA) cooperative MARL algorithms are predominantly evaluated, and possibly overfit, as pointed out in [16], on specific cooperative-competitive benchmarks where a team of cooperative learning agents competes against a team of bots with fixed policies, namely SMAC [11, 43] and GRF [27]. However, we argue that these benchmarks do not allow adequate evaluation of subtle issues involved in *fully* cooperative MARL, including: excessive coordination and exploration capabilities, sharing of capabilities, appropriate timing in the execution of actions, complimentary observability, scaling to large numbers of agents, and possibly with sparse rewards. For instance, tasks from the LBF benchmark [36], and particularly those with large grids, effectively capture more diverse requirements of fully cooperative multi-robot collaboration: They require agents to first engage in *extensive joint exploration* to identify a *certain* food target, followed by coordinated joint actions where all

agents must *simultaneously* consume that target. We conjecture that these aspects are crucial for fully cooperative, real-world tasks. On the other hand, these benchmarks emphasize developing skills that do not play dominant roles for fully cooperative tasks, such as countering the opponent team (e.g., surviving enemy attacks in SMAC or tackling opponents and preventing goals in GRF).

- Most MARL algorithms are evaluated solely on tasks with low-dimensional (mostly tabular) state spaces, and thus their effectiveness on real-world, high-dimensional, image-based observations has not been studied.
- MARL evaluation does not often report the training times of the proposed algorithms, so the results cannot be interpreted as a function of the compute budget used [16].

To address the above challenges, the main contributions of this paper are the following: (1) Our paper investigates the effectiveness of established MARL algorithms and contributes a comprehensive, **updated empirical evaluation and comparison** of MARL algorithms, including algorithms that have demonstrated SoTA performance in SMAC and GRF, across a wide array of complex **fully cooperative** benchmarks. (2) To our knowledge, our work is the first to include tasks with *images* representing *high-dimensional observations* in MARL benchmarking. (3) We contribute an open-source Python MARL framework, namely *PyMARLzoo+*¹, which extends the (E)PyMARL frameworks [36, 43] (widely used in developing established MARL algorithms, such as [21, 28, 40, 53, 54]), facilitating seamless integration with all PettingZoo tasks, thus addressing an open challenge from [49] (EPyMARL supports only the MPE PettingZoo tasks, which were already integrated in [36]). In addition to the already integrated MPE [32, 33], LBF [3] and RWARE [36] benchmarks, our PyMARLzoo+ also integrates the complex fully cooperative Overcooked [7], PressurePlate [1, 17], Capture Target [35, 57] and BoxPushing [58, 59] benchmarks. (4) Our work provides benchmarking of tasks that are indicative of a wide array of real-world applications, and of the evaluation of diverse requirements for fully cooperative MARL, in terms of joint exploration and coordination. (5) Our experimental findings demonstrate that many algorithms, which have been SoTA in SMAC and GRF, may underperform standard MARL algorithms in fully cooperative MARL tasks; thus validating the possible overfitting issues of the current MARL evaluation highlighted in [16]. Furthermore, we point out fully cooperative tasks that are very hard to solve with the existing SoTA methods. (6) Our benchmarking is the first to report the training times of algorithms as a further measure of the algorithm’s performance, so that the reported results are also interpreted as a function of the compute budget used.

2 PRELIMINARIES

2.1 Dec-POMDPs

Fully cooperative MARL is formulated as a Dec-POMDP [34]. A Dec-POMDP for an N -agent task is a tuple $\langle S, A, P, r, F, O, N, \gamma \rangle$, where S is the state space, A is the joint action space $A = A_1 \times \dots \times A_N$, $P(s' | s, a) : S \times A \rightarrow [0, 1]$ is the state transition function, $r(s, a) : S \times A \rightarrow \mathbb{R}$ is the *shared* reward function and $\gamma \in [0, 1)$ is the discount factor.

¹The source code is available at: <https://github.com/AllLabDsUnipi/pymarlzooplus>

Assuming partial observability, each agent at time step t does not have access to the full state, yet it samples observations $o_t^i \in O_i$ according to the observation function $F_i(s) : S \rightarrow O_i$. The joint observations are denoted by $o \in O$ and are sampled according to $F = \prod_i F_i$. The action-observation history for agent i at time t is denoted by $h_t^i \in H_i$, which includes action-observation pairs until $t-1$ and o_t^i , on which the agent can condition its individual stochastic policy $\pi_{\theta_i}^i(a_t^i | h_t^i) : H_i \times A_i \rightarrow [0, 1]$, parameterised by θ_i . The joint policy is denoted by π_θ , with parameters $\theta \in \Theta$. The objective is to find an optimal joint policy which satisfies the optimal value function $V^*(s) = \max_\theta \mathbb{E}_{a \sim \pi_\theta, s' \sim P(\cdot | s, a), o \sim F(s)} [\sum_{t=0}^{\infty} \gamma^t r_t]$.

2.2 Assumptions of interest

In addition to *partial observability*, in our benchmark analysis we consider and study the following assumptions of interest: (a) we utilize the celebrated *CTDE* MARL schema [32], which has been widely adopted by the cooperative MARL community [2, 16] as it enables conditioning approximate value functions on privileged information in a computationally tractable manner, (b) we evaluate *fully cooperative tasks*, that is, tasks where there are no adversarial non-learnable agents (e.g., bots), or teams of opponent agents, interacting with the learnable agents, (c) we also evaluate complex tasks with *sparse-reward settings* that require excessive joint exploration and cooperation to be solved, (d) we also study tasks with *image-based, high-dimensional state and observation spaces* (as in real-world tasks), which have been very frequent in single-agent RL but not in MARL evaluation, and (e) as in [36, 63] we assume that agents *lack explicit communication channels during execution*.

3 ALGORITHMS

In our benchmark analysis, we evaluate and compare a wide range of well-known MARL algorithms (see Table 1). The selection of these algorithms is based on the following important (but not all-encompassing) criteria: (1) they have been widely used as competitive baselines by the MARL community in recent works, (2) they have achieved SoTA performance in cooperative-competitive MARL benchmarks, such as SMAC and GRF, and in our analysis we aim to test their performance in fully cooperative tasks, (3) they are based on improving joint exploration in sparse-reward tasks, and (4) they present diversity in the RL optimization part (e.g., being value-based or actor-critic based, utilizing standard recurrent architectures or transformer-based).

4 FULLY COOPERATIVE MARL BENCHMARKS

In this section, we highlight complex, partial observable, fully cooperative tasks from eight MARL benchmarks that we believe to be of significant interest for MARL research. These tasks represent a wide range of real-world applications and test diverse requirements for fully cooperative MARL (see also Table 2).

4.1 PettingZoo

PettingZoo [49] is a Python library consisting of several MARL tasks. In our analysis, we utilize the following fully cooperative PettingZoo benchmarks: *Entombed Cooperative*, along with *Pistonball* and *Cooperative Pong*. These tasks allow challenging scenarios that provide useful testbeds for fully cooperative MARL, using

Algorithm	Evaluated in	On/Off-Policy	RL Optimization	Network Architectures	Intrinsic Exploration
QMIX [40]	SMAC [43], GRF [27], MPE [33], LBF [3], RWARE [36]	Off-policy	Value-based	RNN & MLP	No
MAA2C [36]	SMAC [43], GRF [27], MPE [33], LBF [3], RWARE [36]	On-policy	Actor-Critic	RNN & MLP	No
COMA [12]	SMAC [43], MPE [33], LBF [3], RWARE [36]	On-policy	Policy Gradient	RNN & MLP	No
MAPPO [62]	SMAC [43], GRF [27], MPE [33], LBF [3], RWARE [36]	On-policy	Actor-Critic	RNN & MLP	No
QPLEX [53]	SMAC [43]	Off-policy	Value-based	RNN & MLP	No
HAPPO [26]	SMAC [43], MA MuJoCo [37]	On-policy	Actor-Critic	RNN & MLP	No
MAT-DEC [55]	SMAC [43], GRF [27], MA MuJoCo [37]	On-policy	Actor-Critic	Transformer & MLP	No
EMC [64]	SMAC [43]	Off-policy	On top of QPLEX	RNN & MLP	<i>Yes:</i> curiosity-driven
MASER [22]	SMAC [43]	Off-policy	On top of QMIX	RNN & MLP	<i>Yes:</i> subgoal generation
EOI [24]	GRF [27], MAgent [65]	Off-policy	On top of MAA2C	RNN & MLP	<i>Yes:</i> individuality
CDS [28]	SMAC [43], GRF [27]	Off-policy	On top of QPLEX	RNN & MLP	<i>Yes:</i> diversity & information sharing

Table 1: Summary of the selected MARL algorithms

high-dimensional, RGB-image-based observation spaces. Indicative references to the PettingZoo benchmark include [8, 47, 49–51].

4.1.1 Entombed Cooperative. Here, agents must extensively coordinate to progress as far as possible into a procedurally generated maze. Each agent needs to quickly navigate down a constantly evolving maze, where only part of the environment is visible. If an agent becomes trapped, they lose. Agents can easily find themselves in dead-ends, only escapable through rare power-ups. A major challenge is that optimal coordination requires agents to position themselves on opposite sides of the map, as power-ups appear on one side or the other but can be used to break through walls symmetrically.

4.1.2 Pistonball. Pistonball is a physics-based fully cooperative game in which the goal is to move a ball to the left boundary of the game area by controlling a set of vertically moving pistons. The main challenge lies in achieving highly coordinated, emergent behavior to optimize performance in the environment. Pistonball uses a realistic physics engine, comparable to the game Angry Birds, adding further complexity to the required agent coordination.

4.1.3 Cooperative Pong. Cooperative Pong is a fully cooperative variant of the classic Pong game where two agents control paddles on opposite sides of the screen, aiming to keep the ball in play. The challenge lies in the asymmetry between the agents—particularly the right paddle, which has a more complex tiered shape—and their limited observation space, restricted to their own half of the screen. This setup requires agents to coordinate excessively and develop cooperative strategies to maximize play time.

4.2 Overcooked

Overcooked [7] is a fully cooperative MARL benchmark, which has been developed to address human-AI coordination. The objective is to deliver soups as quickly as possible, with agents required to place up to three ingredients in a pot, wait for the soup to cook, and then deliver it. The tight space introduces significant challenges in terms of coordination, as agents must split tasks on the fly, avoid collisions, and coordinate effectively in order to achieve high reward. In addition, the rewards are sparse, making the environment even more challenging. Indicative references that use this benchmark include [8, 47, 49]. We utilize the following three different layouts.

4.2.1 Cramped Room. Agents operate in a cramped room, that is, a small room without obstacles. This layout focuses on the agents’ ability to maneuver in a limited space, requiring precise navigation to avoid physical collisions with the other agent.

4.2.2 Asymmetric Advantages. Agents operate in an asymmetric room, where each agent works in its own distinct space. In this layout, the challenge lies in recognizing and exploiting the differing strengths of each agent, such as speed or access to certain kitchen resources. This task tests the agents’ ability to adapt their strategies to complement each other’s capabilities, ensuring that each player’s strengths are used effectively to enhance overall team performance.

4.2.3 Coordination Ring. Agents operate in a room with a ring. They must coordinate their actions to collect onions from the bottom left corner of the room, make soups in the center left, and deliver the dishes to the top right corner of the ring. This task is the most challenging among the three.

MARL Benchmark	Interesting Challenges	Insights for Real-world Applications
<i>Entombed Cooperative</i> (PettingZoo)	<ul style="list-style-type: none"> - RGB observations - dead-ends - excessive coordination to go to opposite sides in order to break through walls symmetrically 	<ul style="list-style-type: none"> - exploration in space missions - search & rescue - collective wildfire movement
<i>Pistonball</i> (PettingZoo)	<ul style="list-style-type: none"> - RGB observations - excessive coordination for emergent behavior and synchronization 	<ul style="list-style-type: none"> - assembly line coordination - dance and performance choreography
<i>Cooperative Pong</i> (PettingZoo)	<ul style="list-style-type: none"> - RGB observations - excessive coordination due to asymmetry between the agents 	<ul style="list-style-type: none"> - collaborative video games - performing arts - sports collaboration
<i>Overcooked</i> (<i>Cramped Room, Asymmetric Advantages, Coordination Ring</i>)	<ul style="list-style-type: none"> - sparse reward - delivering tasks as quickly as possible - agents operating in a confined space - agent collision avoidance - splitting tasks on the fly 	<ul style="list-style-type: none"> - kitchen automation - urgent multi-robot tasks in confined spaces
Pressure Plate	<ul style="list-style-type: none"> - sparse reward - excessive exploration and coordination to sequentially unlock each room 	<ul style="list-style-type: none"> - multi-robot collaboration - escape-rooms-like tasks - team-based problem solving
Spread (MPE)	<ul style="list-style-type: none"> - agent collision avoidance - optimal landmark coverage - very complex if N is large 	<ul style="list-style-type: none"> - traffic management - distributed sensor networks - urban planning
LBF	<ul style="list-style-type: none"> - sparse reward - exploration to identify the food target - excessive coordination to eat the target simultaneously 	<ul style="list-style-type: none"> - multi-robot collaboration - resource management in supply chains - disaster response coordination
RWARE	<ul style="list-style-type: none"> - sparse reward and high-dimensional observations - excessive coordination to execute a specific sequence of actions, at the right time, without immediate feedback 	<ul style="list-style-type: none"> - robot warehousing - logistics management

Table 2: Summary of the fully cooperative, partially observable, benchmark tasks

4.3 PressurePlate

PressurePlate [1] is a fully cooperative environment, with sparse-reward settings, set within a 2D grid-world composed of multiple locked rooms that can be unlocked when an agent stands on the corresponding pressure plate, culminating in a final room containing a goal chest. The primary challenge for agents lies in effectively coordinating their movements and positions to sequentially unlock each room and ultimately reach the chest. Indicative references that use this benchmark task include [1, 17].

4.4 Multi-agent Particle Environment (MPE)

In our analysis, we utilize the well-known *Spread* task of the MPE benchmark [32] that focus on effective navigation of particle agents. Indicative references that use this benchmark task include [10, 36, 41, 66]. In this task, N agents must cover N landmarks while avoiding collisions. Agents are rewarded for staying close to landmarks and penalized for collisions, creating a trade-off between coordination and collision avoidance. The task becomes more complex as N increases. Compared to [36], the evaluated tasks are more challenging by using more agents and fewer training steps.

4.5 Level Based Foraging (LBF)

Level-Based Foraging (LBF) [3] features fully cooperative grid-world environments where agents, assigned levels, must move and

Benchmark	Number of Newly Integrated Tasks
PettingZoo	49
Overcooked	5
Pressure Plate	3
Capture Target	1
Box Pushing	1

Table 3: Newly integrated tasks (in addition to MPE, LBF and RWARE) in PyMARLzoo++.

collect food by coordinating their actions. Agents can collect food only if their combined levels meet or exceed the food's level. The main challenge is the sparse rewards, requiring agents to coordinate closely to collect food simultaneously. Unlike previous work [36], we evaluate more complex LBF tasks, even with a larger number of agents, requiring extensive exploration and coordination. Indicative references that use this benchmark include [7, 14, 19, 60].

4.6 Multi-Robot Warehouse (RWARE)

The Multi-Robot Warehouse (RWARE) environment simulates a fully cooperative, partially observable grid-world where agents must find and deliver shelves to workstations. With limited sight

and partial observations, agents face challenges such as sparse rewards, only given after successful shelf deliveries. This requires precise action sequences, effective exploration, and excessive agent coordination. Unlike previous work [36], this study evaluates more challenging RWARE tasks in hard mode, demanding efficient exploration and excessive cooperation among agents. Indicative references that use RWARE include [49–51].

4.7 Capture Target

Capture Target [35, 57], a fully cooperative multi-agent single-target task, presents significant challenges, aiming multiple agents to locate and capture a flickering target on a grid. Excessive coordination is essential, as the target is only captured when all agents converge on its location simultaneously, despite limited visibility.

4.8 Box Pushing

Box Pushing [58, 59] is a fully cooperative grid environment where two agents must collaborate to move three boxes—two small and one large—to a target area. The main challenge is that the large box, yielding the highest reward, can only be moved if both agents coordinate by positioning themselves in parallel cells and pushing simultaneously, making precise timing and teamwork essential.

5 PYMARLZOO+: PETTINGZOO, OVERCOOKED, PRESSURE PLATE, CAPTURE TARGET AND BOX PUSHING NOW COMPATIBLE WITH (E)PYMARL

To facilitate a comprehensive understanding of MARL algorithms through benchmarking, we open-source *PyMARLzoo+*, a Python framework which extends the widely used (E)PyMARL [36, 43], providing an integration of many SoTA algorithms on a plethora of existing MARL benchmarks under a common framework. The key features of our framework are presented below:

Newly Integrated Benchmarks. Our PyMARLzoo+ integrates 8 MARL benchmarks described in Section 4. More specifically, as we illustrate in Table 3, in addition to the MPE, LBF and RWARE benchmarks (already integrated in EPyMARL [36]), our framework fully supports all tasks from PettingZoo, along with tasks from Overcooked, Pressure Plate, Capture Target, and Box Pushing.

Newly Integrated Algorithms. Our PyMARLzoo+ integrates all the algorithms described in Section 3. We note that, except for the standard baselines (that is, MAA2C, MAPPO and QMIX) already integrated in EPyMARL, all remaining algorithms were integrated as part of our work. Furthermore, the widely used *Prioritized Replay Buffer* [44] has been integrated into the off-policy algorithms.

Image (Observation) encoding. We incorporate three pre-trained architectures as image encoders to transform image-based observations into tabular format for policy learning. Specifically, we integrate the following options: *ResNet18* [18], *CLIP* [38], and *SlimSAM* [9]. *ResNet18*, commonly used in single-agent RL (e.g., in [46]), can capture spatial hierarchies, which are beneficial for extracting relevant features from complex visual data. *CLIP*, also used in single agent RL (e.g., in [15]), utilizes natural language supervision to

produce robust and semantically meaningful representations, facilitating model generalization across diverse tasks by linking visual inputs with textual descriptions. Similarly, *SlimSAM* incorporates a streamlined self-attention mechanism for efficient image encoding. In addition, we offer the option of using *standard CNNs* on raw images for policy training from scratch.

6 BENCHMARK ANALYSIS

6.1 Experimental Setup

In our benchmark analysis, we evaluate MARL algorithms on *Entombed Cooperative*, *Pistonball*, *Cooperative Pong*, *Cramped Room*, *Asymmetric Advantages*, *Coordination Ring*, *Pressure Plate*, *Spread*, *LBF* and *RWARE* benchmarks. To ensure a fair comparison of the selected algorithms, we utilize the same number of training timesteps. Specifically, we use 10 million timesteps for MPE and LBF, 40 million for RWARE, 5 million for PettingZoo; 20 million for PressurePlate, 40 million for the Overcooked’s *Cramped Room*, and 100 million timesteps for the Overcooked’s *Asymmetric Advantages* and *Coordination Ring*. We note that for all PettingZoo tasks, where the observations are RGB images, we have used a pre-trained ResNet18 model as the observation image encoder. Moreover, except HAPPO, we utilize agents to share their policy parameters.

We adopt the experimental setup of [36]: Throughout the training of all algorithms, we conducted 100 test episodes at every 50,000-step interval to evaluate the current policy. During these test episodes, we record the episode returns, i.e., the accumulated rewards per episode, and compute the average return $\frac{1}{N} \sum_{i=1}^N R_{t,i,j}$, where $N = 100$ is the number of test episodes, and $R_{t,i,j}$ is the return of the i -th test episode at timestep t of the j -th experiment.

Our primary metric score is the mean of the unnormalized returns received from the best policy in convergence over five different seeds. As best policy in convergence, we use the policy that received the best return of the last 50 test episodes. We present our results using the mean and the 75% confidence interval.

We used the hyperparameter settings of the newly integrated algorithms, QPLEX, HAPPO, MAT-DEC, CDS, EOI, EMC, and MASER, adhering to configurations suggested by the authors of their original papers for challenging tasks. For MAA2C, MAPPO, COMA and QMIX, we employed the default parameters of EPyMARL. This ensures that the evaluation is as consistent as possible with what is reported in the original papers, allowing for a valid comparison between the algorithms and with results already reported. Our choice to forego extensive tuning is further supported by: (a) preliminary results showing *no significant differences* regarding algorithms’ performance across tasks, and (b) the fact that our results, are based on the mean returns from the *best policy* across different seeds.

All experiments were conducted using CPUs, except for tasks within the PettingZoo environment, where image inputs required the use of GPUs. For these image-based tasks, GPUs were utilized for both the frozen image-encoder and training CNN components when a pre-trained encoder was not available. Specifically, we employed two “g5.16xlarge” AWS instances, each featuring 64 vCPUs, 256 GB RAM, and one A10G GPU with 24 GB of memory to manage the computational demands of image processing.

Tasks\Algorithms	QMIX	QPLEX	MAA2C	MAPPO	HAPPO	MAT-DEC	COMA	EOI	MASER	EMC	CDS
2s-8x8-3p-2f	0.94 ± 0.09	0.63 ± 0.48	0.98 ± 0.02	0.64 ± 0.34	0.00 ± 0.00	0.24 ± 0.22	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
2s-9x9-3p-2f	0.00 ± 0.00	0.60 ± 0.49	0.59 ± 0.38	0.18 ± 0.35	0.00 ± 0.00	0.34 ± 0.24	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.50 ± 0.50	0.00 ± 0.00
2s-12x12-2p-2f	0.89 ± 0.01	0.98 ± 0.00	0.81 ± 0.03	0.77 ± 0.04	0.52 ± 0.26	0.55 ± 0.06	0.03 ± 0.02	0.34 ± 0.23	0.01 ± 0.01	0.87 ± 0.03	0.91 ± 0.02
4s-11x11-3p-2f	0.08 ± 0.19	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
7s-20x20-5p-3f	0.01 ± 0.00	0.00 ± 0.00	0.78 ± 0.02	0.57 ± 0.18	0.00 ± 0.00	0.29 ± 0.09	0.03 ± 0.01	0.03 ± 0.01	0.01 ± 0.00	0.01 ± 0.00	0.00 ± 0.00
8s-25x25-8p-5f	0.02 ± 0.01	0.01 ± 0.00	0.52 ± 0.24	0.41 ± 0.23	0.00 ± 0.00	0.03 ± 0.03	0.03 ± 0.00	0.07 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
7s-30x30-7p-4f	0.06 ± 0.05	0.00 ± 0.00	0.71 ± 0.02	0.57 ± 0.03	0.00 ± 0.00	0.08 ± 0.06	0.02 ± 0.00	0.04 ± 0.00	0.01 ± 0.00	0.00 ± 0.00	0.00 ± 0.00

Table 4: Results in LBF tasks.

Tasks\Algorithms	QMIX	QPLEX	MAA2C	MAPPO	HAPPO	MAT-DEC	COMA	EOI	MASER	EMC	CDS
tiny-2ag-hard	0.00 ± 0.00	0.61 ± 0.45	2.25 ± 0.62	2.91 ± 0.82	1.46 ± 2.06	6.00 ± 8.49	0.02 ± 0.00	6.58 ± 3.92	0.00 ± 0.00	1.66 ± 0.30	4.00 ± 2.30
tiny-4ag-hard	0.00 ± 0.00	11.73 ± 10.81	6.87 ± 7.12	17.15 ± 5.31	24.07 ± 0.71	27.85 ± 19.71	0.02 ± 0.00	12.09 ± 4.59	0.00 ± 0.00	0.00 ± 0.00	27.37 ± 6.88
small-4ag-hard	0.01 ± 0.01	0.94 ± 0.63	1.38 ± 1.26	4.14 ± 2.12	9.69 ± 3.19	0.00 ± 0.00	0.03 ± 0.00	0.17 ± 0.01	0.02 ± 0.00	0.05 ± 0.00	4.11 ± 0.32

Table 5: Results in RWARE tasks.

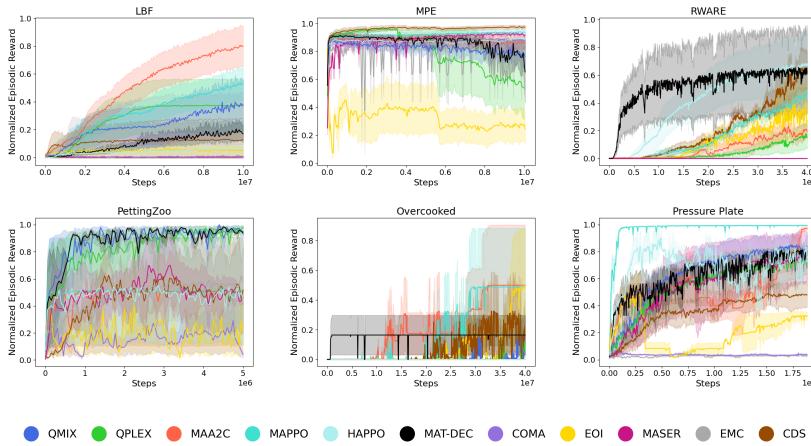


Figure 1: Aggregated Normalized episodic rewards of each benchmark.

Algorithms\Tasks	Spread-4	Spread-5	Spread-8
QMIX	-1278.26 ± 23.13	-2531.17 ± 586.56	-6414.48 ± 27.27
QPLEX	-766.84 ± 14.39	-1800.53 ± 194.49	-13260.36 ± 6200.03
MAA2C	-1190.09 ± 99.93	-2312.56 ± 222.08	-5961.67 ± 66.04
MAPPO	-971.17 ± 124.22	-1910.20 ± 42.86	-5926.39 ± 38.48
HAPPO	-1032.80 ± 45.84	-2000.41 ± 98.20	-6940.61 ± 69.55
MAT-DEC	-1066.62 ± 45.98	-1918.88 ± 15.76	-6843.44 ± 563.39
COMA	-1176.78 ± 33.37	-2003.47 ± 51.18	-6249.07 ± 44.73
EOI	-1963.23 ± 859.27	-5816.66 ± 20.34	-13210.98 ± 6569.53
MASER	-969.06 ± 5.01	-1939.58 ± 113.22	-6242.61 ± 14.40
EMC	-1216.19 ± 10.9	-1961.75 ± 0.71	-6219.14 ± 29.55
CDS	-809.20 ± 47.22	-1641.77 ± 14.61	-6250.72 ± 15.44

Table 6: Results in Spread (MPE) tasks.

Algorithms\Envs	Cramped Room	Asymmetric Advantages	Coordination Ring
QMIX	0.00 ± 0.00	300.00 ± 300.00	0.00 ± 0.00
QPLEX	86.67 ± 122.57	0.00 ± 0.00	0.00 ± 0.00
MAA2C	286.80 ± 9.34	487.80 ± 107.60	0.10 ± 0.10
MAPPO	280.00 ± 0.00	0.30 ± 0.10	0.07 ± 0.09
HAPPO	0.00 ± 0.00	160.10 ± 159.9	0.00 ± 0.00
MAT-DEC	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
COMA	0.20 ± 0.16	0.10 ± 0.10	0.07 ± 0.09
EOI	280.00 ± 0.00	1.60 ± 0.60	0.13 ± 0.09
EMC	0.00 ± 0.00	—	—
MASER	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
CDS	186.67 ± 133.00	70.00 ± 70.00	0.00 ± 0.00

Table 8: Results in Overcooked environments.

Algorithms\Envs	Pistonball	Cooperative Pong	Entompted Cooperative
QMIX	991.59 ± 0.17	199.57 ± 0.61	8.00 ± 0.00
QPLEX	991.46 ± 0.10	197.78 ± 3.14	8.00 ± 0.00
MAA2C	990.83 ± 0.04	-0.33 ± 3.47	6.57 ± 0.05
MAPPO	990.71 ± 0.10	13.22 ± 4.61	6.61 ± 0.05
HAPPO	983.60 ± 0.77	25.62 ± 3.36	7.99 ± 0.01
MAT-DEC	982.57 ± 2.44	200.00 ± 0.00	10.00 ± 0.00
COMA	678.28 ± 324.06	1.10 ± 7.51	7.68 ± 1.7
EOI	948.35 ± 42.74	-1.64 ± 2.66	6.53 ± 0.02
MASER	989.39 ± 0.64	104.25 ± 69.29	8.00 ± 0.00
EMC	265.00 ± 174.58	196.5 ± 2.83	8.00 ± 0.00
CDS	415.82 ± 284.37	197.13 ± 2.26	11.10 ± 1.00

Table 7: Results in PettingZoo using ResNet18.

Algorithms\Envs	4p	6p
QMIX	-210.72 ± 17.84	-3461.77 ± 1020.33
QPLEX	-652.19 ± 9.29	-5183.56 ± 345.46
MAA2C	-281.59 ± 201.77	-547.39 ± 21.00
MAPPO	-135.99 ± 1.32	-494.08 ± 10.54
HAPPO	-258.69 ± 224.93	-584.90 ± 201.68
MAT-DEC	-876.58 ± 1113.53	-2930.77 ± 3652.53
COMA	-4391.79 ± 108.94	-12360.20 ± 314.06
EOI	-3050.64 ± 1125.83	-9221.16 ± 4796.08
MASER	-88.44 ± 2.84	-5257.08 ± 4099.19
EMC	-4518.23 ± 249.68	-12347.60 ± 0.0
CDS	-1926.45 ± 260.85	-8068.23 ± 519.79

Table 9: Results in Pressure Plate tasks.

Algorithms	Episodic Reward	RAM	Training time
MAA2C-ResNet18	990.83 ± 0.04	7GB	4d : 15h
MAA2C-CNN	846.74 ± 19.7	34GB	16d : 0h

Table 10: Results in PettingZoo’s Pistonball comparing ResNet18 as frozen image encoder with trainable CNNs.

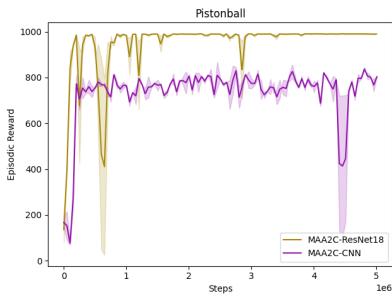


Figure 2: ResNet18 vs Trainable CNNs: MAA2C in PettingZoo’s Pistonball task.

6.2 Main Results and Analysis

In this section, we present our main results and the benchmark analysis. We show the main analytical results in Tables 4, 5, 6, 7, 8 and 9, and the averaged main results in Figure 1. To obtain these average results, we use the *normalized* scores for each task of the benchmark and average over these scores. In the tables, we write in bold the values of the best algorithm in a specific task. If the performance of another algorithm was not statistically significantly different from the best algorithm, the respective value is also in bold. Algorithms with green cell color are the best on average in the corresponding benchmark.

6.2.1 *Key Findings.* Based on results, we highlight the following:

- **Best Algorithms.** The standard MARL algorithms, QPLEX, MAPPO, MAA2C and CDS, demonstrate the most consistent performance across all fully cooperative benchmarks.
- **SoTA exploration-based algorithms mostly underperform.** Exploration-based methods that have reached SoTA performance in SMAC and/or GRF, that is, EOI, EMC and MASER, significantly underperform in most benchmarks compared to the standard methods. In some cases, they even fail entirely, including in tasks with sparse reward settings (e.g., see Tables 4, 5 and 8), in which these methods are expected to improve performance over baselines. The only exception is CDS which is shown to be one of the best algorithms in RWARE and Spread.
- **Value Decomposition.** As demonstrated by the results of QMIX and QPLEX, value decomposition methods tend to converge to suboptimal policies as the number of agents increases, significantly underperforming compared to actor-critic methods (e.g., see the results in Spread, LBF and Pressure Plate). Moreover, our findings contrast with one conclusion of [36] that suggests that value decomposition methods

require sufficiently dense rewards to effectively learn to decompose the value function. We show several cases with sparse reward settings, such as in LBF (e.g., 2s-8x8-3p-2f, 2s-9x9-3p-2f), RWARE (*tiny-4ag-hard*), and Pressure Plate (4p), where QPLEX, or even QMIX, can successfully converge to effective policies.

6.2.2 *Discussion on Algorithms’ Performance.* Next, we analyse the performance of the evaluated MARL algorithms in the selected fully cooperative benchmark tasks.

QMIX, QPLEX. QMIX shows mediocre performance across most tasks, with complete failure in some (e.g., RWARE and many LBF tasks), yet the highest rewards in PettingZoo. Notably, QMIX is the only method to achieve positive rewards in the sparse-reward LBF 4s-11x11-3p-2f task. In contrast, QPLEX generally outperforms QMIX improving performance in cooperative tasks. However, QPLEX fails in the most challenging LBF and Overcooked tasks and struggles in large-scale scenarios such as *Spread-8* and Pressure Plate’s 6p. Interestingly, despite lacking advanced exploration techniques, QPLEX notably improves over QMIX in RWARE.

MAA2C, MAPPO, COMA. MAA2C is one of the most consistent algorithms across all tasks, showing the best performance in complex sparse-reward LBF tasks and in *Spread-8* (along with MAPPO) and Overcooked. Similarly, MAPPO also ranks as one of the most consistent algorithms but significantly outperforms MAA2C in RWARE, *Spread-5*, *Cooperative Pong*, and Pressure Plate’s 6p. MAPPO achieves the highest rewards on average in *Spread* (together with CDS) and Pressure Plate. Interestingly, in the RWARE tasks where MAPPO had been the best option in previous work [36], HAPPO, MAT-DEC, and CDS outperform it. Conversely, COMA has one of the lowest performances overall, with competitive results only in *Spread*, *Entombed Cooperative*, and *Coordination Ring*. However, for the last two tasks, all algorithms perform poorly, so none are truly competitive.

HAPPO, MAT-DEC. Both HAPPO and MAT-DEC do not manage to achieve consistent performance across the evaluated benchmarks. More specifically, both methods are shown to be among the best (together with CDS) in RWARE, while HAPPO also achieving good performance in Pressure Plate and MAT-DEC in the high-dimensional PettingZoo tasks. However, both HAPPO and MAT-DEC perform poorly in tasks where efficient exploration and synchronized coordination are crucial, such as in LBF and Overcooked benchmarks.

MASER, EMC. Despite MASER and EMC achieving strong results in sparse-reward SMAC tasks, they perform poorly in most fully cooperative tasks. The notable exception is PettingZoo tasks, where MASER’s baseline, QMIX, performs better. MASER also achieves the highest rewards in Pressure Plate (4p). Both MASER’s and EMC’s poor overall performance is attributed to its reliance on Q-values to enhance joint exploration: Since agents rarely receive positive rewards, this leads to misleading intrinsic rewards based on irrelevant Q-value-driven sub-goals. However, Q-values prove to be useful for guiding exploration in LBF’s 2s-12x12-2p-2f, as this task is less sparse but also requires good exploration for success.

EOI, CDS. EOI and CDS outperform MASER and EMC in most tasks. EOI achieves the highest rewards in Overcooked’s *Cramped*

Environments\Algorithms	QMIX	QPLEX	MAA2C	MAPPO	HAPPO	MAT-DEC	COMA	EOI	MASER	EMC	CDS
LBF	0d : 8h	0d : 13h	0d : 1h	0d : 2h	0d : 5h	0d : 4h	0d : 1h	0d : 6h	0d : 18h	2d : 4h	0d : 18h
RWARE	1d : 22h	2d : 17h	0d : 9h	0d : 12h	1d : 1h	0d : 20h	0d : 9h	1d : 18h	2d : 16h	19d : 1h	3d : 2h
Spread (MPE)	0d : 15h	0d : 21h	0d : 2h	0d : 3h	0d : 9h	0d : 6h	0d : 2h	0d : 11h	1d : 12h	4d : 9h	4d : 9h
Petting Zoo	1d : 16h	3d : 11h	3d : 23h	3d : 10h	14d : 1h	3d : 6h	3d : 11h	0d : 23h	2d : 1h	3d : 23h	1d : 19h
Overcooked	3d : 14h	5d : 3h	0d : 19h	1d : 4h	1d : 18h	2d : 11h	0d : 21h	3d : 9h	4d : 9h	13d : 14h	2d : 14h
Pressure Plate	0d : 22h	1d : 14h	0d : 4h	0d : 7h	0d : 20h	0d : 12h	0d : 4h	0d : 18h	1d : 6h	9d : 7h	2d : 1h

Table 11: Average (wall-clock) training times over all tasks of each benchmark for all 11 algorithms.

Room (alongside MAA2C and MAPPO) and RWARE’s *tiny-2ag-hard*. However, EOI lacks the consistency of the MAA2C backbone algorithm. We attribute this inconsistency to the emphasis on individuality, which can hinder full cooperation and the high level of coordination required in most fully cooperative tasks. In contrast, CDS is more consistent, showing top performance in Spread and RWARE (alongside MAT-DEC and HAPPO) tasks. CDS’ improved performance is attributed to the fact that although CDS relies on encouraging agents to be more diverse, as EOI does, it also assures sufficient sharing of the most useful experience of the agents.

6.2.3 Pre-trained Image-Encoder vs Trainable CNNs for Image-based Observations. Our experimental results underscore the effectiveness of employing a frozen pre-trained image encoder, such as ResNet18, over trainable CNNs in multi-agent reinforcement learning (MARL) environments. In the specific case of the PettingZoo’s Pistonball task, the use of frozen ResNet18 led to more stable and higher overall policy performance compared to its trainable counterparts. This advantage is clearly illustrated by the smoother convergence curves and the consistently superior performance throughout the training steps, as shown in Figure 2. Moreover, as detailed in Table 10, the non-adaptive nature of frozen ResNet18 highlights its efficiency in managing complex visual contexts without computational overhead. This suggests that pre-trained, static encoders are highly beneficial, providing immediate and reliable performance improvements. The adoption of frozen image encoders significantly improves computational efficiency by eliminating the need for ongoing adjustments during training. This leads to shorter training times for MARL algorithms, as detailed in Table 10. Moreover, by converting image data into compact vector representations, these encoders substantially lower memory requirements, facilitating the execution of complex MARL tasks with RGB-based observations.

6.2.4 Comparison of training times. As can be clearly seen from Table 11, off-policy algorithms generally require longer training due to their use of large replay buffers, which process experiences from multiple past episodes. In contrast, on-policy algorithm, benefit from learning directly from current policy experiences, allowing for faster training. However, in image-based, high-dimensional environments, such as PettingZoo tasks, this parallelism can slow down training, as it does not integrate well with pre-trained models, such as ResNet18, which require GPU resources.

6.2.5 Open Challenges. Based on the results we discussed above, it is evident that fully cooperative MARL tasks need careful algorithmic design, in terms of both excessive coordination and joint exploration, as current SoTA and standard MARL methods are not very effective. Below, we report some significant open challenges arised from our benchmark analysis:

(a) The tasks *Entombed Cooperative* (PettingZoo) and *Coordination*

Ring (Overcooked) are the most challenging, as all evaluated algorithms indeed fail to find any effective policy. Any improvement on these tasks would be of remarkable interest.

- (b) The sparse-reward LBF tasks, with a large grid, three agents and two foods, the sparse-reward Pressure Plate tasks, with more than 4 agents, and the sparse-reward hard RWARE tasks, with larger grids, are quite challenging, as they require excessive coordinated exploration, and any improvement on these is very interesting.
- (c) The Spread tasks, with more than 4 agents are quite challenging, as they require excessive coordination, and any improvement on these without the use of agent communication during execution is very interesting.

7 RELATED WORK

The recent rise in MARL popularity has fragmented community standards and tools, with the frequent introduction of new libraries such as [20, 30, 42, 49]. Among the most popular are PyMARL [43] and EPyMARL [36], both of which have played a crucial role in driving the influx of cooperative MARL algorithms. However, these libraries have somewhat overlooked the integration of fully cooperative MARL environments, pushing researchers to focus on specific benchmarks, such as SMAC [11, 43] and GRF [27], raising concerns about the reliability and generalizability of the proposed algorithms [16]. Despite recent efforts [4, 20, 36, 63] aiming to provide a comprehensive understanding of standard cooperative MARL algorithms through benchmarking, the evaluation of fully cooperative MARL still lacks systematic diversity and reliability.

8 CONCLUSION

In this paper, we highlight and address key concerns in the evaluation of cooperative MARL algorithms by providing an extended benchmarking of well-known MARL methods in fully cooperative tasks. Our extensive evaluations reveal significant discrepancies in the performance of SoTA methods, which can eventually underperform compared to standard baselines. Based on our analysis, as well as by open-sourcing PyMARLzoo+, this paper aims to motivate towards more systematic evaluation of MARL algorithms, encouraging broader adoption of diverse, fully cooperative benchmarks.

ACKNOWLEDGMENTS

The research work contributed by George Papadopoulos was partially supported by the Hellenic Foundation for Research and Innovation (HFRI) under the 5th Call for HFRI PhD Fellowships (Fellowship Number: 20769). The authors also gratefully acknowledge GRNET – National Infrastructures for Research and Technology for providing AWS resources, and the University of Piraeus for funding the participation in the AAMAS 2025 conference and associated travel costs.

REFERENCES

- [1] Ibrahim H Ahmed, Cillian Brewitt, Ignacio Carlucho, Filippos Christianos, Mhairi Dunion, Elliot Fosong, Samuel Garcin, Shangmin Guo, Balint Gyevnar, Trevor McInroe, et al. 2022. Deep reinforcement learning for multi-agent interaction. *Ai Communications* 35, 4 (2022), 357–368.
- [2] Stefano V Albrecht, Filippos Christianos, and Lukas Schäfer. 2024. *Multi-agent reinforcement learning: Foundations and modern approaches*. MIT Press.
- [3] Stefano V. Albrecht and Peter Stone. 2017. Reasoning about Hypothetical Agent Behaviours and their Parameters. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems* (São Paulo, Brazil) (AAMAS '17). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 547–555.
- [4] Matteo Bettini, Amanda Prorok, and Vincent Moens. 2024. Benchmark: Benchmarking multi-agent reinforcement learning. *Journal of Machine Learning Research* 25, 217 (2024), 1–10.
- [5] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym. arXiv:arXiv:1606.01540
- [6] Wolfram Burgard, Mark Moors, Dieter Fox, Reid Simmons, and Sebastian Thrun. 2000. Collaborative multi-robot exploration. In *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings (Cat. No. 00CH37065)*, Vol. 1. IEEE, IEEE, 476–481.
- [7] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. 2019. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems* 32 (2019).
- [8] Changyu Chen, Ramesha Karunasena, Thanh Nguyen, Arunesh Sinha, and Pradeep Varakantham. 2024. Generative modelling of stochastic actions with arbitrary constraints in reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [9] Zigen Chen, Gongfan Fang, Xinyin Ma, and Xinchao Wang. 2024. SlimSAM: 0.1% Data Makes Segment Anything Slim. arXiv:2312.05284 [cs.CV] https://arxiv.org/abs/2312.05284
- [10] Filippos Christianos, Georgios Papoudakis, Muhammad A Rahman, and Stefano V Albrecht. 2021. Scaling multi-agent reinforcement learning with selective parameter sharing. In *International Conference on Machine Learning*. PMLR, 1989–1998.
- [11] Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Mahajan, Jakob Foerster, and Shimon Whiteson. 2024. Smacv2: An improved benchmark for cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [12] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [13] Center for Human-Compatible AI Github contributors. [n.d.]. Github Overcooked. https://github.com/HumanCompatibleAI/overcooked_ai. Overcooked GitHub repository.
- [14] Elliot Fosong, Arrasy Rahman, Ignacio Carlucho, and Stefano V Albrecht. 2024. Learning Complex Teamwork Tasks using a Given Sub-task Decomposition. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*. 598–606.
- [15] Jialu Gao, Kaizhe Hu, Guowei Xu, and Huazhe Xu. 2024. Can pre-trained text-to-image models generate visual goals for reinforcement learning? *Advances in Neural Information Processing Systems* 36 (2024).
- [16] Rihab Gorsane, Omayma Mahjoub, Ruan John de Kock, Roland Dubb, Siddarth Singh, and Arnu Pretorius. 2022. Towards a standardised performance evaluation protocol for cooperative marl. *Advances in Neural Information Processing Systems* 35 (2022), 5510–5521.
- [17] Nikunj Gupta, Somjit Nath, and Samira Ebrahimi Kahou. 2023. CAMMARL: Conformal Action Modeling in Multi Agent Reinforcement Learning. arXiv preprint arXiv:2306.11128 (2023).
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [19] Joey Hong, Sergey Levine, and Anca Dragan. 2024. Learning to influence human behavior with offline reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [20] Siyi Hu, Yifan Zhong, Minquan Gao, Weixun Wang, Hao Dong, Zhihui Li, Xiaodan Liang, Yaodong Yang, and Xiaojun Chang. 2022. Marllib: Extending rllib for multi-agent reinforcement learning. (2022).
- [21] Jeewon Jeon, Woojun Kim, Whiyoung Jung, and Youngchul Sung. 2022. Maser: Multi-agent reinforcement learning with subgoals generated from experience replay buffer. In *International Conference on Machine Learning*. PMLR, 10041–10052.
- [22] Jeewon Jeon, Woojun Kim, Whiyoung Jung, and Youngchul Sung. 2022. Maser: Multi-agent reinforcement learning with subgoals generated from experience replay buffer. In *International Conference on Machine Learning*. PMLR, 10041–10052.
- [23] Chengzhi Jiang and Zhaohan Sheng. 2009. Case-based reinforcement learning for dynamic inventory control in a multi-agent supply-chain system. *Expert Systems with Applications* 36, 3 (2009), 6520–6526.
- [24] Jiechuan Jiang and Zongqing Lu. 2021. The emergence of individuality. In *International Conference on Machine Learning*. PMLR, 4992–5001.
- [25] Andreas Kontogiannis and George A Vouros. 2023. Inherently Interpretable Deep Reinforcement Learning Through Online Mimicking. In *International Workshop on Explainable, Transparent Autonomous Agents and Multi-Agent Systems*. Springer, 160–179.
- [26] Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. 2021. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint arXiv:2109.11251* (2021).
- [27] Karol Kurach, Anton Raichuk, Piotr Steniczyk, Michał Zająć, Olivier Bachem, Lasse Espeholt, Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. 2020. Google research football: A novel reinforcement learning environment. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 4501–4510.
- [28] Chenghao Li, Tonghan Wang, Chengjie Wu, Qianchuan Zhao, Jun Yang, and Chongjie Zhang. 2021. Celebrating diversity in shared multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 34 (2021), 3991–4002.
- [29] Xihan Li, Jia Zhang, Jiang Bian, Yunhai Tong, and Tie-Yan Liu. 2019. A Cooperative Multi-Agent Reinforcement Learning Framework for Resource Balancing in Complex Logistics Network. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*. 980–988.
- [30] Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gonzalez, Michael Jordan, and Ion Stoica. 2018. RLlib: Abstractions for distributed reinforcement learning. In *International conference on machine learning*. PMLR, 3053–3062.
- [31] Bo Liu, Qiang Liu, Peter Stone, Animesh Garg, Yuke Zhu, and Anima Anandkumar. 2021. Coach-player multi-agent reinforcement learning for dynamic team composition. In *International Conference on Machine Learning*. PMLR, 6860–6870.
- [32] Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems* 30 (2017).
- [33] Igor Mordatch and Pieter Abbeel. 2018. Emergence of grounded compositional language in multi-agent populations. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [34] Frans A Oliehoek, Christopher Amato, et al. 2016. *A concise introduction to decentralized POMDPs*. Vol. 1. Springer.
- [35] Shayegan Omidshafiei, Jason Pazis, Christopher Amato, Jonathan P How, and John Vian. 2017. Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In *International Conference on Machine Learning*. PMLR, 2681–2690.
- [36] Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. 2021. Benchmarking Multi-Agent Deep Reinforcement Learning Algorithms in Cooperative Tasks. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS)*.
- [37] Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamenny, Philip Torr, Wendelin Böhmer, and Shimon Whiteson. 2021. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems* 34 (2021), 12208–12221.
- [38] Alex Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [39] Aowabin Rahaman, Arnab Bhattacharya, Thiagarajan Ramachandran, Sayak Mukherjee, Himanshu Sharma, Ted Fujimoto, and Samrat Chatterjee. 2022. Adversar: Adversarial search and rescue via multi-agent reinforcement learning. In *2022 IEEE International Symposium on Technologies for Homeland Security (HST)*. IEEE, 1–7.
- [40] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2020. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research* 21, 178 (2020), 1–51.
- [41] Jingqing Ruan, Yali Du, Xuantang Xiong, Dengpeng Xing, Xiyun Li, Linghui Meng, Haifeng Zhang, Jun Wang, and Bo Xu. 2022. GCS: Graph-based coordination strategy for multi-agent reinforcement learning. *arXiv preprint arXiv:2201.06257* (2022).
- [42] Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Gardar Ingvarsson, Timon Willi, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, et al. 2024. JaxMARL: Multi-Agent RL Environments and Algorithms in JAX. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*. 2444–2446.
- [43] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. 2019. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043* (2019).

- [44] Tom Schaul. 2015. Prioritized Experience Replay. *arXiv preprint arXiv:1511.05952* (2015).
- [45] Sven Seuken and Shlomo Zilberstein. 2007. Improved memory-bounded dynamic programming for decentralized POMDPs. In *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence* (Vancouver, BC, Canada) (*UAI'07*). AUAI Press, Arlington, Virginia, USA, 344–351.
- [46] Rutav M Shah and Vikash Kumar. 2021. RRL: Resnet as representation for Reinforcement Learning. In *International Conference on Machine Learning*. PMLR, 9465–9476.
- [47] Chapman Siu, Jason Traish, and Richard Yi Da Xu. 2021. Dynamic coordination graph for cooperative multi-agent reinforcement learning. In *Asian Conference on Machine Learning*. PMLR, 438–453.
- [48] Jing Sun, Shuo Chen, Cong Zhang, Yining Ma, and Jie Zhang. 2024. Decision-Making With Speculative Opponent Models. *IEEE Transactions on Neural Networks and Learning Systems* (2024).
- [49] Jordan Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. 2021. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 34 (2021), 15032–15043.
- [50] Justin K Terry, Benjamin Black, and Luis Santos. 2020. Multiplayer support for the arcade learning environment. *arXiv preprint arXiv:2009.09341* (2020).
- [51] Justin K Terry, Nathaniel Grammel, Sanghyun Son, Benjamin Black, and Aakriti Agrawal. 2020. Revisiting parameter sharing in multi-agent deep reinforcement learning. *arXiv preprint arXiv:2005.13625* (2020).
- [52] Filippos Christianos Trevor McInroe. [n.d.]. Github pressureplate. <https://github.com/uec-agents/pressureplate/tree/main>. Pressureplate GitHub repository.
- [53] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. 2020. Qplex: Duplex dueling multi-agent q-learning. *arXiv preprint arXiv:2008.01062* (2020).
- [54] Tonghan Wang, Heng Dong, Victor Lesser, and Chongjie Zhang. 2020. ROMA: multi-agent reinforcement learning with emergent roles. In *Proceedings of the 37th International Conference on Machine Learning*, 9876–9886.
- [55] Muning Wen, Jakub Kuba, Runji Lin, Weinan Zhang, Ying Wen, Jun Wang, and Yaodong Yang. 2022. Multi-agent reinforcement learning is a sequence modeling problem. *Advances in Neural Information Processing Systems* 35 (2022), 16509–16521.
- [56] Yuchen Xiao, Joshua Hoffman, and Christopher Amato. 2019. Macro-Action-Based Deep Multi-Agent Reinforcement Learning. In *3rd Annual Conference on Robot Learning*.
- [57] Yuchen Xiao, Joshua Hoffman, and Christopher Amato. 2020. Macro-action-based deep multi-agent reinforcement learning. In *Conference on Robot Learning*. PMLR, 1146–1161.
- [58] Yuchen Xiao, Joshua Hoffman, Tian Xia, and Christopher Amato. 2020. Learning multi-robot decentralized macro-action-based policies via a centralized q-net. In *2020 IEEE International conference on robotics and automation (ICRA)*. IEEE, 10695–10701.
- [59] Yuchen Xiao, Weihao Tan, and Christopher Amato. 2022. Asynchronous actor-critic for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 35 (2022), 4385–4400.
- [60] Mingyu Yang, Yaodong Yang, Zhenbo Lu, Wengang Zhou, and Houqiang Li. 2024. Hierarchical multi-agent skill discovery. *Advances in Neural Information Processing Systems* 36 (2024).
- [61] Yaodong Yang, Guangyong Chen, Weixun Wang, Xiaotian Hao, Jianye Hao, and Pheng-Ann Heng. 2022. Transformer-based working memory for multiagent reinforcement learning with action parsing. *Advances in Neural Information Processing Systems* 35 (2022), 34874–34886.
- [62] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. 2022. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems* 35 (2022), 24611–24624.
- [63] Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. 2020. Benchmarking multi-agent deep reinforcement learning algorithms. (2020).
- [64] Lulu Zheng, Jiarui Chen, Jianhao Wang, Jiamin He, Yujing Hu, Yingfeng Chen, Changjie Fan, Yang Gao, and Chongjie Zhang. 2021. Episodic multi-agent reinforcement learning with curiosity-driven exploration. *Advances in Neural Information Processing Systems* 34 (2021), 3757–3769.
- [65] Lianmin Zheng, Jiacheng Yang, Han Cai, Ming Zhou, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Magent: A many-agent reinforcement learning platform for artificial collective intelligence. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [66] Yifan Zhong, Jakub Grudzien Kuba, Xidong Feng, Siyi Hu, Jiaming Ji, and Yaodong Yang. 2024. Heterogeneous-agent reinforcement learning. *Journal of Machine Learning Research* 25, 1-67 (2024), 1.

A ENVIRONMENT API

The following Python script (Figure 3) provides an example of executing an episode in Butterfly’s Pistonball task of PettingZoo using random actions under our framework. To use a different environment, the user only needs to modify the *args* variable accordingly. Below, we provide more detailed instructions. The script demonstrates: (a) importing necessary packages, (b) initializing the environment, and (c) running an episode by repeatedly rendering the environment and applying random actions until the episode ends.

```

1 # Import packages
2 from envs import REGISTRY as env_REGISTRY
3 import random as rnd
4 # Arguments for PettingZoo
5 args = {
6     "env": "pettingzoo",
7     "env_args": {
8         "key": "pistonball_v6",
9     }
10 }
11 # Initialize environment
12 env = env_REGISTRY[args["env"]](**args["env_args"])
13 n_agnts = env.n_agents
14 n_acts = env.get_total_actions()
15 # Reset the environment
16 obs, state = env.reset()
17 done = False
18 # Run an episode
19 while not done:
20     # Render the environment
21     env.render()
22     # Insert the policy's actions here
23     actions = rnd.choices(range(0, n_acts), k=n_agnts)
24     # Apply an environment step
25     reward, done, info = env.step(actions)
26     obs = env.get_obs()
27     state = env.get_state()
28 # Terminate the environment
29 env.close()
```

Figure 3: Indicative Python script for executing an episode in PettingZoo’s Pistonball.

Following the example presented above (Figure 3), we provide the minimum *args* required for each environment, as well as all of those that can be specified.

PettingZoo args:

```

args = {
    "env": "pettingzoo",
    "env_args": {
        "key": "pistonball_v6",
    }
}
```

Figure 4: Example of minimal *args* for PettingZoo tasks.

```

args = {
    "env": "pettingzoo",
    "env_args": {
        "key": "pistonball_v6",
        "seed": 1,
        "time_limit": 900,
        "partial_observation": False,
        "trainable_cnn": False,
        "image_encoder": "ResNet18",
        "image_encoder_batch_size": 10,
        "image_encoder_use_cuda": True,
        "centralized_image_encoding": True,
        "kwargs": ""
    }
}

```

Figure 5: Example of all *args* for PettingZoo tasks.

where:

- *key* is the selected PettingZoo task. Options: "pistonball_v6", "cooperative_pong_v5", "entombed_cooperative_v3", "space_invaders_v2", "basketball_pong_v3", "boxing_v2", "combat_jet_v1", "combat_tank_v3", "double_dunk_v3", "entombed_competitive_v3", "flag_capture_v2", "foozpong_v3", "ice_hockey_v2", "joust_v3", "mario_bros_v3", "maze_craze_v3", "othello_v3", "pong_v3", "quadrapong_v4", "space_war_v2", "surround_v2", "tennis_v3", "video_checkers_v4", "volleyball_pong_v2", "warlords_v3", "wizard_of_wor_v3", "knights_archers_zombies_v10", "chess_v6", "connect_four_v3", "gin_rummy_v4", "go_v5", "hanabi_v5", "leduc_holdem_v4", "rps_v2", "texas_holdem_no_limit_v6", "texas_holdem_v4", "tictactoe_v3", "simple_v3", "simple_adversary_v3", "simple_crypto_v3", "simple_push_v3", "simple_reference_v3", "simple_speaker_listener_v4", "simple_spread_v3", "simple_tag_v3", "simple_world_comm_v3", "multiwalker_v9", "pursuit_v4", "waterworld_v4".
- *seed* specifies the random sequence of the first environment state which is helpful for reproducible evaluations. The default value is 1.
- *time_limit* is the maximum number of steps before episode termination. The default value is 900 for all Butterfly tasks except from Pistonball, which is 125, and 10000 for all Atari tasks. For all SISL tasks the default value is 500, for all MPE tasks is 25,
- *partial_observation* defines whether to use partial observations or not. Only applicable in "space_invaders_v2" and "entombed_cooperative_v3". The default value is False.
- *trainable_cnn* defines whether to encode the image observations using a pre-trained model or to return the raw images. The default value is False.
- *image_encoder* is the selected image-encoder to use for encoding the image observations. Only applicable if *trainable_cnn* is True. Options: "ResNet18", "SlimSAM", "CLIP". The default value is "ResNet18".
- *image_encoder_batch_size* is the number of images to encode at once with the selected image-encoder. Only applicable if *trainable_cnn* is True. The default value is 10.
- *image_encoder_use_cuda* defines whether to use GPU or CPU for the selected image-encoder. In this way, the policies networks can be in a different device than the image-encoder. The default value is True, but it automatically turns to True if no GPU available is found.
- *centralized_image_encoding* defines whether to encode images in a centralized manner, thus avoiding to load as many image-encoders as the parallel environment processes, reducing the GPU memory required. However, this is only supported by our training API, where multiple parallel processes can be managed. When it is used only with the environment API, the effect is that the raw image observations are returned instead of the encoded representations.
- *kwargs* can consist of specific arguments supported by the selected task². The user can provide the arguments in following format: "('arg1',arg1_value), ('arg2',arg2_value), ...". For instance, we can specify the number of pistons in the *pistonball_v6* task using: "('n_pistons',4)". The default value is "", that is, no arguments provided.

Overcooked args:

```

args = {
    "env": "overcooked",
    "env_args": {
        "key": "cramped_room",
    }
}

```

Figure 6: Example of minimal *args* for Overcooked tasks.

²All the arguments of each task can be found on the official page of PettingZoo: <https://pettingzoo.farama.org/>

```

args = {
    "env": "overcooked",
    "env_args": {
        "key": "cramped_room",
        "seed": 1,
        "time_limit": 500,
        "reward_type": "sparse"
    }
}

```

Figure 7: Example of all *args* for Overcooked tasks.

where:

- *key* is the selected Overcooked scenario. Options: "cramped_room", "asymmetric_advantages", "coordination_ring", "counter_circuit", "forced_coordination".
- *seed* has the same functionality as in PettingZoo. The default value is 1.
- *time_limit* has the same functionality as in PettingZoo. The default value is 500.
- *reward_type* defines whether to provide sparse rewards, that is, only when an order is ready, or shaped rewards, that is, each time a component of an order is ready. Options: "sparse", "shaped". The default value is "sparse".

Pressure Plate args:

```

args = {
    "env": "pressureplate",
    "env_args": {
        "key": "pressureplate-linear-4p-v0",
    }
}

```

Figure 8: Example of minimal *args* for Pressure Plate environment.

```

args = {
    "env": "pressureplate",
    "env_args": {
        "key": "pressureplate-linear-4p-v0",
        "seed": 1,
        "time_limit": 500,
    }
}

```

Figure 9: Example of all *args* for Pressure Plate environment.

where:

- *key* is the selected Pressure Plate scenario. Options: "pressureplate-linear-4p-v0", "pressureplate-linear-5p-v0", "pressureplate-linear-6p-v0".
- *seed* has the same functionality as in PettingZoo. The default value is 1.
- *time_limit* has the same functionality as in PettingZoo. The default value is 500.

Capture Target args:

```

args = {
    "env": "capturetarget",
    "env_args": {
        "key": "CaptureTarget-6x6-1t-2a-v0",
    }
}

```

Figure 10: Example of minimal *args* for Capture Target environment.

```

args = {
    "env": "capturetarget",
    "env_args": {
        "key": "CaptureTarget-6x6-1t-2a-v0",
        "seed": 1,
        "time_limit": 60,
        "obs_one_hot": False,
        "target_flick_prob": 0.3,
        "tgt_avoid_agent": True,
        "tgt_trans_noise": 0.0,
        "agent_trans_noise": 0.1,
    }
}

```

Figure 11: Example of all *args* for Capture Target environment.

where:

- *key* is the selected Capture Target scenario. Options: "CaptureTarget-6x6-1t-2a-v0".
- *seed* has the same functionality as in PettingZoo. The default value is 1.
- *time_limit* has the same functionality as in PettingZoo. The default value is 60.
- *obs_one_hot* defines whether to return the observations in one-hot encoding format. The default value is False.
- *target_flick_prob* specifies the probability of not observing the target. The default value is 0.3.
- *tgt_avoid_agent* defines whether the target keeps moving away from the agents or not. The default value is True.
- *tgt_trans_noise* specifies the target's transition probability for arriving an unintended adjacent cell. The default value is 0.0.
- *agent_trans_noise* specifies the agent's transition probability for arriving an unintended adjacent cell. The default value is 0.1.

Box Pushing *args*:

```

args = {
    "env": "boxpushing",
    "env_args": {
        "key": "BoxPushing-6x6-2a-v0",
    }
}

```

Figure 12: Example of minimal *args* for Box pushing environment.

```

args = {
    "env": "boxpushing",
    "env_args": {
        "key": "BoxPushing-6x6-2a-v0",
        "seed": 1,
        "time_limit": 60,
        "random_init": True
    }
}

```

Figure 13: Example of all *args* for Box Pushing environment.

where:

- *key* is the selected Box Pushing scenario. Options: "BoxPushing-6x6-2a-v0".
- *seed* has the same functionality as in PettingZoo. The default value is 1.
- *time_limit* has the same functionality as in PettingZoo. The default value is 60.
- *random_init* shows whether random initial position of agents should be applied. The default value is True.

LBF, RWARE, MPE *args*:

```

args = {
    "env": "gymma",
    "env_args": {
        "key": "mpe:SimpleSpeakerListener-v0",
    }
}

```

Figure 14: Example of minimal *args* for LBF, RWARE, and MPE environments.

```

args = {
    "env": "gymma",
    "env_args": {
        "key": "mpe:SimpleSpeakerListener-v0",
        "seed": 1,
        "time_limit": 500,
    }
}

```

Figure 15: Example of all *args* for LBF, RWARE, and MPE environments.

where:

- *key* is the selected scenario of LBF, RWARE, or MPE environments.

Options for LBF:

"lbforaging:Foraging-4s-11x11-3p-2f-coop-v2", "lbforaging:Foraging-2s-11x11-3p-2f-coop-v2", "lbforaging:Foraging-2s-8x8-3p-2f-coop-v2", "lbforaging:Foraging-2s-9x9-3p-2f-coop-v2", "lbforaging:Foraging-7s-20x20-5p-3f-coop-v2", "lbforaging:Foraging-2s-12x12-2p-2f-coop-v2", "lbforaging:Foraging-8s-25x25-8p-5f-coop-v2", "lbforaging:Foraging-7s-30x30-7p-4f-coop-v2".

Options for RWARE:

"rware:rware-small-4ag-hard-v1", "rware:rware-tiny-4ag-hard-v1", "rware:rware-tiny-2ag-hard-v1".

Options for MPE:

"mpe:SimpleSpread-3-v0", "mpe:SimpleSpread-4-v0", "mpe:SimpleSpread-5-v0", "mpe:SimpleSpread-8-v0", "mpe:SimpleSpeakerListener-v0", "mpe:MultiSpeakerListener-v0".

- *seed* has the same functionality as in PettingZoo. The default value is 1.
- *time_limit* has the same functionality as in PettingZoo. The default value is 50 for LBF, 500 for RWARE, and 25 for MPE.

B INSTALLATION

In this section, we provide the bash command needed for installing the PyMARLzoo+ environment below in Figure 16, assuming that Python ≥ 3.8 is installed on the machine. This minimal installation step is intended to facilitate a smooth setup process, thereby enabling users to access and utilize the framework’s features efficiently and with ease.

```
$ pip install pymarlzooplus
```

Figure 16: Bash command to install PyMARLzoo+ package using the pip tool.

C TRAINING API

The training framework of (E)PyMARL has been enhanced to support the integration of new algorithms and environments. We introduced a novel category of algorithm modules, termed *explorers*, which integrate the exploration components of the CDS and EOI algorithms. This integration facilitates a structured and transparent code repository that accommodates both general-purpose and exploration-specific algorithms. Furthermore, the widely used *Prioritized Replay Buffer* [44] has been universally integrated into all off-policy algorithms, promoting efficient experimentation without necessitating code modifications. Figure 17 illustrates command examples for training an algorithm (*<algo>*) on each environment task, demonstrating the minimal user effort required. The available options for each *<algo>* and *<key>* and *<time_limit>* are described in Appendix A. Examples of usage of the extra arguments for PettingZoo and Overcooked are provided in Figure 18, in which 10 pistons and sparse reward type has been used for each environment, respectively.

LBF, RWARE, MPE

```
$ python3 src/main.py --config=<algo> --env-config=gymma with env_args.time_limit=<time_limit> env_args.key=<key>
PettingZoo
$ python3 src/main.py --config=<algo> --env-config=pettingzoo with env_args.time_limit=<time_limit> env_args.key=<key> env_args.kwargs=<kwargs>
Overcooked
$ python3 src/main.py --config=<algo> --env-config=overcooked with env_args.time_limit=<time_limit> env_args.key=<key> env_args.reward_type=<reward_type>
PressurePlate
$ python3 src/main.py --config=<algo> --env-config=pressureplate with env_args.key=<key> env_args.time_limit=<time_limit>
Capture Target
$ python3 src/main.py --config=<algo> --env-config=capturagetarget with env_args.key=<key> env_args.time_limit=<time_limit>
Box Pushing
$ python3 src/main.py --config=<algo> --env-config=boxpushing with env_args.key=<key> env_args.time_limit=<time_limit>
```

Figure 17: Example Bash commands for training with different algorithms in different environment tasks. The `<algo>`, and `<key>` correspond to the available options as mentioned in Appendix A. The `<time_limit>` field corresponds to the maximum number of steps before episode termination, therefore it should be an integer.

PettingZoo with 10 pistons in the Pistonball task for training QMIX with time limit 900 steps.

```
$ python3 src/main.py --config=qmix --env-config=pettingzoo with env_args.time_limit=900 env_args.key="pistonball_v6" env_args.kwargs="{'n_pistons':10},"
Overcooked with sparse reward type in Cramped room scenario for training QMIX with time limit 900 steps.
$ python3 src/main.py --config=qmix --env-config=overcooked with env_args.time_limit=900 env_args.key="cramped_room" env_args.reward_type="sparse"
```

Figure 18: Example Bash commands for training QMIX algorithm in PettingZoo and Overcooked environment tasks.

D MULTI-AGENT REINFORCEMENT LEARNING ENVIRONMENTS

D.1 Multi-Agent Particle Environment Details

For our experiments in the MPE environment, we focused on the **Spread** environment, which corresponds to the cooperative navigation environment, as described in the paper by Lowe et al. [32]. An example of the environment task is visualized in Figure 19.

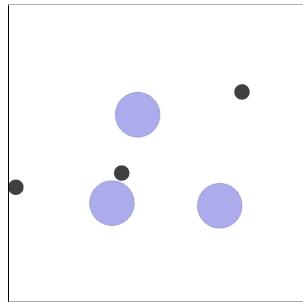


Figure 19: Sample snapshot of the Spread MPE environment.

D.1.1 Observations.

The environment involves N agents and N landmarks, where the agents are tasked with navigating to and covering the landmarks while preventing collisions with each other. Each agent observes its velocity and position, the relative position of other landmarks and agents, and lastly, the communication from all other agents.

$$O^i = (velocity_x, velocity_y, agent_position_x, agent_position_y, N \times (relative_position_x_landmarks, relative_position_y_landmarks), (N - 1) \times (relative_position_x_other_agent, relative_position_y_other_agent), (N + 1) \times one_hot_value)$$

An example of an agent's observation space is provided below for the "SimpleSpread-3-v0" task:

```
array([-0.5          0.           0.8486883   -0.03767432  -0.44009647  0.24901201 ,
       -1.0116129  -0.4476194  -1.7624141   0.89627826  0.09657926  -0.60416883 ,
      -1.5796437  -0.8132916  0.           0.           0.           0.        ], dtype=float32),
```

This example observation vector consists of a vector, with 18 values, as it corresponds to N = 3 agents and landmarks. For N = 4 agents and landmarks, the observation vector would contain 24 values per agent; for N = 5, it would contain 30 values, and so on.

D.1.2 Actions.

The agent's action space includes the five standard actions that an agent i can choose from:

$$A^i = (\text{No action, Move left, Move right, Move down, Move up})$$

D.1.3 Rewards.

Each agent is awarded jointly based on the distance the agent is closest to every landmark. The distance is calculated based on the sum of minimum distances. When agents run into one another, they are penalized locally with -1 for each collision.

D.2 Multi-Robot Warehouse (RWARE) Environment Details

The Multi-Robot Warehouse (RWARE) environment [36], designed based on real-world scenarios, depicts a warehouse that includes robots operating and gathering shelves, transferring them and the ordered items to a specific location, the workstation. Robots move the items back to empty shelf locations after humans have obtained the contents of a shelf. This cooperative environment has different settings that can be configured, like the size, the number of agents, the communication, and reward options. In our experiments, the tiny (10x11), and small (10x20) warehouse sizes were used, leveraging two and four agents on the 'hard' difficulty level. A sample snapshot of the "rware-small-4ag-hard-v1" task is given in Figure 20.

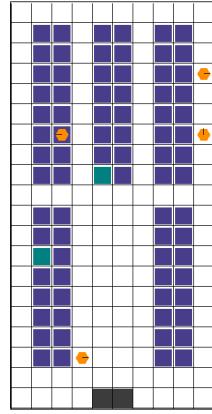


Figure 20: Sample snapshot of small size RWARE environment with four agents.

D.2.1 Observations.

In this environment, each agent observation space is partially observable and includes information about the 3×3 grid, centered around itself, that an agent can observe. The information that is observed incorporates the agent's own state (position, rotation, load) and the relative state of nearby agents and shelves. Specifically, the observation space for each agent i can be described as a vector containing the following elements:

$$O^i = (\text{position_x}, \text{position_y}, \text{carrying_shelf}, \text{direction_one_hot}, \text{path_restricted}, \text{grid_info})$$

where:

- **position_x, position_y:** The agent's current position in x and y coordinates.
- **carrying_shelf:** A binary value (1 or 0) showing whether the agent is carrying a shelf.
- **direction_one_hot:** Four values of the agent's current direction (up, down, left, right) for each item in one-hot encoding.
- **path_restricted:** A binary value (1 or 0) indicating if the agent is on a path where shelves cannot be placed.
- **grid_info:** Information about the 3×3 grid surrounding the agent, split into 9 groups (one for each cell). The visibility range can be adjusted as well. Each group contains 7 elements:
 - A binary value indicating if an agent exists in the cell.
 - Four values in one-hot encoding, representing the direction of any agent in the cell.
 - A binary value showing the presence of a shelf.
 - A binary value describing if the shelf is requested for delivery.

An example of an agent's observation array on the "rware-small-4ag-hard-v1" task with a 3×3 grid is provided below:

```
array([5.  5.  1.  1.  0.  0.  0.  0.  1.  0.  0.  0.  1.  0.  0.  0.  1.  0.  0.  1.  0.  0.  0.  1. ,  
     0.  0.  0.  0.  0.  1.  0.  0.  0.  1.  0.  1.  1.  0.  0.  0.  1.  0.  0.  0.  1.  0.  0.  0. ,  
     0.  0.  0.  1.  0.  0.  0.  1.  0.  0.  0.  0.  1.  0.  0.  0.  1.  0.  0.  0.  0.  0. ], dtype=float32),
```

D.2.2 Actions.

The action space of this environment for each agent is comprised of four discrete possible options. For each agent i , the action space is the following:

$$A^i = (\text{Turn left}, \text{Turn right}, \text{Forward}, \text{Load/Unload shelf})$$

Each agent can only rotate and move forward with the first three actions. The last action, relating to loading and unloading a shelf, can be performed only when a robot is in a specific location, below a shelf.

D.2.3 Rewards.

In every timestep, a fixed number of shelves (R) is requested to be transferred to a target location, and new shelves are randomly selected and included in the current requests. The reward value of an agent is 1 when it correctly delivers the ordered shelf. A difficulty that arises for the agents is to both complete the deliveries and find an empty spot to return the previously delivered shelf. Since multiple actions are needed between the shelf delivery by the agents, the reward signal is sparse.

D.3 PettingZoo Environment Details

For the PettingZoo environments[49], we have mainly used the tasks from Atari and Butterfly environments, the Entombed Cooperative, the Pistonball, and the Cooperative Pong. In the exploration game Entombed Cooperative, the agents must assist each other cooperatively to get through the evolving maze as far as they can. In the cooperative Pistonball physics game, the goal is to activate the pistons (agents) that move vertically so as to navigate the ball to the left wall of the game border. The Cooperative Pong game is a collaborative version of a pong game, where two agents (paddles) work together to continue playing with the ball for the longest period. Snapshot visualizations of these three fully cooperative environment games are illustrated in Figure 21.

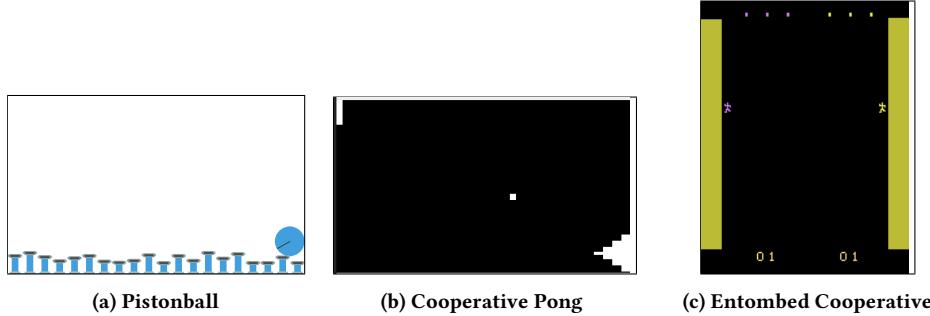


Figure 21: Sample snapshots of the three PettingZoo game environments.

D.3.1 Observations.

Pistonball The agent's observation space is composed of RGB values representing an image that illustrates the area occupied by the two pistons or the wall adjacent to the agent. An example of image-based observation of an agent (piston) can be found in Figure 22 which is composed of RGB values and is of total size (457, 120, 3). In case the image encoding technique is used, as described in the paper, the image observations are transformed into vectors. In that case, an example observation vector of a piston in the environment with 10 agents is given as a vector of size (512,) below:

```
array([7.91885853e-01  0.0000000e+00  7.57015705e-01  7.02866971e-01 ,
       1.37508601e-01  0.0000000e+00  9.95369107e-02  9.06493664e-02,
       ...
       1.01920377e-04  2.64243525e-03  2.42827028e-01  1.12921096e-01,
       1.05303168e+00  3.09428126e-01  1.34141669e-01  2.01561041e-02], dtype=float32)
```

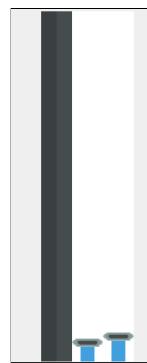


Figure 22: Sample image-based observation of a piston in the Pistonball environment of PettingZoo with 10 agents.

Cooperative Pong The observation space of each agent in this game is the entire screen of the game. An example of image-based observation of an agent (paddle) can be seen in Figure 21b which is composed of RGB values and is of total size (280, 480, 3).

In case the image encoding technique is utilised, the image observations are transformed into vectors. An example observation vector of one of the paddles in the environment is given as a vector of size (512,) below:

```
array([4.37711984e-01  5.60766220e-01  1.59607649e+00  1.11085367e+00,
       3.02959293e-01  3.26711655e-01  5.47930487e-02  9.35874805e-02,
       ...
       3.57522756e-01  3.15122038e-01  3.19070041e-01  4.04478341e-01,
       9.99535978e-01  1.10822819e-01  2.30887890e-01  4.81555223e-01], dtype=float32)
```

Entombed Cooperative In this game, each agent's observations are the RGB values - the entire screen of the game. An example of image-based partial observations of the two agents can be seen in Figure 23 which are composed of RGB values and are of total size (210, 160, 3).

In case the image encoding technique is utilized, the image observations are transformed into vectors. An example observation vector of one of the agents in the environment is given as a vector of size (512,) below:

```
array([1.31586099e+00  4.78404522e-01  1.17896628e+00  4.77674603e-01,
       2.60396451e-01  8.15561116e-02  2.70262837e-01  3.13673377e-01,
       ...
       1.10996671e-01  1.7864831e-01  7.65883446e-01  8.93630207e-01,
       1.32799280e+00  8.94335330e-01  6.74108326e-01  4.32223547e-03], dtype=float32)
```

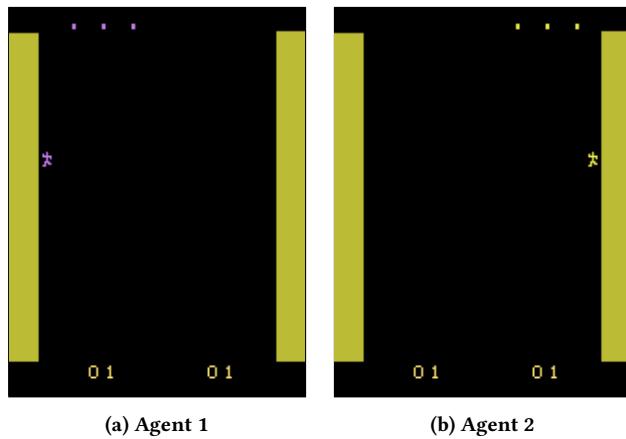


Figure 23: Sample image-based partial observations of each agent in the Entombed Cooperative environment of PettingZoo.

Note that in our enhanced version, we implemented a method to mitigate flickering issues, which involves combining images from sequential frames. Specifically, our wrapper merges images from two consecutive frames to stabilize the representation of agent positions, thereby addressing the flickering problem. This process is repeated for four sequential steps to ensure consistency across frames. Our method identifies the position of Agent 1 in each frame by detecting specific RGB values and determines which frame depicts the agent more prominently. It then creates a combined image where the confirmed position of Agent 1 is reinforced. Additionally, for scenarios requiring partial observations, our wrapper generates separate observations for each agent by removing the visual representation of the other agent from the combined image. This approach not only stabilizes visual inputs but also supports scenarios where agents receive only their individual perspectives, enhancing the task's applicability for studies on cooperative perception and action under partial observability conditions.

D.3.2 Actions.

Pistonball The actions can be given either in discrete or continuous mode. In our experiments, the discrete action space has been used in which three actions are available for each agent i :

$$A^i = (\text{Move down}, \text{Stay still}, \text{Move up})$$

If the actions are continuous, then the range value is $[-1, 1]$, scaled by four, and is linked to how much the pistons are moved up or down. In that way, action 1 will raise a piston by four pixels and -1 lower by four pixels.

Cooperative Pong The discrete action space of each agent is composed by two actions. Each agent can be moved either up or down. Therefore, the action space for each agent i is the following:

$$A^i = (\text{Move up}, \text{Move down})$$

Entombed Cooperative The action space of each agent is composed of 18 options. In an agent's turn, the actions that can be taken are the following:

$$A^i = (\text{No operation}, \text{Fire}, \text{Move up}, \text{Move right}, \text{Move left}, \text{Move down}, \text{Move upright}, \text{Move up left}, \text{Move downright}, \text{Move down left}, \text{Fire up}, \text{Fire right}, \text{Fire left}, \text{Fire down}, \text{Fire up right}, \text{Fire up left}, \text{Fire down right}, \text{Fire down left})$$

D.3.3 Rewards.

Pistonball The reward of each agent is composed of the overall amount of leftward movement of the ball and the leftward movement if it was close to the piston. In that way, each piston can take a global reward, while also a local reward is applied to those agents that are close to the ball, meaning positioned under any area of the ball. The local reward is computed as:

$$\text{local_reward} = 0.5 \times \Delta x_{\text{ball}}$$

where Δx_{ball} represents the change in the ball's x-position. The global reward is calculated as:

$$\text{global_reward} = \left(\frac{\Delta x_{\text{ball}}}{x_{\text{start}}} \right) \times 100 + \text{time_penalty}$$

where x_{start} is the starting position of the ball and time_penalty is a default value of -0.1. Overall, the reward for each agent is calculated as :

$$R^i = \text{local_ratio} * \text{local_reward} + (1 - \text{local_ratio}) * \text{global_reward}.$$

Cooperative Pong At each time step, the reward for each agent i when the ball remains successfully within the bounds is calculated as:

$$R^i = \frac{\text{max_reward}}{\text{max_cycles}}$$

where the default value is set to 0.11. Alternatively, if the ball goes out of bounds, the agent receives a penalty determined by the parameter off_screen_penalty, which by default is set to -10. Upon this event, the game terminates.

Entombed Cooperative The reward, identical to every single agent, in this environment, is given immediately after altering a section or resetting a stage, where each stage is composed of 5 sections. The reward could be received only when an agent loses a life, but not when it loses the last life as the game ends without the option of stage resetting.

D.4 Level-Based Foraging (LBF) Environment Details

The Level-Based Foraging [3] environment is a multi-agent reinforcement learning environment built on OpenAI's RL framework [5]. It features multiple agents operating within a grid-world setting, where their shared objective is to collect as much food as possible. The success of each agent in picking up food is governed by its level of compatibility with the task. The nature of the agents can be either competitive or cooperative in order to reach the goal. Success requires a balance of both to achieve optimal rewards. For our experiments, we analyzed different configurations, varying the amount of players, food, grid size as well as observation radius.

D.4.1 Observations.

In this environment, each agent has to navigate to the food scattered in the grid-world map and collect as many as possible. However, an agent can only pick up one food item, if its level is less than or equal to the combined levels of all collaborating agents. The game is succeeded, if all such food items are collected. Each agent has knowledge either the full state or, if specified, a partial state of the environment. This includes the positions and levels of all entities in the map. The observation array contains N subarrays, each being the observation array of agent i respectfully. The observation array of each agent i contains info on all existing food instances and agents, including themselves. The info for each observed instance is a triplet of the form:

$$O^i = (\text{absolute_position_x}, \text{absolute_position_y}, \text{level})$$

The final observation space for each agent is of size: $(x, y, \text{level}) \times \text{food_count} + (x, y, \text{level}) \times \text{player_count} = 3 \times (\text{food_count} + \text{player_count})$

If a food is picked up or another player is not found yet, the triplet value for that object is set as (-1, -1, 0).

Below is an example of an observation vector of an agent, for the environment "Foraging-4s-11x11-3p-2f-coop-v2":

```
array([ 2.,  0.,  4., -1., -1.,  0.,  2.,  4.,  2.,
       2.,  3.,  1., -1., -1.,  0.], dtype=float32),
```

The three values in the array correspond to the triplet (position_x, position_y, level) of the first food item on the map, as observed by the agent. In this specific configuration, the environment consists of 2 food items and 3 players. Therefore, the first two triplets correspond to the food, while the rest represent the players. The second triplet of the array in particular is out of sight, as indicated by the initialized empty values (-1, -1, 0).

D.4.2 Actions.

For this environment, the agent can perform 6 discrete actions. At each time step, the environment returns an action list, which contains all actions taken by the agents. Each action an agent i can take is one of:

$$A^i = (\text{Noop}, \text{Move north}, \text{Move south}, \text{Move west}, \text{Move east}, \text{Pick-up})$$

From the actions above, "Noop" action signifies the agent remaining idle, while "Pick-up" action allows the agent to collect food from an adjacent cell. The remaining actions result in the agent's moving in the corresponding direction on the 2D grid.

D.4.3 Rewards.

An agent is rewarded only if they have collected at least one food item. When a food item is picked up, all agents adjacent to the food cell are rewarded accordingly, thus underlying the importance of cooperation in the game.

The reward given to each depends on the level of each participant agent so that higher-level agents get higher rewards for the same food picked up. The reward i for each participant agent i is calculated by the level of the food item picked up weighted by the level of the agent. This is then normalized by a factor of the sum of the other adjacent agent levels multiplied by the total number of foods that spawned in the current game. All rewards are normalized so that all individual agent's rewards sum to one.

D.5 Overcooked Environment Details

In our experiments, we evaluated the performance of the proposed algorithms using the Overcooked environment, as described in [7]. The environment is a two-agents cooperative game, where the goal is to prepare and deliver various types of soups as efficiently as possible. The game takes place in a grid-based world, where each cell represents a distinct entity. The challenge is the ability of the agents to coordinate effectively, manage resources and optimize task execution. The agents' actions include navigating the grid, collecting ingredients, preparing the soups, and finally delivering them to a designated service area. The recipes available in the game involve two different ingredients, which can be combined in varying orders depending on the specific soup being prepared.

D.5.1 Observations.

In our experiments, we utilized a variety of environmental layouts to test different cooperation challenges. Specifically, we evaluated agent performance across three distinct layouts: "Cramped room", "Asymmetric advantages", and "Coordination ring". Each room focuses on different cooperation challenges.

For the "**Cramped room**", the agents are positioned within a single shared space. However, the grid layout is highly restrictive, so the environment tests the agents' ability to coordinate in confined spaces.

For the "**Asymmetric advantages**", the layout is divided into separate rooms with each agent positioned in each one. This configuration tests the agents' high-level strategy decision-making ability. Successful coordination in this scenario relies on each agent's capacity to optimize its actions based on the position and current actions of its co-players, highlighting the importance of inter-agent awareness and strategic planning.

In the "**Coordination ring**", similar to the cramped room layout, both agents are placed in a shared space but in a more spacious room. However, the pot (element needed for main interaction) is positioned at the top right corner of the layout, while all other interactions are positioned in the left bottom corner. The agents are tested whether they can coordinate and navigate the environment, ensuring they can execute distant but interconnected actions in a timely and efficient manner.

The three layouts can be seen in Figure 24, referenced by [7].

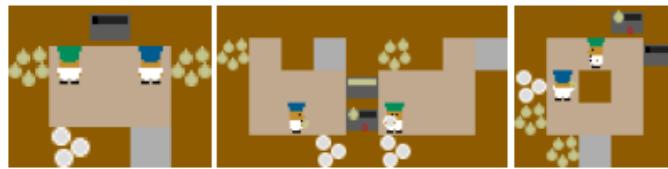


Figure 24: Sample snapshot of the Overcooked layouts. On the left is, "Cramped room", followed by the "Asymmetric advantages" and on the right is the "Coordination ring" layout.

In each timestep, both agents have full knowledge of the state of the environment. Each agent's observation space contains their position as well as the encoding of the next state. An overview of the observation matrix of an agent i is :

$$O^i = (\text{player_i_features}, \text{other_player_features}, \text{player_i_dist_to_other_players}, \text{player_i_position})$$

As observed from the above, their observation matrix contains both their own and their co-players' features, which is mostly their positional information to all objects in their environment. Specifically, according to the source code [13], each player's features is a flattened list which contains:

- orientation: one-hot-encoding of the direction the player is currently facing. Possibilities are north, south, east, and west. This vector length is 4.
- obj: one-hot-encoding of the object currently being held. If nothing held, all 0s. Possibilities are onion, soup, dish, and tomato. This vector length is 4.
- wall_j: 0, 1 boolean value of whether the player has a wall immediately in direction j. This vector length is 4.
- closest_{onion|tomato|dish|soup}serving|empty_counter}: (dx, dy) the distance to the item on both axes. (0, 0) if it is currently held. This vector length is 6 * 2.
- closest_soup_n_{onions|tomatoes}: an int indicating how many of the specified ingredients are inside the closest soup. This vector length is 2.
- closest_pot_j_exists: 0, 1, boolean value if there is a pot in the j direction. If not, then all the next pot features are 0. This vector length is 1.
- closest_pot_j_{is_empty|is_full|is_cooking|is_ready}: 0, 1 boolean value for each pot statuses. 0 if there is no pot in j direction. This vector length is 4.
- closest_pot_j_{num_onions|num_tomatoes}: int value indicating the number of each ingredient in jth closest pot. 0 if there is no pot in j direction. This vector length is 2.
- closest_pot_j_cook_time: int value indicating seconds remaining until soup is ready. -1 if there is no soup cooking. This vector length is 1.
- closest_pot_j: (dx, dy) distance on both axes to the pot in the j direction. (0, 0) of no pot on j direction. This vector length is 2.

Note that the final number of closest_pot features is multiplied by the number of total pots in the environment, which is 2. So the total

$$\begin{aligned}\|O^i\| &= \text{player_i_features} + \text{other_player_features} + \text{player_i_relative_position} + \text{player_i_absolute_position} \\ &= (2 * \text{pot_features} + \text{player_i_other_features}) + (2 * \text{pot_features} + \text{other_player_other_features}) + 2 + 2 \\ &= 2 * 20 + 56 = 96\end{aligned}$$

An example of the total observation array of a sampled time step is given below:

```
array([ 0.,  1.,  0.,  0.,  0.,  0.,  0.,  0., -1., -1.,  0.,  0.,  0.,
       1.,  0.,  0.,  0.,  0.,  2.,  1.,  0.,  0.,  1.,  1.,  0.,  0.,
      0.,  0.,  0.,  0.,  1., -2.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,
      0.,  0.,  0.,  0.,  1.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,
      0.,  0.,  1.,  0.,  0., -2.,  2.,  0.,  0.,  0.,  0.,  0.,  0.,
      2.,  0.,  0.,  1.,  1.,  0.,  0.,  0.,  0.,  0.,  0., -1., -1.,
      0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  1.,
      0.,  2., -1.,  1.,  2.], dtype=float32),
```

This array contains the observation space, which is of size (96), that belongs to a participating agents. All the values of the observation array also correspond to the fields explained earlier, which correspond to their features, the other player's features, their relative position, and lastly their absolute position on the grid world. For example, the first 4 values {0, 1, 0, 0} represent agent's orientation in one-hot encoding. This means the agent is facing south. The next 4 numbers {0, 0, 0, 0} correspond to the one-hot encoding of the object the agent holds, and since all are zero values, the agent does not hold any object. The rest of the values can be interpreted as explained in the individual fields of the player's features. Finally, the last 4 values correspond to the relative and absolute position. So (2, -1) means that the 2nd agent is two positions on the x-axis further, but 1 position on the y-axis below. The last 2 values (1, 2) show that the agent is in position (1, 2) of the grid.

D.5.2 Actions.

At each time step, the actions an agent can take correspond to the available directions they can move to. The agent takes 6 discrete actions as seen below:

$$A^i = (\text{Move north}, \text{Move south}, \text{Move east}, \text{Move west}, \text{Noop}, \text{Interact})$$

The last action is not related to navigation of the agent, but if selected it interacts with a nearby object.

D.5.3 Rewards.

The reward types available are shaped or sparse and are jointly distributed to both agents. In our experiments, we have selected the sparse type, which is given to the agents only when a soup is delivered. Once it is delivered, an agent is rewarded 0 points, if the soup is not made according to the available recipes for the current game. If the recipe of the soup is in the bonus orders, it receives a weighted reward of the base score value from the recipe multiplied by the additional bonus score of 2. Otherwise, the agent is rewarded with the calculated base score of 20. While the option for bonus orders is available, our selected layouts do not include any. So the final reward for each agent is 0 if the soup was not created with the accepted number and order of ingredients, otherwise, the agent is rewarded 20 points.

D.6 Pressureplate Environment Details

In our experiments, we utilized the Pressureplate environment [52], which is a grid-world environment requiring agent cooperation. To progress and succeed, agents must collaborate by ensuring that the agent with the corresponding ID remains on a pressure plate in each

room, which unlocks the door for the remaining agents to move to the next room. This process continues until all agents are on their respective plates, and the final agent reaches the treasure chest, which signifies the goal in the last room. A sample snapshot of the Pressureplate environment can be seen in Figure 25.

D.6.1 *Observations.*

In each new level, the grid map is divided into rooms by strategically placed walls, with each room containing a pressure plate and a corresponding locked door. At the start of each game, agents spawn in the southernmost room, with each assigned a unique ID that matches one of the plates. Agents have partial observability of the environment, depending on its sensor range, specified by the parameter "sight". An example of an agent's observation array at a sampled time-step is given as:

The observation array of each agent contains a combination of 5 flattened observation sub-arrays of dimensions $sight \times sight$ each. Each observation sub-array represents the separate layout of a different element in the environment and the grid that is observable by the agent around them. These elements include the positions of other agents, walls, plates, doors, and goal (treasure chest), giving a total of 5 distinct observation sub-arrays per agent. If an element of the selected layout exists in one of these grid cells, then the value of that cell position on the respective sub-array will be 1, otherwise 0. After the layout sub-arrays are constructed, they are flattened to a final observation array of

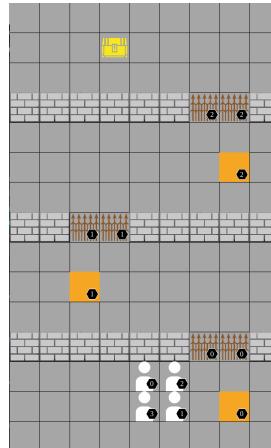


Figure 25: Snapshot of the Pressureplate environment.

the agent. Finally, the agent's position coordinates (x, y) are concatenated in the end, and thus form the final observation array.

The total observation array's dimension is :

$$\|O^i\| = N * (5 * (sight \times sight) + 2)$$

where N is the number of agents.

D.6.2 Actions.

At each time step, each agent i can select from 5 possible discrete actions, corresponding to their navigation on the grid-world map. These possible actions are:

$A^i \equiv (\text{Move up}, \text{Move down}, \text{Move left}, \text{Move right}, \text{Noop})$

From the above, the last action represents idleness while the rest results in moving the agent by one cell in the chosen direction.

D6.3 Rewards

Reward is independent of the actions of other participating agents. At each time step, the distance between each agent's position and their corresponding pressure plate is calculated and their reward is then assigned to them, based on this distance. If the agent is in the same room as their respective plate, the reward will be the negative normalized Manhattan distance between the agent and the plate. However, if the agent is in a different room, the reward is calculated based on the number of rooms separating the agent from the room that contains their plate. Maximum rewards are given when all agents are in the correct room and the final agent has reached the treasure chest.

D.7 Capture Target Environment Details

Agents in the Capture Target environment must learn the grid's transition dynamics to achieve simultaneous target capture, as quickly as possible. Both agents and targets move inside the grid. Agents observe the targets with a flickering probability at each time step. In the default Capture Target setup, which is presented in Figure 26, there are 2 agents and 1 target, moving inside a 6X6 grid. The target's observation flickering probability is set to 0.3.

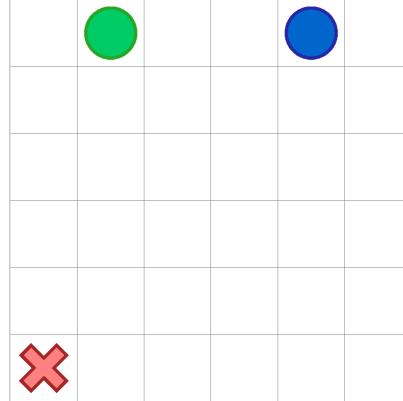


Figure 26: CaptureTarget environment. Two agents learn to capture the target simultaneously.

D.7.1 Observations.

By default, the observations' array consists of the agents' and their targets' normalized positions on the grid. There is also an option for one-hot encoded observations. Targets' positions may differ from the actual ones, due to the flickering probability applied. In the 2 agents, 1 target scenario, the observation array per time step is of size 2x4. For each agent i , the observation vector is

$$O^i = (\text{x_position}, \text{y_position}, \text{target_x_position}, \text{target_y_position})$$

An example of an agent's observation array of a sample time step in a 6x6 grid is given below:

```
array([ 0.4  0.4  0.4  1.6], dtype=float32)
```

D.7.2 Actions.

At each time step an agent can move in one of the available directions. An agent i can pick one of the following actions:

$$A^i = (\text{Move north, Move south, Move west, Move east, No move})$$

D.7.3 Rewards.

In the successful scenario, where the agents capture the targets simultaneously, they are given +1 reward. At any other time step, the reward received is 0.

D.8 Box Pushing Environment Details

The Box Pushing[56] is a cooperative robotics problem introduced by Seuken and Zilberstein [45] in which the main objective is two robots (agents) should cooperate and push a big box to the yellow space to get higher reward than pushing the other two small boxes that exist in the environment. In the default Box Pushing Target setup, which is depicted in Figure 27, two agents, one big box and two small boxes exist and can move in a 6x6 grid.

D.8.1 Observations.

The observation space of each agent consists of 5 elements in which the possible values are 0 and 1 indicating whether the element is in front of an agent (1) or not (0). An overview of the observation vector of an agent i is:

$$O^i = (\text{small_box}, \text{large_box}, \text{empty}, \text{wall}, \text{teammate})$$

An example of an agent's observation array of a sample time step in a 6x6 grid is given below:

```
array([0.  0.  1.  0.  0.], dtype=float32)
```

The values in this observation space show that the agent has an empty cell in front of them.

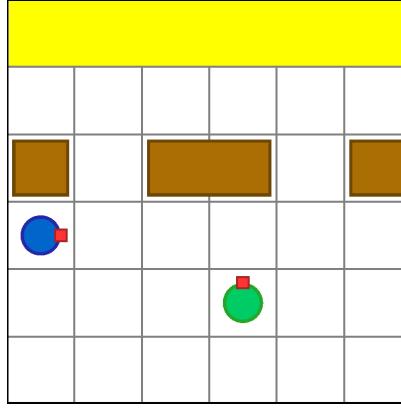


Figure 27: Snapshot of the Box Pushing environment.

D.8.2 Actions.

In every time step, an agent i can choose one of the four available actions:

$$A^i = (\text{Move forward}, \text{Turn left}, \text{Turn right}, \text{Stay})$$

The big box can be shifted when two robots are positioned facing it in adjacent cells and move forward simultaneously. The small boxes can be moved by one grid cell when a robot is positioned toward it and performs the move forward action.

D.8.3 Rewards.

At every time step, the agents are rewarded by -0.1, and a successful push of the large box to the yellow space gives a +100 reward. In case a small box is pushed to the goal area, a +10 reward is returned to the agent. Unsuccessful scenarios that give a -5 penalty, are considered when an agent hits the boundary or when an agent alone tries to push the big box.

E ADDITIONAL RESULTS

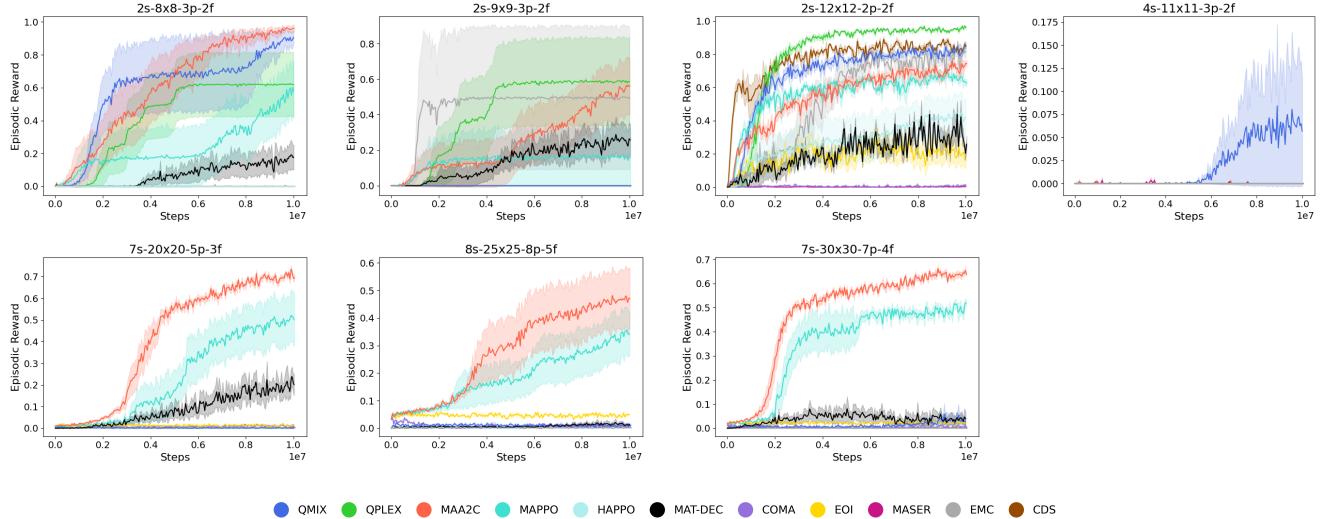


Figure 28: Episodic rewards of all 11 algorithms in LBF tasks rendering the mean and the 75% confidence interval over 5 different seeds.

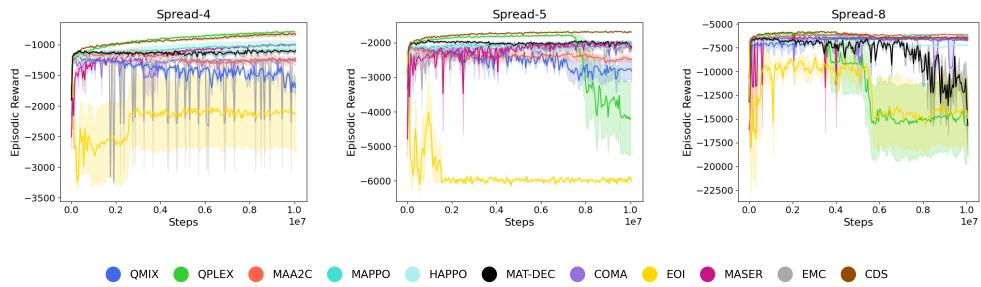


Figure 29: Episodic rewards of all 11 algorithms in MPE tasks rendering the mean and the 75% confidence interval over 5 different seeds.

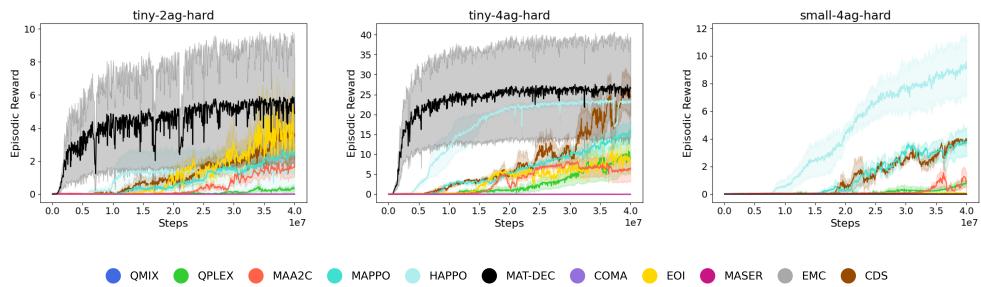


Figure 30: Episodic rewards of all 11 algorithms in RWARE tasks rendering the mean and the 75% confidence interval over 5 different seeds.

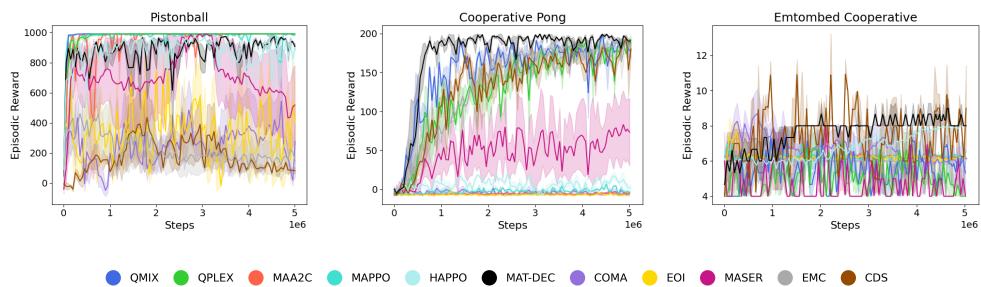


Figure 31: Episodic rewards of all 11 algorithms in PettingZoo environments rendering the mean and the 75% confidence interval over 5 different seeds.

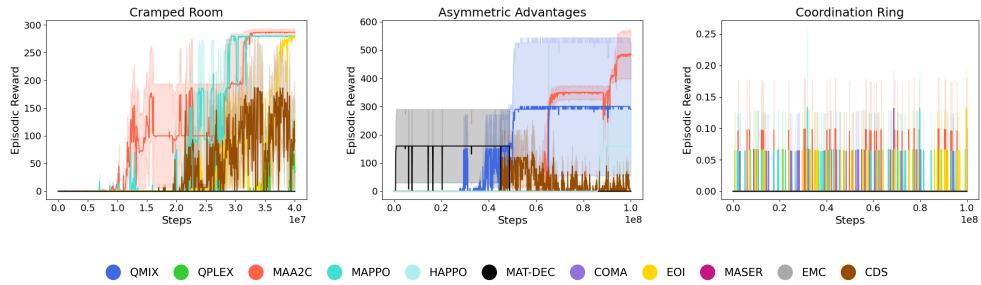


Figure 32: Episodic rewards of all 11 algorithms in Overcooked environments rendering the mean and the 75% confidence interval over 5 different seeds.

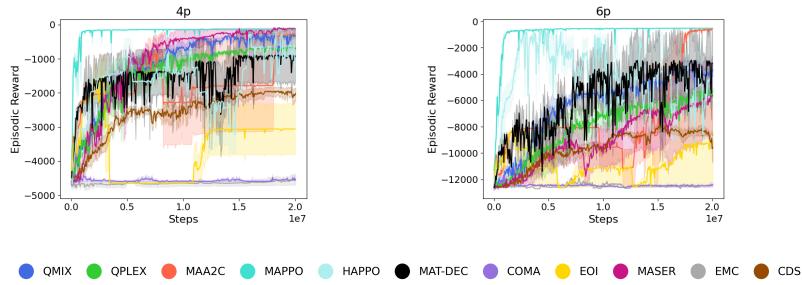


Figure 33: Episodic rewards of all 11 algorithms in Pressure Plate tasks rendering the mean and the 75% confidence interval over 5 different seeds.

Algorithms\Environments	LBF	RWARE	Spread (MPE)	Petting Zoo	Overcooked	Pressure Plate
QMIX	0.29 ± 0.18	0.00 ± 0.00	0.79 ± 0.12	0.40 ± 0.28	0.21 ± 0.19	0.86 ± 0.11
QPLEX	0.32 ± 0.17	0.16 ± 0.12	0.64 ± 0.30	0.40 ± 0.28	0.06 ± 0.06	0.77 ± 0.15
MAA2C	0.64 ± 0.13	0.13 ± 0.06	0.81 ± 0.11	0.34 ± 0.31	0.53 ± 0.27	0.97 ± 0.01
MAPPO	0.46 ± 0.11	0.29 ± 0.15	0.83 ± 0.11	0.34 ± 0.31	0.19 ± 0.18	0.98 ± 0.01
HAPPO	0.08 ± 0.08	0.42 ± 0.22	0.80 ± 0.14	0.34 ± 0.30	0.11 ± 0.10	0.97 ± 0.01
MAT-DEC	0.22 ± 0.08	0.41 ± 0.29	0.80 ± 0.14	0.40 ± 0.28	0.00 ± 0.00	0.85 ± 0.07
COMA	0.02 ± 0.01	0.00 ± 0.00	0.81 ± 0.12	0.23 ± 0.21	0.00 ± 0.00	0.32 ± 0.26
EOI	0.07 ± 0.05	0.23 ± 0.12	0.50 ± 0.25	0.32 ± 0.30	0.19 ± 0.18	0.51 ± 0.20
MASER	0.01 ± 0.00	0.00 ± 0.00	0.82 ± 0.12	0.37 ± 0.30	0.00 ± 0.00	0.79 ± 0.17
EMC	0.23 ± 0.16	0.02 ± 0.02	0.81 ± 0.12	0.16 ± 0.07	0.00 ± 0.00	0.32 ± 0.26
CDS	0.13 ± 0.14	0.42 ± 0.26	0.83 ± 0.13	0.21 ± 0.11	0.18 ± 0.10	0.60 ± 0.20

Table 12: Mean average episodic rewards and 75% confidence interval over each benchmark tasks for each algorithm.

Tasks\Algorithms	QMIX	QPLEX	MAA2C	MAPPO	HAPPO	MAT-DEC	COMA	EOI	MASER	EMC	CDS
2s-8x8-3p-2f	0d : 7h	0d : 11h	0d : 1h	0d : 2h	0d : 4h	0d : 3h	0d : 1h	0d : 6h	0d : 14h	1d : 8h	0d : 13h
2s-9x9-3p-2f	0d : 7h	0d : 11h	0d : 1h	0d : 2h	0d : 4h	0d : 3h	0d : 1h	0d : 6h	0d : 14h	2d : 2h	0d : 14h
2s-12x12-2p-2f	0d : 6h	0d : 11h	0d : 1h	0d : 2h	0d : 3h	0d : 3h	0d : 1h	0d : 5h	0d : 11h	1d : 6h	0d : 10h
4s-11x11-3p-2f	0d : 7h	0d : 12h	0d : 1h	0d : 2h	0d : 4h	0d : 3h	0d : 1h	0d : 6h	0d : 14h	2d : 0h	0d : 13h
7s-20x20-5p-3f	0d : 9h	0d : 14h	0d : 1h	0d : 2h	0d : 6h	0d : 4h	0d : 1h	0d : 6h	0d : 20h	1d : 12h	0d : 20h
8s-25x25-8p-5f	0d : 12h	0d : 20h	0d : 2h	0d : 3h	0d : 10h	0d : 6h	0d : 1h	0d : 8h	1d : 5h	4d : 12h	1d : 7h
7s-30x30-7p-4f	0d : 11h	0d : 18h	0d : 2h	0d : 3h	0d : 9h	0d : 6h	0d : 1h	0d : 7h	1d : 2h	2d : 18h	1d : 3h
Spread 4	0d : 10h	0d : 15h	0d : 2h	0d : 2h	0d : 6h	0d : 4h	0d : 2h	0d : 7h	1d : 1h	3d : 20h	0d : 19h
Spread 5	0d : 12h	0d : 18h	0d : 2h	0d : 3h	0d : 8h	0d : 5h	0d : 2h	0d : 9h	1d : 6h	4d : 5h	1d : 0h
Spread 8	0d : 23h	1d : 8h	0d : 4h	0d : 5h	0d : 14h	0d : 9h	0d : 4h	0d : 18h	2d : 5h	5d : 3h	1d : 19h
tiny-2ag-hard	1d : 14h	2d : 3h	0d : 8h	0d : 11h	0d : 17h	0d : 15h	0d : 8h	1d : 12h	1d : 19h	18d : 11h	2d : 5h
tiny-4ag-hard	2d : 2h	3d : 1h	0d : 10h	0d : 13h	1d : 6h	0d : 23h	0d : 10h	1d : 22h	2d : 12h	19d : 3h	3d : 13h
small-4ag-hard	2d : 2h	3d : 1h	0d : 10h	0d : 14h	1d : 6h	0d : 22h	0d : 11h	1d : 21h	2d : 12h	19d : 15h	3d : 13h
Pistonball	2d : 7h	4d : 4h	4d : 15h	4d : 13h	32d : 6h	3d : 4h	4d : 5h	1d : 9h	3d : 4h	6d : 15h	4d : 9h
Cooperative Pong	1d : 15h	3d : 17h	2d : 20h	1d : 13h	5d : 15h	1d : 0h	2d : 5h	0d : 14h	1d : 5h	2d : 20h	2d : 0h
Entomped Cooperative	1d : 4h	2d : 14h	4d : 10h	4d : 5h	4d : 8h	3d : 6h	4d : 1h	0d : 23h	1d : 18h	2d : 12h	1d : 11h
Cramped Room	1d : 18h	2d : 11h	0d : 10h	0d : 14h	0d : 21h	0d : 18h	0d : 10h	1d : 15h	2d : 0h	13d : 14h	2d : 12h
Asymmetric Advantages	4d : 13h	6d : 16h	1d : 1h	1d : 12h	2d : 6h	2d : 0h	1d : 3h	4d : 7h	5d : 5h	—	6d : 11h
Coordination Ring	4d : 13h	6d : 6h	1d : 0h	1d : 11h	2d : 3h	1d : 22h	1d : 2h	4d : 7h	5d : 22h	—	6d : 7h
4p	0d : 20h	1d : 10h	0d : 4h	0d : 6h	0d : 14h	0d : 10h	0d : 4h	0d : 17h	1d : 2h	9d : 1h	1d : 16h
6p	1d : 1h	1d : 18h	0d : 4h	0d : 8h	0d : 20h	0d : 14h	0d : 5h	0d : 20h	1d : 10h	9d : 13h	2d : 10h

Table 13: Training times of all 11 algorithms in all benchmarks.

F HYPERPARAMETERS

In this section, the hyperparameters used for each algorithm are presented in separate tables. We note that only in *Spread-4* and *Spread-5*, QMIX, MAA2C, COMA and MAPPO use non-sharing policy parameters.

Table 14: Hyperparameters for MAPPO

Name	Value
agent runner	parallel(10)
optimizer	Adam
batch size	10
hidden dimension	64
learning rate	0.0005
reward standardisation	True
network type	GRU
entropy coefficient	0.01
target update	200
buffer size	10
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
n-step	5
epochs	4
clip	0.2

Table 15: Hyperparameters for MAA2C

Name	Value
agent runner	parallel(10)
optimizer	Adam
batch size	10
hidden dimension	64
learning rate	0.0005
reward standardisation	True
network type	GRU
entropy coefficient	0.01
target update	200
buffer size	10
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
n-step	5

Table 16: Hyperparameters for QMIX

Name	Value
agent runner	episode
batch size	32
optimizer	Adam
hidden dimension	64
learning rate	0.0005
reward standardisation	True
network type	GRU
evaluation epsilon	0.0
epsilon anneal	50000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	5000
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
mixing network hidden dimension	32
hypernetwork dimension	64
hypernetwork number of layers	2

Table 17: Hyperparameters for QPLEX

Name	Value
agent runner	episode
optimizer	RMSProp
batch size	32
hidden dimension	64
learning rate	0.0005
reward standardisation	True
network type	GRU
evaluation epsilon	0.0
epsilon anneal	200000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	5000
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
mixing network hidden dimension	32
hypernetwork dimension	64
hypernetwork number of layers	2

Table 18: Hyperparameters for MAT-DEC

Name	Value
agent runner	parallel(10)
optimizer	Adam
batch size	10
hidden dimension	64
learning rate	0.0005
reward standardisation	False
value standardisation	True
actor network type	MLP
critic network type	Transformer
evaluation epsilon	0.0
epsilon anneal	50000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	10
γ (discounted factor)	0.99
entropy coefficient	0.01
value loss coefficient	1
observation agent id	True
observation last action	False
number of attention heads	1
number of blocks	1
epoch	15
clip	0.2
λ_{GAE}	0.95
use Huber loss	True
delta coefficient of Huber loss	10.0
max norm of gradients	10.0
number of mini-batches	1

Table 19: Hyperparameters for CDS

Name	Value
agent runner	episode
optimizer	Adam
batch size	5000
hidden dimension	128
learning rate	0.0005
reward standardisation	True
network type	GRU
evaluation epsilon	0.0
epsilon anneal	50000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	5000
γ (discounted factor)	0.99
observation agent id	False
observation last action	True
n-step	10
entropy coefficient	0.001
number of attention heads	4
attention regulation coefficient	0.001
mixing network hidden dimension	32
hypernetwork dimension	64
hypernetwork number of layers	2
β_1	0.5
β_2	1.0
β	0.1
α	0.6
λ	0.1

Table 20: Hyperparameters for COMA

Name	Value
agent runner	parallel(10)
optimizer	Adam
batch size	10
hidden dimension	64
learning rate	0.0003
reward standardisation	True
network type	GRU
entropy coefficient	0.001
target update	200
buffer size	10
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
n-step	10

Table 21: Hyperparameters for EOI

Name	Value
agent runner	episode
optimizer	Adam
batch size	10
hidden dimension	128
learning rate	0.0005
reward standardisation	True
network type	GRU
target update	200
buffer size	10
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
n-step	5
entropy coefficient	0.01
classifier (ϕ) learning	0.0001
classifier (ϕ) batch size	256
classifier (ϕ) β_2	0.1

Table 22: Hyperparameters for EMC

Name	Value
agent runner	episode
optimizer	RMSProp
batch size	32
hidden dimension	64
learning rate	0.0005
reward standardisation	True
network type	GRU
evaluation epsilon	0.0
epsilon anneal	50000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	5000
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
episodic memory capacity	1000000
episodic latent dimension	4
soft update weight	0.005
weighting term λ of episodic loss	0.1
curiosity decay rate (η_t)	0.9
number of attention heads	4
attention regulation coefficient	0.001
mixing network hidden dimension	32
hypernetwork dimension	64
hypernetwork number of layers	2

Table 23: Hyperparameters for MASER

Name	Value
agent runner	episode
optimizer	RMSProp
batch size	1
hidden dimension	64
learning rate	0.0005
reward standardisation	False
network type	GRU
evaluation epsilon	0.0
epsilon anneal	50000
epsilon start	1.0
epsilon finish	0.05
target update	200
buffer size	5000
γ (discounted factor)	0.99
observation agent id	True
observation last action	True
mixing network hidden dimension	32
representation network dimension	128
α	0.5
λ	0.03
λ_I	0.0008
λ_E	0.00006
λ_D	0.00014

Table 24: Hyperparameters for HAPPO

Name	Value
agent runner	parallel(10)
optimizer	Adam
batch size	10
hidden dimension	64
learning rate	0.0005
reward standardisation	False
value standardisation	True
network type	GRU
entropy coefficient	0.01
buffer size	10
γ (discounted factor)	0.99
observation agent id	False
observation last action	False
epochs	5
n-step	1
max norm of gradients	10
number of mini-batches	1
λ_{GAE}	0.95
clip	0.2
use huber loss	True
delta coefficient of huber loss	10.0