
Benchmarking Multi-Agent Deep Reinforcement Learning Algorithms in Cooperative Tasks

Georgios Papoudakis *
School of Informatics
University of Edinburgh
g.papoudakis@ed.ac.uk

Filippos Christianos *
School of Informatics
University of Edinburgh
f.christianos@ed.ac.uk

Lukas Schäfer
School of Informatics
University of Edinburgh
l.schaefer@ed.ac.uk

Stefano V. Albrecht
School of Informatics
University of Edinburgh
s.albrecht@ed.ac.uk

Abstract

Multi-agent deep reinforcement learning (MRL) suffers from a lack of commonly-used evaluation tasks and criteria, making comparisons between approaches difficult. In this work, we provide a systematic evaluation and comparison of three different classes of MRL algorithms (independent learning, centralised multi-agent policy gradient, value decomposition) in a diverse range of cooperative multi-agent learning tasks. Our experiments serve as a reference for the expected performance of algorithms across different learning tasks, and we provide insights regarding the effectiveness of different learning approaches. We open-source EPy-MRL, which extends the PyMRL codebase to include additional algorithms and allow for flexible configuration of algorithm implementation details such as parameter sharing. Finally, we open-source two environments for multi-agent research which focus on coordination under sparse rewards.

1 Introduction

Multi-agent reinforcement learning (MRL) algorithms use RL techniques to co-train a set of agents in a multi-agent system. Recent years have seen a plethora of new MRL algorithms which integrate deep learning techniques [Papoudakis et al., 2019, Hernandez-Leal et al., 2019]. However, comparison of MRL algorithms is difficult due to a lack of established benchmark tasks, evaluation protocols, and metrics. While several comparative studies exist for single-agent RL [Duan et al., 2016, Henderson et al., 2018, Wang et al., 2019], we are unaware of such comparative studies for recent MRL algorithms. Albrecht and Ramamoorthy [2012] compare several MRL algorithms but focus on the application of classic (non-deep) approaches in simple matrix games. Such comparisons are crucial in order to understand the relative strengths and limitations of algorithms, which may guide practical considerations and future research.

We contribute a **comprehensive empirical comparison of nine MRL algorithms in a diverse set of cooperative multi-agent tasks**. We compare three classes of MRL algorithms: independent learning, which applies single-agent RL algorithms for each agent without consideration of the multi-agent structure [Tan, 1993]; centralised multi-agent policy gradient [Lowe et al., 2017, Foerster et al., 2018, Yu et al., 2021]; and value decomposition [Sunehag et al., 2018, Rashid et al., 2018] algorithms. The two latter classes of algorithms follow the Centralised Training Decentralised

*Equal Contribution

Execution (CTDE) paradigm. These algorithm classes are frequently used in the literature either as baselines or building blocks for more complex algorithms [He et al., 2016, Sukhbaatar et al., 2016, Foerster et al., 2016, Raileanu et al., 2018, Jaques et al., 2019, Iqbal and Sha, 2019, Du et al., 2019, Ryu et al., 2020]. We evaluate algorithms in two matrix games and four multi-agent environments, in which we define a total of 25 different cooperative learning tasks. Hyperparameters of each algorithm are optimised separately in each environment using a grid-search, and we report the maximum and average evaluation returns during training. We run experiments with shared and non-shared parameters between agents, a common implementation detail in MARL that has been shown to affect converged returns [Christianos et al., 2021]. In addition to reporting detailed benchmark results, we analyse and discuss insights regarding the effectiveness of different learning approaches.

To facilitate our comparative evaluation, we created the open-source codebase **EPyMARL** (Extended PyMARL)², an extension of PyMARL [Samvelyan et al., 2019] which is commonly used in MARL research. EPyMARL implements additional algorithms and allows for flexible configuration of different implementation details, such as whether or not agents share network parameters. Moreover, we have implemented and open-sourced **two new multi-agent environments**: Level-Based Foraging (LBF) and Multi-Robot Warehouse (RWARE). With these environments we aim to test the algorithms’ ability to learn coordination tasks under sparse rewards and partial observability.

2 Algorithms

2.1 Independent Learning (IL)

For IL, each agent is learning independently and perceives the other agents as part of the environment.

IQL: In Independent Q-Learning (IQL) [Tan, 1993], each agent has a decentralised state-action value function that is conditioned only on the local history of observations and actions of each agent. Each agent receives its local history of observations and updates the parameters of the Q-value network [Mnih et al., 2015] by minimising the standard Q-learning loss [Watkins and Dayan, 1992].

IA2C: Independent synchronous Advantage Actor-Critic (IA2C) is a variant of the commonly-used A2C algorithm [Mnih et al., 2016, Dhariwal et al., 2017] for decentralised training in multi-agent systems. Each agent has its own actor to approximate the policy and critic network to approximate the value-function. Both actor and critic are trained, conditioned on the history of local observations, actions and rewards the agent perceives, to minimise the A2C loss.

IPPO: Independent Proximal Policy Optimisation (IPPO) is a variant of the commonly-used PPO algorithm [Schulman et al., 2017] for decentralised training in multi-agent systems. The architecture of IPPO is identical to IA2C. The main difference between PPO and A2C is that PPO uses a surrogate objective which constrains the relative change of the policy at each update, allowing for more update epochs using the same batch of trajectories. In contrast to PPO, A2C can only perform one update epoch per batch of trajectories to ensure that the training batch remains on-policy.

2.2 Centralised Training Decentralised Execution (CTDE)

In contrast to IL, CTDE allows sharing of information during training, while policies are only conditioned on the agents’ local observations enabling decentralised execution.

Centralised policy gradient methods One category of CTDE algorithms are centralised policy gradient methods in which each agent consists of a decentralised actor and a centralised critic, which is optimised based on shared information between the agents.

MADDPG: Multi-Agent DDPG (MADDPG) [Lowe et al., 2017] is a variation of the DDPG algorithm [Lillicrap et al., 2015] for MARL. The actor is conditioned on the history of local observations, while critic is trained on the joint observation and action to approximate the joint state-action value function. Each agent individually minimises the deterministic policy gradient loss [Silver et al., 2014]. A common assumption of DDPG (and thus MADDPG) is differentiability of actions with respect to the parameters of the actor, so the action space must be continuous. Lowe et al. [2017] apply the Gumbel-Softmax trick [Jang et al., 2017, Maddison et al., 2017] to learn in discrete action spaces.

²<https://github.com/uoeeagents/epymar1>

Table 1: Overview of algorithms and their properties.

	Centr. Training	Off-/On-policy	Value-based	Policy-based
IQL	✗	Off	✓	✗
IA2C	✗	On	✓	✓
IPPO	✗	On	✓	✓
MADDPG	✓	Off	✓	✓
COMA	✓	On	✓	✓
MAA2C	✓	On	✓	✓
MAPPO	✓	On	✓	✓
VDN	✓	Off	✓	✗
QMIX	✓	Off	✓	✗

COMA: In Counterfactual Multi-Agent (COMA) Policy Gradient, Foerster et al. [2018] propose a modification of the advantage in the actor’s loss computation to perform counterfactual reasoning for credit assignment in cooperative MARL. The advantage is defined as the discrepancy between the state-action value of the followed joint action and a counterfactual baseline. The latter is given by the expected value of each agent following its current policy while the actions of other agents are fixed. The standard policy loss with this modified advantage is used to train the actor and the critic is trained using the TD-lambda algorithm [Sutton, 1988].

MAA2C: Multi-Agent A2C (MAA2C) is an actor-critic algorithm in which the critic learns a joint state value function (in contrast, the critics in MADDPG and COMA are also conditioned on actions). It extends the existing on-policy actor-critic algorithm A2C by applying centralised critics conditioned on the state of the environment rather than the individual history of observations. It is often used as a baseline in MARL research and is sometimes referred to as Central-V, because it computes a centralised state value function. However, MAPPO also computes a centralised state value function, and in order to avoid confusion we refer to this algorithm as MAA2C.

MAPPO: Multi-Agent PPO (MAPPO) [Yu et al., 2021] is an actor-critic algorithm (extension of IPPO) in which the critic learns a joint state value function, similarly to MAA2C. In contrast to MAA2C, which can only perform one update epoch per training batch, MAPPO can utilise the same training batch of trajectories to perform several update epochs.

Value Decomposition Another recent CTDE research direction is the decomposition of the joint state-action value function into individual state-action value functions.

VDN: Value Decomposition Networks (VDN) [Sunehag et al., 2018] aim to learn a linear decomposition of the joint Q-value. Each agent maintains a network to approximate its own state-action values. VDN decomposes the joint Q-value into the sum of individual Q-values. The joint state-action value function is trained using the standard DQN algorithm [Watkins and Dayan, 1992, Mnih et al., 2015]. During training, gradients of the joint TD loss flow backwards to the network of each agent.

QMIX: QMIX [Rashid et al., 2018] extends VDN to address a broader class of environments. To represent a more complex decomposition, a parameterised mixing network is introduced to compute the joint Q-value based on each agent’s individual state-action value function. A requirement of the mixing function is that the optimal joint action, which maximises the joint Q-value, is the same as the combination of the individual actions maximising the Q-values of each agent. QMIX is trained to minimise the DQN loss and the gradient is backpropagated to the individual Q-values.

3 Multi-Agent Environments

We evaluate the algorithms in two finitely repeated matrix games and four multi-agent environments within which we define a total of 25 different learning tasks. All tasks are fully-cooperative, i.e. all agents receive identical reward signals. These tasks range over various properties including the degree of observability (whether agents can see the full environment state or only parts of it), reward density (receiving frequent/dense vs infrequent/sparse non-zero rewards), and the number of agents involved. Table 2 lists environments with properties, and we give brief descriptions below. We

Table 2: Overview of environments and properties.

	Observability	Rew. Sparsity	Agents	Main Difficulty
Matrix Games	Full	Dense	2	Sub-optimal equilibria
MPE	Partial / Full	Dense	2-3	Non-stationarity
SMAC	Partial	Dense	2-10	Large action space
LBF	Partial / Full	Sparse ³	2-4	Coordination
RWARE	Partial	Sparse	2-4	Sparse reward

believe each of the following environments addresses a specific challenge of MARL. Full details of environments and learning tasks are provided in Appendix C.

3.1 Repeated Matrix Games

We consider two cooperative matrix games proposed by Claus and Boutilier [1998]: the *climbing* and *penalty* game. The common-payoff matrices of the climbing and penalty game, respectively, are:

$$\begin{bmatrix} 0 & 6 & 5 \\ -30 & 7 & 0 \\ 11 & -30 & 0 \end{bmatrix} \quad \begin{bmatrix} k & 0 & 10 \\ 0 & 2 & 0 \\ 10 & 0 & k \end{bmatrix}$$

where $k \leq 0$ is a penalty term. We evaluate in the penalty game for $k \in \{-100, -75, -50, -25, 0\}$. The difficulty of this game strongly correlates with k : the smaller k , the harder it becomes to identify the optimal policy due to the growing risk of penalty k . Both games are applied as repeated matrix games with an episode length of 25 and agents are given constant observations at each timestep. These matrix games are challenging due to the existence of local minima in the form of sub-optimal Nash equilibria [Nash, 1951]. Slight deviations from optimal policies by one of the agents can result in significant penalties, so agents might get stuck in risk-free (deviations from any agent does not significantly impede payoff) local optima.

3.2 Multi-Agent Particle Environment

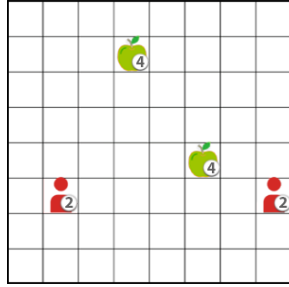
The Multi-Agent Particle Environments (MPE) [Mordatch and Abbeel, 2017] consists of several two-dimensional navigation tasks. We investigate four tasks that emphasise coordination: Speaker-Listener, Spread, Adversary⁴, and Predator-Prey⁴. Agent observations consist of high-level feature vectors including relative agent and landmark locations. The actions allow for two-dimensional navigation. All tasks but Speaker-Listener, which also requires binary communication, are fully observable. MPE tasks serve as a benchmark for agent coordination and their ability to deal with non-stationarity [Papoudakis et al., 2019] due to significant dependency of the reward with respect to joint actions. Individual agents not coordinating effectively can severely reduce received rewards.

3.3 StarCraft Multi-Agent Challenge

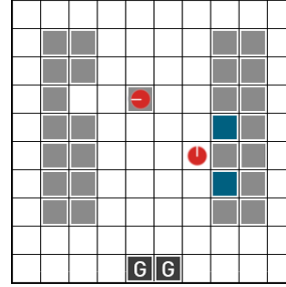
The StarCraft Multi-Agent Challenge (SMAC) [Samvelyan et al., 2019] simulates battle scenarios in which a team of controlled agents must destroy an enemy team using fixed policies. Agents observe other units within a fixed radius, and can move around and select enemies to attack. We consider five tasks in this environment which vary in the number and types of units controlled by agents. The primary challenge within these tasks is the agents’ ability to accurately estimate the value of the current state under partial observability and a growing number of agents of diverse types across tasks. Latter leads to large action spaces for agents which are able to select other agents or enemy units as targets for healing or attack actions, respectively, depending on the controlled unit.

³Rewards in LBF are sparser compared to MPE and SMAC, but not as sparse as in RWARE.

⁴Adversary and Predator-Prey are originally competitive tasks. The agents controlling the adversary and prey, respectively, are controlled by a pretrained policy obtained by training all agents with the MADDPG algorithm for 25000 episodes (see Appendix C for details).



(a) Level-Based Foraging (LBF)



(b) Multi-Robot Warehouse (RWARE)

Figure 1: Illustrations of the open-sourced multi-agent environments [Christianos et al., 2020].

3.4 Level-Based Foraging

In Level-Based Foraging (LBF) [Albrecht and Ramamoorthy, 2013, Albrecht and Stone, 2017] agents must collect food items which are scattered randomly in a grid-world. Agents and items are assigned levels, such that a group of one or more agents can collect an item if the sum of their levels is greater or equal to the item’s level. Agents can move in four directions, and have an action that attempts to load an adjacent item (the action will succeed depending on the levels of agents attempting to load the particular item). LBF allows for many different tasks to be configured, including partial observability or a highly cooperative task where all agents must simultaneously participate to collect the items. We define seven distinct tasks with a variable world size, number of agents, observability, and cooperation settings indicating whether all agents are required to load a food item or not. We implemented the LBF environment which is publicly available on GitHub, under the MIT licence: <https://github.com/ue-agents/lb-foraging>.

3.5 Multi-Robot Warehouse

The Multi-Robot Warehouse environment (RWARE) represents a cooperative, partially-observable environment with sparse rewards. RWARE simulates a grid-world warehouse in which agents (robots) must locate and deliver requested shelves to workstations and return them after delivery. Agents are only rewarded for completely delivering requested shelves and observe a 3×3 grid containing information about the surrounding agents and shelves. The agents can move forward, rotate in either direction, and load/unload a shelf. We define three tasks which vary in world size, number of agents and shelf requests. The sparsity of the rewards makes this a hard environment, since agents must correctly complete a series of actions before receiving any reward. Additionally, observations are sparse and high-dimensional compared to the other environments. RWARE is the second environment we designed and open-source under the MIT licence: <https://github.com/ue-agents/robotic-warehouse>.

We have developed and plan to maintain the LBF and RWARE environments as part of this work. They have already been used in other multi-agent research [Christianos et al., 2020, 2021, Rahman et al., 2021, Papoudakis et al., 2021]. For more information including installation instructions, interface guides with code snippets and detailed descriptions, see Appendix A.

4 Evaluation

4.1 Evaluation Protocol

To account for the improved sample efficiency of off-policy over on-policy algorithms and to allow for fair comparisons, we train on-policy algorithms for a factor of ten more samples than off-policy algorithms. In MPE and LBF we train on-policy algorithms for 20 million timesteps and off-policy algorithms for two million timesteps, while in SMAC and RWARE, we train on-policy and off-policy algorithms for 40 and four million timesteps, respectively. By not reusing samples through an experience replay buffer, on-policy algorithms are less sample efficient, but not generally slower (if the simulator is reasonably fast) and thus this empirical adjustment is fair. We perform in total 41 evaluations of each algorithm at constant timestep intervals during training, and at each evaluation

point we evaluate for 100 episodes. In matrix games, we train off-policy algorithms for 250 thousand timesteps and on-policy algorithms for 2.5 million timesteps and evaluate every 2.5 thousand and 25 thousand timesteps, respectively, for a total of 100 evaluations.

4.2 Parameter Sharing

Two common configurations for training deep MARL algorithms are: without and with parameter sharing. Without parameter sharing, each agent uses its own set of parameters for its networks. Under parameter sharing, all agents share the same set of parameters for their networks. In the case of parameter sharing, the policy and critic (if there is one) additionally receive the identity of each agent as a one-hot vector input. This allows for each agent to develop a different behaviour. The loss is computed over all agents and used to optimise the shared parameters. In the case of varying input sizes across agents, inputs are zero-padded to ensure identical input dimensionality. Similarly if agents have varying numbers of actions, action selection probabilities for invalid actions are set to 0.

4.3 Hyperparameter Optimisation

Hyperparameter optimisation was performed for each algorithm separately in each environment. From each environment, we selected one task and optimised the hyperparameters of all algorithms in this task. In the MPE environment, we perform the hyperparameter optimisation in the Speaker-Listener task, in the SMAC environment in the “3s5z” task, in the LBF environment in the “15x15-3p-5f” task, and in the RWARE environment in the “Tiny 4p” task. We train each combination of hyperparameters using three different seeds and compare the maximum evaluation returns. The best performing combination on each task is used for all tasks in the respective environment for the final experiments. In Appendix I, we present the hyperparameters that were used in each environment and algorithm.

4.4 Performance Metrics

Maximum returns: For each algorithm, we identify the evaluation timestep during training in which the algorithm achieves the highest average evaluation returns across five random seeds. We report the average returns and the 95% confidence interval across five seeds from this evaluation timestep.

Average returns: We also report the average returns achieved throughout all evaluations during training. Due to this metric being computed over all evaluations executed during training, it considers learning speed besides final achieved returns.

4.5 Computational Requirements

All experiments presented in this work were executed purely on CPUs. The experiments were executed in compute clusters that consist of several nodes. The main types of CPU models that were used for this work are Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz and AMD EPYC 7502 32-Core processors. All but the SMAC experiments were executed using a single CPU core. All SMAC experiments were executed using 5 CPU cores. The total number of CPU hours that were spent for executing the experiments in this work (excluding the hyperparameter search) are 138,916.

4.6 Extended PyMARL

Implementation details in reinforcement learning significantly affect the returns that each algorithm achieves [Andrychowicz et al., 2021]. To enable consistent evaluation of MARL algorithms, we open-source the Extended PyMARL (EPyMARL) codebase. EPyMARL is an extension of the PyMARL codebase [Samvelyan et al., 2019]. PyMARL provides implementations for IQL, COMA, VDN and QMIX. We increase the scope of the codebase to include five additional policy gradients algorithms: IA2C, IPPO, MADDPG, MAA2C and MAPPO. The original PyMARL codebase implementation assumes that agents share parameters and that all the agents’ observation have the same shape. In general, parameter sharing is a commonly applied technique in MARL. However, it was shown that parameter sharing can act as an information bottleneck, especially in environments with heterogeneous agents [Christianos et al., 2021]. EPyMARL allows training MARL algorithms without parameter sharing, training agents with observations of varying dimensionality, and tuning several implementation details such as reward standardisation, entropy regularisation, and the use of

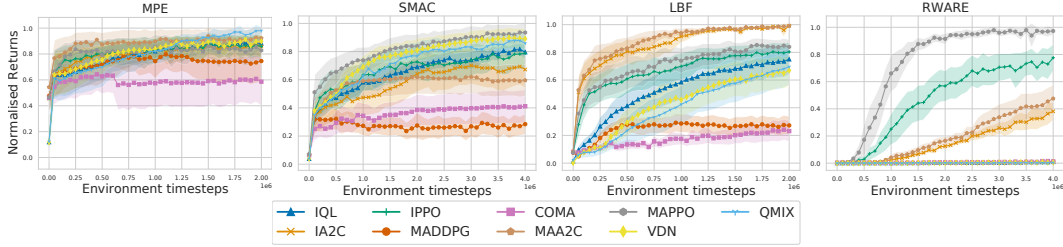


Figure 2: Normalised evaluation returns averaged over the tasks in the all environments except matrix games. Shaded part represents the 95% confidence interval.

Table 3: Maximum returns and 95% confidence interval over five seeds for all nine algorithms with parameter sharing in all 25 tasks. The highest value in each task is presented in bold. Asterisks denote the algorithms that are not significantly different from the best performing algorithm in each task.

Tasks \ Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX
Matrix Games									
Climbing	195.00 ± 67.82	175.00 ± 0.00	175.00 ± 0.00	170.00 ± 10.00	185.00 ± 48.99	175.00 ± 0.00	175.00 ± 0.00	175.00 ± 54.77	175.00 ± 54.77
Penalty k=0	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	249.98 ± 0.04	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00
Penalty k=25	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.97 ± 0.02	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00
Penalty k=50	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.98 ± 0.02	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00
Penalty k=75	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.97 ± 0.02	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00
Penalty k=100	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.97 ± 0.03	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00
MPE									
Speaker-Listener	-18.36 ± 4.67	-12.60 ± 3.62 *	-13.10 ± 3.50	-13.56 ± 1.73	-30.40 ± 5.18	-10.71 ± 0.38 *	-10.68 ± 0.30	-15.95 ± 2.48	-11.56 ± 0.53
Spread	-132.63 ± 2.22	-134.43 ± 1.15	-133.86 ± 3.67	-141.70 ± 1.74	-204.31 ± 6.30	-129.90 ± 1.63 *	-133.54 ± 3.08	-131.03 ± 1.85	-126.62 ± 2.96
Adversary	9.38 ± 0.91	12.12 ± 0.44 *	12.17 ± 0.32	8.97 ± 0.89	8.05 ± 0.89	12.06 ± 0.45 *	11.30 ± 0.38	9.28 ± 0.90	9.67 ± 0.66
Tag	22.18 ± 2.83	17.44 ± 1.31	19.44 ± 2.94	12.50 ± 6.30	8.72 ± 4.42	19.95 ± 7.15 *	18.52 ± 5.64	24.50 ± 2.19	31.18 ± 3.81
SMAC									
2s_vs_1sc	16.72 ± 0.38	20.24 ± 0.00	20.24 ± 0.01	13.14 ± 2.01	11.04 ± 7.21	20.20 ± 0.05 *	20.25 ± 0.00	18.04 ± 0.33	19.01 ± 0.40
3s5z	16.44 ± 0.15	18.56 ± 1.31 *	13.36 ± 2.08	12.04 ± 0.82	18.90 ± 1.01 *	19.95 ± 0.05 *	20.39 ± 1.14	19.57 ± 0.20 *	19.66 ± 0.14 *
corridor	15.72 ± 1.77	18.59 ± 0.62	17.97 ± 3.44 *	5.85 ± 0.58	7.75 ± 0.19	8.97 ± 0.29	17.14 ± 4.39 *	15.25 ± 4.18 *	16.45 ± 3.54 *
MMM2	13.69 ± 1.02	10.70 ± 2.77	11.37 ± 1.15	3.96 ± 0.32	6.95 ± 0.27	10.37 ± 1.95	17.78 ± 0.44	18.49 ± 0.31	18.40 ± 0.24 *
3s_vs_5z	21.15 ± 0.41	4.42 ± 0.02	19.36 ± 6.15 *	5.99 ± 0.58	3.23 ± 0.05	6.68 ± 0.55	18.17 ± 4.17 *	19.03 ± 5.77 *	16.04 ± 2.87
LBF									
8x8-2p-2f-c	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.46 ± 0.02	0.61 ± 0.30	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.96 ± 0.07 *
8x8-2p-2f-2s-c	1.00 ± 0.00	1.00 ± 0.00	0.78 ± 0.05	0.70 ± 0.04	0.45 ± 0.15	1.00 ± 0.00	0.85 ± 0.06	1.00 ± 0.00	1.00 ± 0.00
10x10-3p-3f	0.93 ± 0.02	1.00 ± 0.00	0.98 ± 0.01	0.24 ± 0.04	0.19 ± 0.06	1.00 ± 0.00	0.99 ± 0.01	0.84 ± 0.08	0.84 ± 0.08
10x10-3p-3f-2s	0.86 ± 0.01	0.94 ± 0.03 *	0.70 ± 0.03	0.41 ± 0.03	0.29 ± 0.12	0.96 ± 0.02	0.72 ± 0.03	0.90 ± 0.03	0.90 ± 0.01
15x15-3p-5f	0.17 ± 0.08	0.89 ± 0.04	0.77 ± 0.08	0.10 ± 0.02	0.08 ± 0.04	0.87 ± 0.06 *	0.77 ± 0.02	0.15 ± 0.02	0.69 ± 0.04
15x15-4p-3f	0.54 ± 0.18	0.99 ± 0.01 *	0.98 ± 0.01	0.17 ± 0.03	0.17 ± 0.04	1.00 ± 0.00	0.96 ± 0.02	0.38 ± 0.13	0.15 ± 0.06
15x15-4p-5f	0.22 ± 0.04	0.93 ± 0.03 *	0.67 ± 0.22	0.12 ± 0.06	0.12 ± 0.06	0.95 ± 0.01	0.70 ± 0.25 *	0.30 ± 0.04	0.25 ± 0.09
RWARE									
Tiny 4p	0.72 ± 0.37	26.34 ± 4.60	31.82 ± 10.71	0.54 ± 0.10	1.16 ± 0.15	32.50 ± 9.79	49.42 ± 1.22	0.80 ± 0.28	0.30 ± 0.19
Small 4p	0.14 ± 0.28	6.54 ± 1.15	19.78 ± 3.12	0.18 ± 0.12	0.16 ± 0.16	10.30 ± 1.48	27.00 ± 1.80	0.18 ± 0.27	0.06 ± 0.08
Tiny 2p	0.28 ± 0.38	8.18 ± 1.25	20.22 ± 1.76 *	0.44 ± 0.34	0.48 ± 0.34	8.38 ± 2.59	21.16 ± 1.50	0.12 ± 0.07	0.14 ± 0.19

recurrent or fully-connected networks. EPyMARL is publicly available on GitHub and distributed under the Apache License: <https://github.com/uoε-agents/epymarl>.

5 Results

In this section we compile the results across all environments and algorithms. Figure 2 presents the normalised evaluation returns in all environments, except matrix games. We normalise the returns of all algorithms in each task in the $[0, 1]$ range using the following formula: $\text{norm}_G^a = (G_t^a - \min(G_t)) / (\max(G_t) - \min(G_t))$, where G_t^a is the return of algorithm a in task t , and G_t is the returns of all algorithms in task t . Table 3 presents the maximum returns for the nine algorithms in all 25 tasks with parameter sharing. The maximum returns without parameter sharing, as well as the average returns both with and without parameter sharing are presented in Appendix E. In Table 3, we highlight the highest mean in bold. We performed two-sided t-tests with a significance threshold of 0.05 between the highest performing algorithm and each other algorithm in each task. If an algorithm’s performance was *not* statistically significantly different from the best algorithm, the respective value is annotated with an asterisk (i.e. bold or asterisks in the table show the best performing algorithms per task). In the SMAC tasks, it is a common practice in the literature to report the win-rate as a percentage and not the achieved returns. However, we found it more informative to report the achieved returns since it is the metric that the algorithms aim to optimise. Moreover, higher returns do not always correspond to higher win-rates which can make the interpretation of the performance metrics more difficult. For completeness, we report the win-rates achieved by all algorithms in Appendix G.

5.1 Independent Learning

We find that IL algorithms perform adequately in all tasks despite their simplicity. However, performance of IL is limited in partially observable SMAC and RWARE tasks, compared to their CTDE counterparts, due to IL algorithms’ inability to reason over joint information of agents.

IQL: IQL performs significantly worse than the other IL algorithms in the partially-observable Speaker-Listener task and in all RWARE tasks. IQL is particularly effective in all but three LBF tasks, where relatively larger grid-worlds are used. IQL achieves the best performance among all algorithms in the “3s_vs_5z” task, while it performs competitively in the rest of the SMAC tasks.

IA2C: The stochastic policy of IA2C appears to be particularly effective on all environments except in a few SMAC tasks. In the majority of tasks, it performs similarly to IPPO with the exception of RWARE and some SMAC tasks. However, it achieves higher returns than IQL in all but two SMAC tasks. Despite its simplicity, IA2C performs competitively compared to all CTDE algorithms, and significantly outperforms COMA and MADDPG in the majority of the tasks.

IPPO: IPPO in general performs competitively in all tasks across the different environments. On average (Figure 2) it achieves higher returns than IA2C in MPE, SMAC and RWARE tasks, but lower returns in the LBF tasks. IPPO also outperforms MAA2C in the partially-observable RWARE tasks, but in general it performs worse compared to its centralised MAPPO version.

5.2 Centralised Training Decentralised Execution

Centralised training aims to learn powerful critics over joint observations and actions to enable reasoning over a larger information space. We find that learning such critics is valuable in tasks which require significant coordination under partial observability, such as the MPE Speaker-Listener and harder SMAC tasks. In contrast, IL is competitive compared to CTDE algorithms in fully-observable tasks of MPE and LBF. Our results also indicate that in most RWARE tasks, MAA2C and MAPPO significantly improve the achieved returns compared to their IL (IA2C and IPPO) versions. However, training state-action value functions appears challenging in RWARE tasks with sparse rewards, leading to very low performance of the remaining CTDE algorithms (COMA, VDN and QMIX).

Centralised Multi-Agent Policy Gradient Centralised policy gradient methods vary significantly in performance.

MADDPG: MADDPG performs worse than all the other algorithms except COMA, in the majority of the tasks. It only performs competitively in some MPE tasks. It also exhibits very low returns in discrete grid-world environments LBF and RWARE. We believe that these results are a direct consequence of the biased categorical reparameterisation using Gumbel-Softmax.

COMA: In general, COMA exhibits one of the lowest performances in most tasks and only performs competitively in one SMAC task. We found that COMA suffers very high variance in the computation of the counterfactual advantage. In the Speaker-Listener task, it fails to find the sub-optimal local minima solution that correspond to returns around to -17. Additionally, it does not exhibit any learning in the RWARE tasks in contrast to other on-policy algorithms.

MAA2C: MAA2C in general performs competitively in the majority of the tasks, except a couple of SMAC tasks. Compared to MAPPO, MAA2C achieves slightly higher returns in the MPE and the LBF tasks, but in most cases significantly lower returns in the SMAC and RWARE tasks.

MAPPO: MAPPO achieves high returns in the vast majority of tasks and only performs slightly worse than other algorithms in some MPE and LBF tasks. Its main advantage is the combination of on-policy optimisation with its surrogate objective which significantly improves the sample efficiency compared to MAA2C. Its benefits can be observed in RWARE tasks where its achieved returns exceed the returns of all other algorithms (but not always significantly).

Value Decomposition Value decomposition is an effective approach in most environments. In the majority of tasks across all environments except RWARE, VDN and QMIX outperform or at least match the highest returns of any other algorithm. This suggests that VDN and QMIX share the major advantages of centralised training. In RWARE, VDN and QMIX do not exhibit any learning, similar to IQL, COMA and MADDPG, indicating that value decomposition methods require sufficiently dense rewards to successfully learn to decompose the value function into the individual agents.

VDN: While VDN and QMIX perform similarly in most environments, the difference in performance is most noticeable in some MPE tasks. It appears VDN’s assumption of linear value function decomposition is mostly violated in this environment. In contrast, VDN and QMIX perform similarly in most SMAC tasks and across all LBF tasks, where the global utility can apparently be represented by a linear function of individual agents’ utilities.

QMIX: Across almost all tasks, QMIX achieves consistently high returns, but does not necessarily achieve the highest returns among all algorithms. Its value function decomposition allows QMIX to achieve slightly higher returns in some of the more complicated tasks where the linear value decomposition of VDN in is not sufficient.

5.3 Parameter Sharing

Figure 3 presents the normalised maximum returns averaged over the nine algorithms and tasks with and without parameter sharing in all environments. We observe that in all environments except the matrix games, parameter sharing improves the returns over no parameter sharing. While the average values presented in Figure 3 do not seem statistically significant, by looking closer in Tables 3 and 7 we observe that in several cases of algorithm-task pairs the improvement due to parameter sharing seems significant. Such improvements can be observed for most algorithms in MPE tasks, especially in Speaker-Listener and Tag. For SMAC, we observe that parameter sharing improves the returns in harder tasks. Similar observations can be made for LBF and RWARE. In these environments, the return improvement of parameter sharing appears to correlate with the sparsity of rewards. For tasks with larger grid-worlds or fewer agents, where the reward is more sparse, parameter sharing leads to large increases in returns compared to simpler tasks. This does not come as a surprise since parameter sharing uses a larger number of trajectories to train the same shared architecture to improve sample efficiency compared to no parameter sharing.

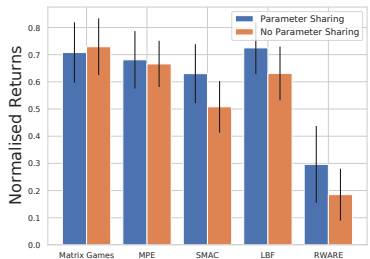


Figure 3: Normalised maximum returns averaged over all algorithms with/without parameter sharing (with standard error).

6 Analysis

Independent learning can be effective in multi-agent systems. Why and when? It is often stated that IL is inferior to centralised training methods due to the environment becoming non-stationary from the perspective of any individual agent. This is true in many cases and particularly crucial when IL is paired with off-policy training from an experience replay buffer, as pointed out by Lowe et al. [2017]. In our experiments, IQL trains agents independently using such a replay buffer and is thereby limited in its performance in tasks that require extensive coordination among the agents. There, agents depend on the information about other agents and their current behaviour to choose well-coordinated actions and hence learning such policies from a replay buffer (where other agents differ in their behaviour) appears infeasible. However, this is not a concern to multi-agent environments in general. In smaller SMAC tasks and most LBF tasks, each agent can independently learn a policy that achieves relatively high returns by utilising only its local observation history, and without requiring extensive coordination with other agents. E.g. in LBF, agents "only" have to learn to pick up the food until it is collected. Of course, they will have to coordinate such behaviour with other agents, but naively going to food (especially when others are also close) and attempting to pick it up can be a viable local minima policy, and hard to improve upon. Whenever more complicated coordination is required, such as simultaneously picking up an item with higher level, exploring and learning those joint actions becomes difficult. IA2C and IPPO on the other hand learn on-policy, so there is no such training from a replay buffer. These algorithms should be expected to achieve higher returns in the majority of the tasks as they learn given the currently demonstrated behaviour of other agents. As long as effective behaviour is eventually identified through time and exploration, IA2C and IPPO can learn more effective policies than IQL despite also learning independently.

Centralised information is required for extensive coordination under partial-observability. We note that the availability of joint information (observations and actions) over all agents serves as a

powerful training signal to optimise individual policies whenever the full state is not available to individual agents. Comparing the performance in Table 3 of IA2C and MAA2C, two almost identical algorithms aside from their critics, we notice that MAA2C achieves equal or higher returns in the majority of the tasks. This difference is particularly significant in tasks where agent observations lack important information about other agents or parts of the environment outside of their receptive field due to partial observability. This can be observed in RWARE tasks with 4 agents which require extensive coordination so that agents are not stuck in the narrow passages. However, in RWARE Tiny 2p task, the performance of IA2C and MAA2C is similar as only two agents rarely get stuck in the narrow passages. Finally, IA2C and MAA2C have access to the same information in fully-observable tasks, such as most LBF and MPE tasks, leading to similar returns. A similar pattern can be observed for IPPO and MAPPO. However, we also observe that centralised training algorithms such as COMA, MADDPG, VDN and QMIX are unable to learn effective behaviour in the partially-observable RWARE. We hypothesise that training larger networks over the joint observation- and action-space, as required for these algorithms, demands sufficient training signals. However, rewards are sparse in RWARE and observations are comparably large.

Value decomposition – VDN vs QMIX. Lastly, we address the differences observed in value decomposition applied by VDN and QMIX. Such decomposition offers an improvement in comparison to the otherwise similar IQL algorithm across most tasks. Both VDN and QMIX are different in their decomposition. QMIX uses a trainable mixing network to compute the joint Q-value compared to VDN which assumes linear decomposition. With the additional flexibility, QMIX aims to improve learnability, i.e. it simplifies the learning objective for each agent to maximise, while ensuring the global objective is maximised by all agents [Agogino and Tumer, 2008]. Such flexibility appears to mostly benefit convergence in harder MPE tasks, such as Speaker-Listener and Tag, but comes at additional expense seen in environments like LBF, where the decomposition did not have to be complex. It appears that the dependency of rewards with respect to complicated interactions between agents in MPE tasks and the resulting non-stationarity benefits more complex decomposition. Finally, VDN and QMIX perform significantly worse than the policy gradient methods (except COMA) in the sparse-reward RWARE tasks. This does not come as a surprise, since the utility of the agents is rarely greater than 0, which makes it hard to successfully learn the individual utilities.

7 Conclusion

We evaluated nine MARL algorithms in a total of 25 cooperative learning tasks, including combinations of partial/full observability, sparse/dense rewards, and a number of agents ranging from two to ten. We compared algorithm performance in terms of maximum and average returns. Additionally, we further analysed the effectiveness of independent learning, centralised training, and value decomposition and identify types of environments in which each strategy is expected to perform well. We created EPyMARL, an open-source codebase for consistent evaluation of MARL algorithms in cooperative tasks. Finally, we implement and open-source LBF and RWARE, two new multi-agent environments which focus on sparse-reward exploration which previous environments do not cover. Our work is limited to cooperative environments and commonly-used MARL algorithms. Competitive environments as well as solutions to a variety of MARL challenges such as exploration, communication, and opponent modelling require additional studies in the future. We hope that our work sheds some light on the relative strengths and limitations of existing MARL algorithms, to provide guidance in terms of practical considerations and future research.

Funding Disclosure

This research was in part financially supported by the UK EPSRC Centre for Doctoral Training in Robotics and Autonomous Systems (G.P., F.C.), and the Edinburgh University Principal’s Career Development Scholarship (L.S.).

References

Adrian K Agogino and Kagan Tumer. Analyzing and visualizing multiagent rewards in dynamic and stochastic domains. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2008.

- Stefano V. Albrecht and Subramanian Ramamoorthy. Comparative evaluation of MAL algorithms in a diverse set of ad hoc team problems. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2012.
- Stefano V. Albrecht and Subramanian Ramamoorthy. A game-theoretic model and best-response learning method for ad hoc coordination in multiagent systems. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2013.
- Stefano V. Albrecht and Peter Stone. Reasoning about hypothetical agent behaviours and their parameters. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2017.
- Marcin Andrychowicz, Anton Raichuk, Piotr Stańczyk, Manu Orsini, Sertan Girgin, Raphael Marinier, Léonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, et al. What matters in on-policy reinforcement learning? a large-scale empirical study. *International Conference on Learning Representations*, 2021.
- Nolan Bard, Jakob N. Foerster, Sarath Chandar, Neil Burch, Marc Lanctot, H. Francis Song, Emilio Parisotto, Vincent Dumoulin, Subhodeep Moitra, Edward Hughes, et al. The hanabi challenge: A new frontier for ai research. *Artificial Intelligence*, page 103216, 2020.
- Charles Beattie, Joel Z. Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew LeFrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016.
- Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, pages 253–279, 2013.
- Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Shared experience actor-critic for multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, 2020.
- Filippos Christianos, Georgios Papoudakis, Arrasy Rahman, and Stefano V. Albrecht. Scaling multi-agent reinforcement learning with selective parameter sharing. In *International Conference on Machine Learning*, 2021.
- Caroline Claus and Craig Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI Conference on Artificial Intelligence*, 1998.
- Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. OpenAI baselines, 2017.
- Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. LIIR: Learning individual intrinsic reward in multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 2019.
- Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. *International Conference on Machine Learning*, 2016.
- Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation matters in deep RL: A case study on PPO and TRPO. *International Conference on Learning Representations*, 2019.
- Jakob N. Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 2016.
- Jakob N. Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. *AAAI Conference on Artificial Intelligence*, 2018.
- Hado Hasselt. Double q-learning. *Advances in Neural Information Processing Systems*, 2010.
- He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep reinforcement learning. *International Conference on Machine Learning*, 2016.

- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. *AAAI Conference on Artificial Intelligence*, 2018.
- Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor. A survey and critique of multiagent deep reinforcement learning. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2019.
- Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. *International Conference on Machine Learning*, 2019.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *International Conference on Learning Representations*, 2017.
- Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro A. Ortega, DJ Strouse, Joel Z. Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. *International Conference on Machine Learning*, 2019.
- Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The malmo platform for artificial intelligence experimentation. In *International Joint Conference on Artificial Intelligence*, 2016.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, page 293–321, 1992.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information Processing Systems*, 2017.
- Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *International Conference on Learning Representations*, 2017.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, pages 529–533, 2015.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *International Conference on Machine Learning*, 2016.
- Igor Mordatch and Pieter Abbeel. Emergence of grounded compositional language in multi-agent populations. *arXiv preprint arXiv:1703.04908*, 2017.
- John Nash. Non-cooperative games. *Annals of mathematics*, pages 286–295, 1951.
- Georgios Papoudakis, Filippos Christianos, Arrasy Rahman, and Stefano V. Albrecht. Dealing with non-stationarity in multi-agent deep reinforcement learning. *arXiv preprint arXiv:1906.04737*, 2019.
- Georgios Papoudakis, Filippos Christianos, and Stefano V. Albrecht. Local information agent modelling in partially-observable environments. In *Advances in Neural Information Processing Systems*, 2021.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019.
- Arrasy Rahman, Niklas Höpner, Filippos Christianos, and Stefano V. Albrecht. Towards open ad hoc teamwork using graph-based policy learning. In *International Conference on Machine Learning*, 2021.

- Roberta Raileanu, Emily Denton, Arthur Szlam, and Rob Fergus. Modeling others using oneself in multi-agent reinforcement learning. *International Conference on Machine Learning*, 2018.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. *International Conference on Machine Learning*, 2018.
- Cinjon Resnick, Wes Eldridge, David Ha, Denny Britz, Jakob Foerster, Julian Togelius, Kyunghyun Cho, and Joan Bruna. Pommerman: A multi-agent playground. *arXiv preprint arXiv:1809.07124*, 2018.
- Heechang Ryu, Hayong Shin, and Jinkyoo Park. Multi-agent actor-critic with hierarchical graph attention network. *AAAI Conference on Artificial Intelligence*, 2020.
- 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. The StarCraft multi-agent challenge. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2019.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. *International Conference on Machine Learning*, 2014.
- Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation. *Advances in Neural Information Processing Systems*, 2016.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *International Conference on Autonomous Agents and Multi-Agent Systems*, 2018.
- Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 1988.
- Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. *International Conference on Machine Learning*, 1993.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. StarCraft II: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- Christopher J. C. H. Watkins and Peter Dayan. Q-learning. *Machine Learning*, 1992.
- Chao Yu, Akash Velu, Eugene Vinitzky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of mappo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.

Responsibility Statement

The authors bear all responsibility in case of violation of rights in the proposed environments and the included benchmark. All the proposed environments use an MIT license for their source code.

A New Multi-Agent Reinforcement Learning Environments

As part of this work, we developed and open-sourced two novel MARL environments focusing on cooperation and sparse-rewards. In this supplementary section, we will provide details about both added environments including our motivation for creating them, accessibility and licensing, installation information, a description of their interface and code snippets to get started as well as further details on their observations, reward functions and dynamics. For high-level descriptions, see Sections 3.4 and 3.5 and see Appendix C for more information on the specific tasks used for the benchmark.

A.1 Motivation

Environments for MARL evaluation are scattered and few environments have been established as standard benchmarking problems. Most notably, the Starcraft Multi-Agent Challenge (SMAC) [Samvelyan et al., 2019] and Multi-Agent Particle Environment (MPE) [Mordatch and Abbeel, 2017] are prominent examples for MARL environments. While more environments exist, they often represent different types of games such as turn-based board games [Bard et al., 2020] or contain image frames as observations [Johnson et al., 2016, Beattie et al., 2016, Resnick et al., 2018]. Environments with these properties often require further solutions not specific to MARL. Additionally, we found that the core challenge of exploration, for which a multitude of environments exist for single-agent RL research (e.g. Atari games such as Montezuma’s Revenge [Bellemare et al., 2013]), is underrepresented in MARL environments. The Level-Based Foraging and Multi-Robot Warehouse environments aim to represent sparse-reward hard exploration problems which require significant cooperation across agents. Both environments are flexible in their configurations to enable partial- or full-observability, fully-cooperative or mixed reward settings and allow faster training than a multitude of other environments (see Appendix A.8 for a comparison of simulation speeds across all environments used in the benchmark).

A.2 Accessibility and Licensing

Both environments are publicly available on GitHub under the following links.

Table 4: GitHub repositories for Level-Based Foraging and Multi-Robot Warehouse environments.

Level-Based Foraging	https://github.com/uoel-agents/lb-foraging
Multi-Robot Warehouse	https://github.com/uoel-agents/robotic-warehouse

The environments are licensed under the MIT license which can be found in the LICENSE file within respective repositories. Both environments will be supported and maintained by the authors as needed.

A.3 Installation

Our novel environments can be installed as Python packages. We recommend users to setup a virtual environment to manage packages and dependencies for individual projects, e.g. using `venv` or `Anaconda`. Then the code repository can be cloned using `git` and installed as a package using the Python package manager `pip` as follows at the example of the Multi-Robot Warehouse environment:

```
$git clone git@github.com:uoel-agents/robotic-warehouse.git
$cd robotic-warehouse
$pip install -e .
```

In order to install the Level-Based Foraging environment, execute the following, similar code:

```
$git clone git@github.com:uoe-agents/lb-foraging.git
$cd lb-foraging
$pip install -e .
```

A.4 Environment Interface

Both environments follow the interface framework of OpenAI's Gym. Below, we will piece-by-piece explain the interface and commands needed to interact with the installed environment within Python.

Package import First, the installed package and gym (installed above as dependency) need to be imported (shown at the example of the Multi-Robot Warehouse):

```
import gym
import robotic_warehouse
```

Environment creation Following the import of the required packages, the environment can be instantiated. A large selection of configurations for both environments are already registered to gym upon importing the package. These can simply be created:

```
env = gym.make("rware-tiny-2ag-v1")
```

For an overview over the naming of these environment configuration names, see Appendix A.5.

Start an episode A new episode within the environment can be started with the following command:

```
obs = env.reset()
```

This function returns the initial state or observation of the environment, indicating s_0 .

Environment steps In order to interact with the environment, an action for each agent should be provided. In the case of the "rware-tiny-2ag-v1", there are two agents in the environment. Therefore, the environment expects to receive an action for all agents. In order to gain insight into the shape of expected actions or received observations, the respective spaces of the environment can be inspected:

```
env.action_space # Tuple(Discrete(5), Discrete(5))
env.observation_space # Tuple(Box(XX,), Box(XX,))
```

Using the action space of the environment, we can sample random valid actions and take a step in the environment:

```
actions = env.action_space.sample() # the action space can be sampled
print(actions) # e.g. (1, 0)
next_obs, reward, done, info = env.step(actions)

print(done) # [False, False]
print(reward) # [0.0, 0.0]
```

Note, that for these multi-agent environments, actions are a **list** or **tuple** containing individual actions for all agents. Similarly, the received values are also lists of observations at the next timestep (`next_obs`), rewards for the completed transition (`reward`), flags indicating whether the episode has terminated (`done`) and an additional dictionary (`info`) which may contain meta-information on the transition or episode.

Rendering Both novel environments support rendering to visualise states and inspect agents' behaviour. See Figure 4 for exemplary visualisations of these environments.

```
env.render()
```

Single episode We can put all these steps together for a script which executes a single episode using random actions and rendering the environment:

```
import gym
import robotic_warehouse

env = gym.make("rware-tiny-2ag-v1")

obs = env.reset()
done = [False] * env.n_agents

while not all(done):
    actions = env.action_space.sample()
    next_obs, reward, done, info = env.step(actions)
    env.render()

env.close()
```

A.5 Environment Naming

In this section, we will briefly outline the configuration possibilities for Level-Based Foraging and Multi-Robot Warehouse environments.

Level-Based Foraging For the Level-Based Foraging environment, we can create an environment as follows

```
env = gym.make("Foraging<obs>-<x_size>x<y_size>-<n_agents>p-<food>f<force_c>-v1")
```

with the following options for each field:

- `<obs>`: This optional field can either be empty (`""`) or indicate a partially observable task with visibility radius of two fields (`"-2s"`).
- `<x_size>`: This field indicates the horizontal size of the environment map and can by default take any values between 5 and 20.
- `<y_size>`: This field indicates the vertical size of the environment map and can by default take any values between 5 and 20. It should be noted, that upon import only environments with square dimensions (`<x_size> = <y_size>`) are registered and ready for creation.
- `<n_agents>`: This field indicates the number of agents within the environment. By default, any values between 2 and 5 are automatically registered.
- `<food>`: This field indicates the number of food items scattered within the environment. It can take any values between 1 and 10 by default.
- `<force_c>`: This optional field can either be empty (`""`) or indicate a task with only "cooperative food" (`"-coop"`). In the latter case, the environment will only contain food of a level such that all agents have to cooperate in order to pick the food up. This mode should only be used with up to four agents.

In order to register environments of more different configurations, see the current registration.

Multi-Robot Warehouse For the Multi-Robot Warehouse environment, we can create an environment as follows

```
env = gym.make("rware-<size>-<num_agents>ag<diff>-v1")
```

with the following options for each field:

- `<size>`: This field represents the size of the warehouse. By default the size can take on four values: `"tiny"`, `"small"`, `"medium"` and `"large"`. These size identifiers define the number of rows and columns of groups of shelves within the warehouse and set these to be (1, 3), (2, 3), (2, 5) and (3, 5) respectively. By default, each group of shelves consists of 16 shelves organised in a 8×2 grid.

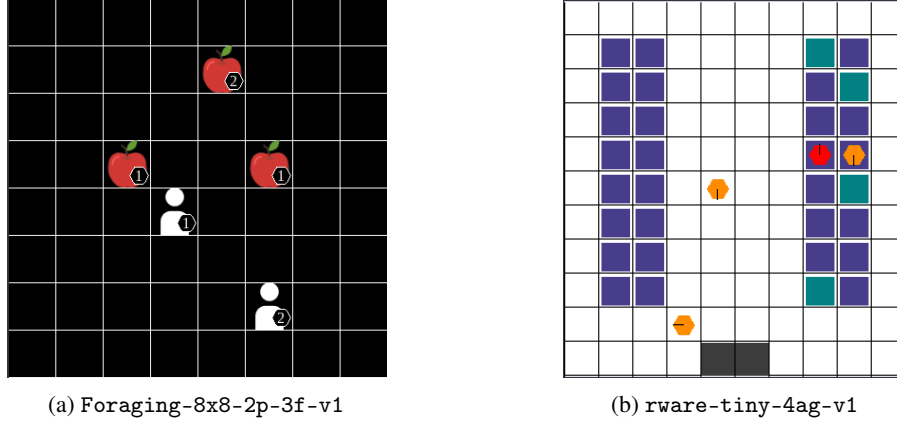


Figure 4: Environment renderings matching observations for (a) Level-Based Foraging and (b) Multi-Robot Warehouse.

- `<num_agents>`: This field indicates the number of agents and can by default take any values between 1 and 20.
- `<diff>`: This optional field can indicate changes in the difficulty of the environment given by the total number of requests at a time. Agents have to collect and deliver specific requested shelves. By default, there are N requests at each point in time with N being the number of agents. With this field, the number of requests can be set to half ("`-hard`") or double ("`-easy`") the number of agents.

Additionally, by default agents observe only fields within immediate grids next to their location and episodes terminate after 500 steps. For a more extensive set of configurations, including variations of visibility, see the full registration.

A.6 Environment Details - Level-Based Foraging

Below, we will provide additional details to observations, actions and dynamics for the Level-Based Foraging environment.

Observations As seen above, agents receive observations at each timestep which correspond to the full state of the environment or a partial view in the case of partial observability. Below, we will outline a realistic observation at the example of the `Foraging-8x8-2p-3f-v1` task visualised in Figure 4a:

```
(
  array([1., 4., 2., 3., 2., 1., 3., 5., 1., 6., 6., 2., 4., 4., 1.],
        dtype=float32),
  array([1., 4., 2., 3., 2., 1., 3., 5., 1., 4., 4., 1., 6., 6., 2.],
        dtype=float32)
)
```

The observation consists of two arrays, each corresponding to the observation of one of the two agents within the environment. Within that vector, triplets of the form $(x, y, level)$ are written. Specifically, the first three (number of food items in the environment) triplets for a total of 9 elements contain the x and y coordinates and level of each food item, and the following two (number of agents) triplets have the respective values for each agent. The coordinates always start from the bottom left square in the observability radius of the agent. When food items or agents are not visible, either because they are outside of the observable radius or the food has been picked up, then the respective values are replaced with $(-1, -1, 0)$.

Actions In Level-Based Foraging environments, each agent has six possible discrete actions to choose at each timestep:

$$A^i = \{\text{Noop, Move North, Move South, Move West, Move East, Pickup}\}$$

While the first action corresponds to the action simply staying within its grid and the last action being used to pickup nearby food, the remaining four actions encode discrete 2D navigation. Upon choosing these actions, the agent moves a single cell in the chosen direction within the grid.

Rewards Agents only receive non-zero reward for picking up food within the Level-Based Foraging environment. The reward for picking up food depends on the level of the collected food as well as the level of each contributing agent. The reward of agent i for picking up food is defined as

$$r^i = \frac{FoodLevel * AgentLevel}{\sum FoodLevels \sum LoadingAgentsLevel}$$

with normalisation. Rewards are normalised to ensure that the sum of all agents' returns on a solved episode equals to one.

Dynamics The transition function of the Level-Based Foraging domain is fairly straightforward. Agents transition to cells following their movement based on chosen actions. Furthermore, agents successfully collect food as long as the sum of the levels of all loading agents is greater or equal to the level of the food.

A.7 Environment Details - Multi-Robot Warehouse

Now, we will provide additional details to observations, actions and dynamics for the Multi-Robot Warehouse environment.

Observations Agent observations for the Multi-Robot Warehouse environment are defined to be partially observable and contain all information about cells in the immediate proximity of an agent. By default, each agent observes information within a 3×3 square centred on the agent, but the visibility range can be modified as an environment parameter. Agents observe their own location, rotation and load, the location and rotation of other observed agents as well as nearby shelves with information whether those are requested or not. For the precise details, see the example below matching the state visualised in Figure 4b:

```
(
  array([8., 3., 0., 0., 1., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0., 0., 1.,
        0., 0., 0., 1., 1., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0.,
        1., 0., 1., 0., 1., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 1.,
        1., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 1., 0., 0., 1., 0., 0.,
        0., 0., 0.], dtype=float32),
  array([4., 9., 0., 0., 0., 1., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 0., 1.,
        0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
        0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0.,
        1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0.,
        0., 0., 0.], dtype=float32),
  ...
)
```

Again, each element of the list is the observation that corresponds to each of the agents. The first three values of the vectors correspond to the agent itself with the x and y coordinates, and a value of "1" or "0" depending on whether the agent is currently carrying a shelf. The next four values are a one-hot encoding of the direction the current agent is facing (up/down/left/right respectively for each item in the one-hot encoding). Then, a single value of "1" or "0" represents whether the agent is currently in a location that acts as a path and therefore not allowed to place shelves on that location. The rest of the values can be split into 9 groups of 7 elements, each group corresponding to a square in the observation radius (3x3 centred around the agent - can be increased from the default for more visibility). In this group, the elements are only ones and zeros and correspond to: agent exists in the square, one-hot encoding (4 elements) of direction of agent (if exists), shelf exists in the square, shelf (if exists) is currently requested to be delivered to the goal location.

Actions The action space within the Multi-Robot Warehouse environment is very similar to the Level-Based Foraging domain with four discrete actions to choose from:

$$A^i = \{\text{Turn Left, Turn Right, Move Forward, Load/ Unload Shelf}\}$$

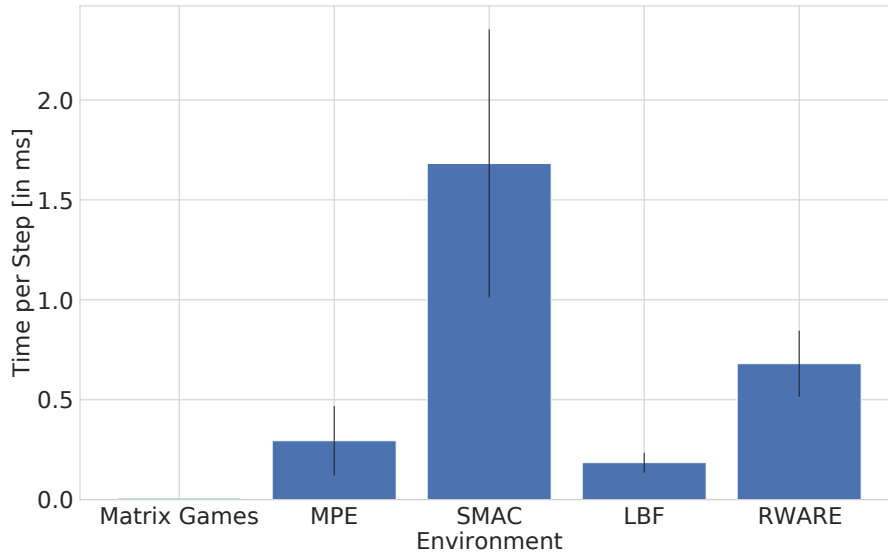


Figure 5: Mean simulation time per step for all environments. Bars indicate standard deviations of simulation speed across all tasks within the environments.

Agents also have to navigate a 2D grid-world, but they are only able to move forward or rotate to change their orientation in 90° steps in either direction. Agents are unable to move upon cells which are already occupied. Besides movement, agents are only able to load and unload shelves. Loading shelves is only possible when the agent is located at the location of an unloaded shelf. Similarly, agents are only able to unload a currently loaded shelf on a location where no shelf is located but is within a group where initially a shelf has been stored.

Rewards At each time, a set number of shelves R is requested. Agents are rewarded with a reward of 1 for successfully delivering a requested shelf to a goal location at the bottom of the warehouse. A significant challenge in this environment is for agents to successfully deliver requested shelves but also finding an empty location to return the previously delivered shelf. Without unloading the previously delivered shelf, the agent is unable to collect new requested shelves. Having multiple steps between deliveries leads to a very sparse reward signal.

Dynamics Agents move through the grid-world as expected given their chosen actions. Whenever multiple agents collide, i.e. they attempt to move to the same location, movement is resolved in a way to maximise mobility. When two or more agents attempt to move to the same location, we prioritise agents to move which also blocks other agents. Otherwise, the selection is done arbitrarily. Additionally, it is worth noting that a shelf is uniformly sampled and added to the list of currently requested shelves whenever a previously requested shelf is delivered to a goal location. Therefore, R requested shelves are constantly ensured. Note that R directly affects the difficulty of the environment. A small R , especially on a larger grid, dramatically affects the sparsity of the reward and thus makes exploration more difficult as randomly delivering a requested shelf becomes increasingly improbable.

A.8 Environments Simulation Speed Comparison

In order to demonstrate and compare simulation speed of all considered environments within the conducted benchmark, we simulate 10,000 environments steps using random action selection. Environments were simulated without rendering to represent conditions at training or evaluation time. Total time and time per step are reported in Table 5. Additionally, mean simulation time per step with standard deviations across for all environments are reported in Figure 5.

Table 5: Simulation time for 10,000 steps and time per step in all 25 tasks.

	Tasks	Number of agents	Total time [in s]	Time per step [in ms]
Matrix Games	Climbing	2	0.079	0.008
	Penalty k=0	2	0.079	0.008
	Penalty k=-25	2	0.078	0.008
	Penalty k=-50	2	0.080	0.008
	Penalty k=-25	2	0.081	0.008
	Penalty k=-25	2	0.078	0.008
MPE	Speaker-Listener	2	0.852	0.085
	Spread	3	4.395	0.439
	Adversary	3	1.651	0.165
	Tag	4	4.911	0.491
SMAC	2s_vs_1sc	2	7.798	0.780
	3s5z	8	18.656	1.866
	MMM2	10	22.480	2.248
	corridor	6	24.890	2.489
	3s_vs_5z	3	10.266	1.027
LBF	15x15-4p-3f	4	2.516	0.252
	8x8-2p-2f-2s-c	2	1.210	0.121
	10x10-3p-3f-2s	3	1.662	0.166
	8x8-2p-2f-c	2	1.253	0.125
	15x15-4p-5f	4	2.559	0.256
	15x15-3p-5f	3	1.939	0.194
	10x10-3p-3f	3	1.849	0.185
RWARE	Tiny 2p	2	4.469	0.447
	Tiny 4p	4	7.976	0.798
	Small 4p	4	7.974	0.797

While matrix games are unsurprisingly the fastest environments to simulate due to their simplicity, Level-Based Foraging tasks are faster to simulate compared to all other environments aside simplest Multi-Agent Particle environment tasks. Despite their arguably more complex environments, the Multi-Robot Warehouse environment is only marginally more expensive to simulate compared to the Multi-Agent Particle environment. It is also worth noting that the cost of simulation of Level-Based Foraging and Multi-Robot Warehouse tasks mostly depends on the number of agents. This allows to simulate large warehouses without additional cost. Unsurprisingly, the Starcraft Multi-Agent Challenge is by far the most expensive environment to simulate as it requires to run the complex game of StarCraft II.

Speed comparisons are conducted on a personal computer running Ubuntu 20.04.2 with a Intel Core i7-10750H CPU (six cores at 2.60GHz) and 16GB of RAM. Simulation was executed on a single CPU thread.

B The EPyMARL Codebase

As part of this work we extended the well-known PyMARL codebase [Samvelyan et al., 2019] to include more algorithms, support more environments as well as allow for more flexible tuning of the implementation details.

B.1 Motivation

Implementation details in Deep RL algorithms significantly affect their achieved returns [Engstrom et al., 2019]. This problem becomes even more apparent in deep MARL research, where the existence of multiple agents significantly increases the amount of implementation details that affect the performance. Therefore, it is often difficult to measure the benefits of newly proposed algorithms compared to existing ones. We believe that a unified codebase, which implements the majority of MARL algorithms used as building blocks in research, would significantly benefit the community. Finally, EPyMARL allows for tuning of additional implementation details compared to PyMARL, including parameter sharing, reward standardisation and entropy regularisation.

B.2 Accessibility and Licensing

All code for EPyMARL is publicly available open-source on GitHub under the following link: <https://github.com/uoε-agents/epymar1>.

All source code that has been taken from the PyMARL repository was licensed (and remains so) under the Apache License v2.0 (included in LICENSE file). Any new code is also licensed under the Apache License v2.0. The NOTICE file in the GitHub repository contains information about the files that have been added or modified compared to the original PyMARL codebase.

B.3 Installation

In the GitHub repository of EPyMARL, the file *requirements.txt* contains all Python packages that have to be install as dependencies in order to use the codebase. We recommend using a Python virtual environment for training MARL algorithm using EPyMARL. After activating the virtual environment, the following commands will install the required packages

```
$ git clone git@github.com:uoε-agents/epymar1.git
$ cd epymar1
$ pip install -r requirements.txt
```

B.4 Execution

EPyMARL is written in Python 3. The neural networks and their operations are implemented using the Pytorch framework [Paszke et al., 2019]. To train an algorithm (QMIX in this example) in a Gym-based environment, execute the following command from the home folder of EPyMARL:

```
$ python3 src/main.py --config=qmix --env-config=gymma with
  ↪ env_args.time_limit=50 env_args.key="lbforaging:Foraging-8x8-2p-3f-v0"
```

where *config* is the configuration of the algorithm, *gymma* is a Gym compatible wrapper, *env_args.time_limit* is the time limit of the task (number of steps before the episode ends), and *env_args.key* is the name of the task. Default configuration for all algorithms can be found in the *src/config/algs* folder of the codebase.

C Task Specifications

Below, we provide details and descriptions of the environments and tasks used for the evaluation.

C.1 Multi-Agent Particle Environment [Mordatch and Abbeel, 2017]

This environment consists of multiple tasks involving the cooperation and competition between agents. All tasks involve particles and landmarks in a continuous two-dimensional environment. Observations consist of high-level feature vectors and agents are receiving dense reward signals. The action space among all tasks and agents is discrete and usually includes five possible actions corresponding to no movement, move right, move left, move up or move down. All experiments in this environment are executed with a maximum episode length of 25, i.e. episodes are terminated after 25 steps and a new episode is started. All considered tasks are visualised in Figure 6.

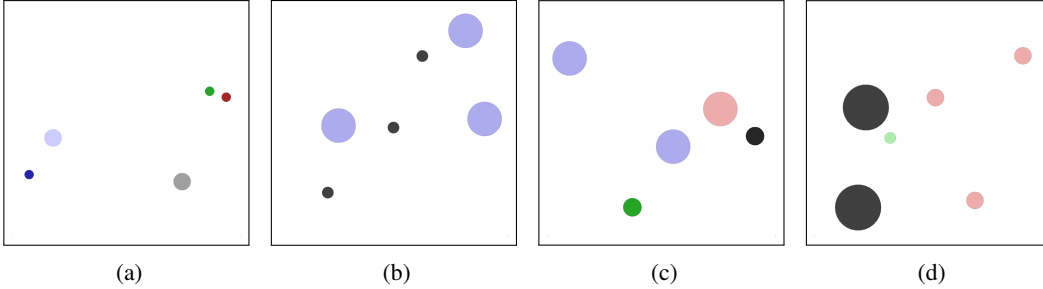


Figure 6: Illustration of MPE tasks (a) Speaker-Listener, (b) Spread, (c) Adversary and (d) Predator-Prey.

MPE Speaker-Listener: In this task, one static speaker agent has to communicate a goal landmark to a listening agent capable of moving. There are a total of three landmarks in the environment and both agents are rewarded with the negative Euclidean distance of the listener agent towards the goal landmark. The speaker agent only observes the colour of the goal landmark. Meanwhile, the listener agent receives its velocity, relative position to each landmark and the communication of the speaker agent as its observation. As actions, the speaker agent has three possible options which have to be trained to encode the goal landmark while the listener agent follows the typical five discrete movement actions of MPE tasks.

MPE Spread: In this task, three agents are trained to move to three landmarks while avoiding collisions with each other. All agents receive their velocity, position, relative position to all other agents and landmarks. The action space of each agent contains five discrete movement actions. Agents are rewarded with the sum of negative minimum distances from each landmark to any agent and an additional term is added to punish collisions among agents.

MPE Adversary: In this task, two cooperating agents compete with a third adversary agent. There are two landmarks out of which one is randomly selected to be the goal landmark. Cooperative agents receive their relative position to the goal as well as relative position to all other agents and landmarks as observations. However, the adversary agent observes all relative positions without receiving information about the goal landmark. All agents have five discrete movement actions. Agents are rewarded with the negative minimum distance to the goal while the cooperative agents are additionally rewarded for the distance of the adversary agent to the goal landmark. Therefore, the cooperative agents have to move to both landmarks to avoid the adversary from identifying which landmark is the goal and reaching it as well. For this competitive scenario, we use a fully cooperative version where the adversary agent is controlled by a pretrained model obtained by training all agents using the MADDPG algorithm for 25,000 episodes.

MPE Predator-Prey: In this task, three cooperating predators hunt a fourth agent controlling a faster prey. Two landmarks are placed in the environment as obstacles. All agents receive their own velocity and position as well as relative positions to all other landmarks and agents as observations. Predator agents also observe the velocity of the prey. All agents choose among five movement actions. The agent controlling the prey is punished for any collisions with predators as well as for leaving the observable environment area (to prevent it from simply running away without needing to learn to evade). Predator agents are collectively rewarded for collisions with the prey. We employ a fully cooperative version of this task with a pretrained prey agent. Just as for the Adversary task, the model for the prey is obtained by training all agents using the MADDPG algorithm for 25,000 episodes.

C.2 StarCraft Multi-Agent Challenge [Samvelyan et al., 2019]

The StarCraft Multi-Agent Challenge is a set of fully cooperative, partially observable multi-agent tasks. This environment implements a variety of micromanagement tasks based on the popular real-time strategy game StarCraft II⁵ and makes use of the StarCraft II Learning Environment (SC2LE) [Vinyals et al., 2017]. Each task is a specific combat scenario in which a team of agents, each agent controlling an individual unit, battles against an army controlled by the centralised built-in AI of

⁵StarCraft II is a trademark of Blizzard Entertainment™.

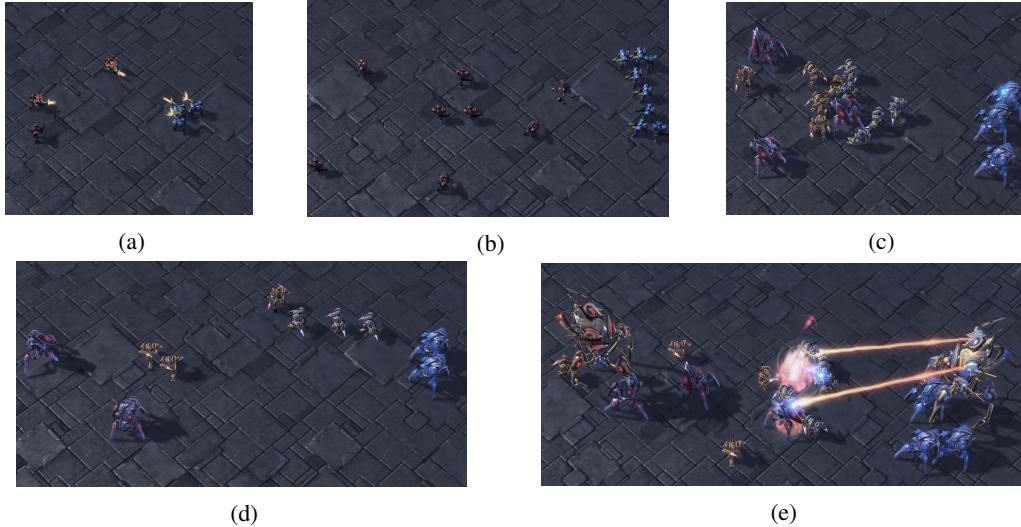


Figure 7: Examples of SMAC tasks with various team configurations and unit types.

the StarCraft game. These tasks require agents to learn precise sequences of actions to enable skills like *kiting* as well as coordinate their actions to focus their attention on specific opposing units. All considered tasks are symmetric in their structure, i.e. both armies consist of the same units. Figure 7 visualises each considered task in this environment.

SMAC 2s_vs_1sc: In this scenario, agents control two stalker units and defeat the enemy team consisting of a single, game-controlled spine crawler.

SMAC 3s5z: In this symmetric scenario, each team controls three stalkers and five zerglings for a total of eight agents.

SMAC MMM2: In this symmetric scenario, each team controls seven marines, two marauders, and one medivac unit. The medivac unit assists other team members by healing them instead of inflicting damage to the enemy team.

SMAC corridor: In this asymmetric scenario, agents control six zealots fighting an enemy team of 24 zerglings controlled by the game. This tasks requires agents to make effective use of terrain features and employ certain game-specific tricks to win.

SMAC 3s_vs_5z: Finally, in this scenario a team of three stalkers is controlled by agents to fight against a team of five game-controlled zerglings.

C.3 Level-Based Foraging [Albrecht and Ramamoorthy, 2013]

The Level-Based Foraging environment consists of tasks focusing on the coordination of involved agents. The task for each agent is to navigate the grid-world map and collect items. Each agent and item is assigned a level and items are randomly scattered in the environment. In order to collect an item, agents have to choose a certain action next to the item. However, such collection is only successful if the sum of involved agents’ levels is equal or greater than the item level. Agents receive reward equal to the level of the collected item. Figure 8 shows the tasks used for our experiments. Unless otherwise specified, every agent can observe the whole map, including the positions and levels of all the entities and can choose to act by moving in one of four directions or attempt to pick up an item.

The tasks that were selected are denoted first by the grid-world size (e.g. 15×15 means a 15 by 15 grid-world). Then, the number of agents is shown (e.g. “4p” means four agents/players), and the number of items scattered in the grid (e.g. “3f” means three food items). We also have some special flags, the “2s” which denotes partial observability with a range of two squares. In those tasks, agents can only observe items or other agents as long as they are located in a square of size 5x5 centred

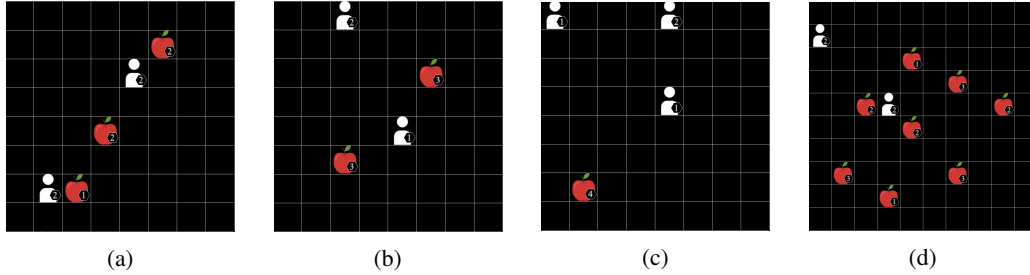


Figure 8: Examples of LBF tasks with variable agents, grid-sizes and food items.

around them. Finally, the flag “c” means a cooperative-only variant where all items can only be picked up if all agents in the level attempt to load it simultaneously.

There are many aspects of LBF that are interesting. An increased number of agents and grid-world size naturally makes the problem harder, requiring agents to coordinate to cover a larger area. In addition, partial observability tasks could benefit from agents modelling other agents since it can improve coordination. A task that we were unable to get any algorithm to learn is the cooperative-only variant with three or more agents. In that case, many agents are required to coordinate to be able to gather any reward, making the task very hard to solve given the sparsity of the rewards. This last task could on its own provide an interesting challenge for research on multi-agent intrinsic exploration.

Notably, LBF is not only designed for a cooperative reward. In other settings, it can work as a mixed competitive/cooperative environment, where agents must switch between cooperating (for gathering items that they can only pick with others) and competing for items that they can load on their own (without sharing the item reward).

C.4 Multi-Robot Warehouse

The multi-robot warehouse environment is a set of collaborative, partially observable multi-agent tasks simulating a warehouse operated by robots. Each agent controls a single robot aiming to collect requested shelves. At all times, N shelves are requested and each timestep a request is delivered to the goal location, a new (currently unrequested) shelf is uniformly sampled and added to the list of requests. Agents observe a 3×3 grid including information about potentially close agents, given by their location and rotation, as well as information on surrounding shelves and a list of requests. The action space is discrete and contains of four actions corresponding to turning left or right, moving forward and loading or unloading a shelf. Agents are only rewarded whenever an agent is delivering a requested shelf to a goal position. Therefore, a very specific and long sequence of actions is required to receive any non-zero rewards, making this environment very sparsely rewarded. We use multi-robot warehouse tasks with warehouses of varying size and number of agents N (which is also equal to the number of requested shelves). Figure 9 illustrates the tiny and small warehouses with 2 agents.

For this environment, we have defined three tasks. The “tiny” map is a grid-world of 11 by 11 squares and the “small” is 11 by 20 squares. The “2p” and “4p” signify two and four robots (agents) respectively. The larger map is considerably harder given the increased sparsity of the rewards.

There are many challenges in the multi-robot warehouse domain as well. The tasks we test in the benchmark paper are limited because of the algorithms’ inability to learn and decompose the reward under such sparse reward settings. However, if this challenge is surpassed, either by using a non-cooperative setting or a new algorithm, then the domain offers additional challenges by requiring algorithms to scale to a large number of agents, communicate intentions to increase efficiency and more. This environment is based on a real-world problem, and scaling to scenarios with large maps and hundreds of agents still requires a considerable research effort.

D Computational Cost

Approximately 138,916 CPU hours were spent for executing the experiments presented in the paper without considering the CPU hours required for the hyperparameter search. Figure 10 presents the

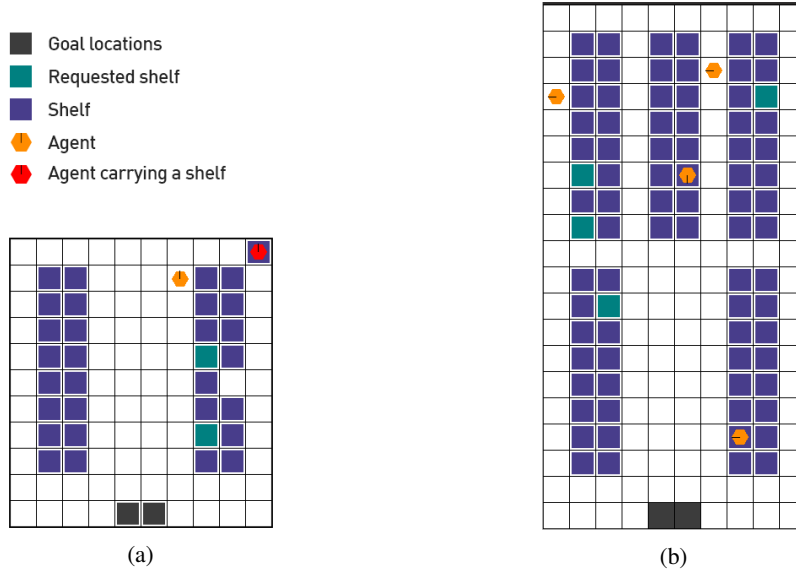


Figure 9: Illustrations of (a) Tiny 2p and (b) Small 4p.

cumulative CPU hours required to train each algorithm in each environments (summed over the different tasks and seeds) with and without parameter sharing using the best identified hyperparameter configurations reported in Appendix I. We observe that the computational cost of running experiments in SMAC is significantly higher compared to any other environment. Finally, the CPU hours required for training the algorithms without sharing is slightly higher compared to training with parameter sharing.

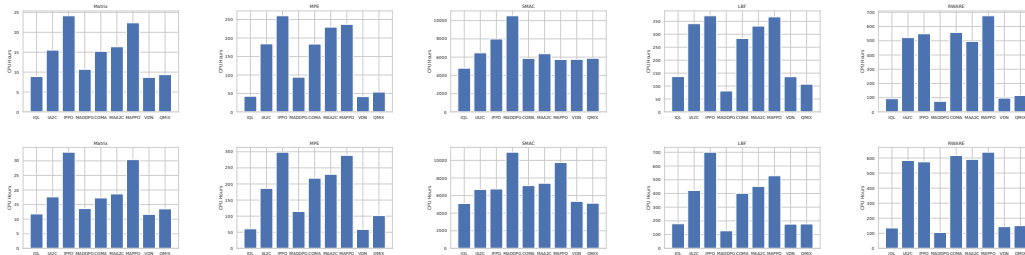


Figure 10: CPU hours required to execute the experiments for each algorithm and environment with (top row) and without (bottom row) parameter sharing.

E Additional Results

Table 6 presents the average returns over training of the nine algorithms in the 25 different tasks with parameter sharing. Tables 7 and 8 present the maximum and the average returns respectively of the nine algorithms in the 25 different tasks without parameter sharing.

Table 6: Average returns and 95% confidence interval over five seeds for all nine algorithms with parameter sharing in all 25 tasks.

Tasks \Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX	
Matrix Games	Climbing	134.65 ± 0.63	169.34 ± 1.09	170.70 ± 1.77	156.45 ± 8.09	177.31 ± 49.52	167.89 ± 4.36	170.76 ± 1.79	125.50 ± 0.54	125.50 ± 0.54
	Penalty k=0	245.64 ± 1.70	244.27 ± 1.13	247.44 ± 0.59	246.39 ± 0.45	245.63 ± 0.64	244.78 ± 0.87	247.69 ± 0.09	239.80 ± 2.93	243.76 ± 2.36
	Penalty k=25	44.65 ± 2.67	48.29 ± 0.30	48.44 ± 0.21	48.59 ± 0.05	48.33 ± 0.20	48.45 ± 0.12	48.46 ± 0.22	43.32 ± 5.98	45.99 ± 3.27
	Penalty k=50	39.70 ± 5.15	46.56 ± 0.57	46.79 ± 0.48	47.22 ± 0.17	46.61 ± 0.44	46.81 ± 0.28	46.82 ± 0.49	39.70 ± 9.63	42.28 ± 6.31
	Penalty k=75	34.75 ± 7.62	44.82 ± 0.84	45.15 ± 0.75	45.83 ± 0.28	44.88 ± 0.68	45.16 ± 0.43	45.19 ± 0.76	34.75 ± 14.26	38.56 ± 9.34
Penalty k=100	29.80 ± 10.10	43.08 ± 1.13	43.50 ± 1.02	44.42 ± 0.36	43.15 ± 0.93	43.51 ± 0.59	43.55 ± 1.04	29.80 ± 18.89	34.85 ± 12.37	
MPE	Speaker-Listener	-27.64 ± 3.90	-17.61 ± 2.99	-17.42 ± 3.23	-18.46 ± 0.68	-38.20 ± 6.29	-15.17 ± 0.44	-15.01 ± 0.64	-27.41 ± 3.11	-21.29 ± 2.79
	Spread	-155.81 ± 1.50	-152.72 ± 0.96	-149.89 ± 2.91	-157.10 ± 2.30	-245.22 ± 84.46	-144.73 ± 4.09	-149.26 ± 0.94	-148.57 ± 1.67	-154.70 ± 4.90
	Adversary	7.58 ± 0.14	10.18 ± 0.05	10.21 ± 0.16	7.80 ± 1.43	6.12 ± 0.35	10.11 ± 0.14	9.61 ± 0.07	7.64 ± 0.21	8.11 ± 0.37
Tag	13.70 ± 1.97	12.43 ± 1.05	13.60 ± 2.95	6.65 ± 3.90	5.11 ± 0.58	11.93 ± 2.09	13.78 ± 4.40	15.24 ± 1.59	15.00 ± 2.73	
SMAC	2s_vs_1sc	14.76 ± 0.45	19.74 ± 0.02	19.44 ± 0.29	10.15 ± 1.32	9.04 ± 0.83	17.89 ± 0.85	19.67 ± 0.09	16.11 ± 0.23	15.98 ± 0.77
	3s5z	14.09 ± 0.28	14.84 ± 1.29	11.80 ± 1.51	8.60 ± 2.35	15.51 ± 0.98	18.82 ± 0.14	19.09 ± 0.38	17.85 ± 0.25	18.36 ± 0.07
	corridor	10.91 ± 0.82	13.14 ± 1.24	14.60 ± 3.43	5.15 ± 0.25	7.00 ± 0.15	7.89 ± 0.28	13.20 ± 2.98	11.14 ± 1.66	11.17 ± 1.88
	MMM2	10.11 ± 0.32	7.31 ± 1.89	9.97 ± 1.33	3.42 ± 0.05	6.50 ± 0.17	9.07 ± 1.35	15.39 ± 0.16	15.93 ± 0.23	15.63 ± 0.32
3s_vs_5z	17.35 ± 0.23	4.32 ± 0.04	13.38 ± 4.36	5.34 ± 0.47	1.15 ± 1.35	6.17 ± 0.39	13.09 ± 2.63	14.72 ± 4.01	9.68 ± 1.87	
LBF	8x8-2p-2f-c	0.75 ± 0.04	0.97 ± 0.00	0.94 ± 0.02	0.32 ± 0.02	0.32 ± 0.12	0.97 ± 0.00	0.95 ± 0.01	0.64 ± 0.09	0.39 ± 0.10
	8x8-2p-2f-2s-c	0.86 ± 0.01	0.97 ± 0.00	0.50 ± 0.01	0.54 ± 0.05	0.24 ± 0.08	0.97 ± 0.00	0.77 ± 0.02	0.83 ± 0.01	0.77 ± 0.03
	10x10-3p-3f	0.54 ± 0.02	0.95 ± 0.01	0.90 ± 0.02	0.20 ± 0.06	0.15 ± 0.05	0.95 ± 0.01	0.91 ± 0.01	0.40 ± 0.05	0.32 ± 0.07
	10x10-3p-3f-2s	0.69 ± 0.02	0.84 ± 0.01	0.62 ± 0.01	0.27 ± 0.02	0.23 ± 0.06	0.85 ± 0.02	0.66 ± 0.01	0.64 ± 0.02	0.67 ± 0.01
	15x15-3p-5f	0.09 ± 0.02	0.61 ± 0.06	0.41 ± 0.09	0.08 ± 0.00	0.06 ± 0.03	0.59 ± 0.09	0.43 ± 0.09	0.08 ± 0.01	0.04 ± 0.01
	15x15-4p-3f	0.24 ± 0.05	0.89 ± 0.03	0.82 ± 0.06	0.13 ± 0.01	0.12 ± 0.03	0.92 ± 0.01	0.79 ± 0.03	0.16 ± 0.03	0.08 ± 0.01
15x15-4p-5f	0.15 ± 0.03	0.59 ± 0.06	0.40 ± 0.13	0.13 ± 0.01	0.07 ± 0.02	0.73 ± 0.02	0.39 ± 0.14	0.15 ± 0.02	0.09 ± 0.02	
RWARE	Tiny 2p	0.04 ± 0.03	2.91 ± 0.45	12.63 ± 1.38	0.11 ± 0.07	0.13 ± 0.05	3.20 ± 0.41	15.42 ± 1.20	0.03 ± 0.01	0.03 ± 0.03
	Tiny 4p	0.33 ± 0.13	10.30 ± 0.93	22.68 ± 7.40	0.28 ± 0.03	0.39 ± 0.06	14.39 ± 4.01	40.17 ± 1.42	0.29 ± 0.13	0.10 ± 0.09
	Small 4p	0.03 ± 0.04	2.45 ± 0.18	9.19 ± 2.36	0.06 ± 0.02	0.08 ± 0.01	3.48 ± 0.42	18.12 ± 1.11	0.02 ± 0.03	0.01 ± 0.01

Table 7: Maximum returns and 95% confidence interval over five seeds for all nine algorithms without parameter sharing in all 25 tasks.

Tasks \Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX	
Matrix Games	Climbing	150.00 ± 0.00	175.00 ± 0.00	175.00 ± 0.00	170.00 ± 10.00	195.00 ± 40.00	175.00 ± 0.00	175.00 ± 0.00	150.00 ± 0.00	155.00 ± 10.00
	Penalty k=0	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	249.94 ± 0.08	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00	250.00 ± 0.00
	Penalty k=25	50.00 ± 0.00	50.00 ± 0.00	90.00 ± 80.00	89.99 ± 80.01	50.00 ± 0.00	130.00 ± 97.98	90.00 ± 80.00	50.00 ± 100.00	50.00 ± 100.00
	Penalty k=50	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.98 ± 0.01	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 100.00	50.00 ± 100.00
	Penalty k=75	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.98 ± 0.01	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 100.00	50.00 ± 100.00
Penalty k=100	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	49.98 ± 0.01	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 0.00	50.00 ± 100.00	50.00 ± 100.00	
MPE	Speaker-Listener	-18.61 ± 5.65	-17.08 ± 3.45	-15.56 ± 4.40 *	-12.73 ± 0.73 *	-26.50 ± 0.50	-13.66 ± 3.67 *	-14.35 ± 3.56 *	-15.47 ± 1.26	-11.59 ± 0.67
	Spread	-141.87 ± 1.68	-131.74 ± 4.33 *	-132.46 ± 3.54 *	-136.73 ± 0.83	-169.04 ± 2.72	-130.88 ± 2.44 *	-128.64 ± 2.83	-142.13 ± 1.86	-139.97 ± 2.51 *
	Adversary	9.09 ± 0.52	10.80 ± 1.97 *	11.17 ± 0.85 *	8.81 ± 0.61	9.18 ± 0.43	10.88 ± 2.43 *	12.04 ± 0.53	9.34 ± 0.57	11.32 ± 0.78 *
	Tag	19.18 ± 2.30	16.04 ± 8.08 *	18.46 ± 5.19 *	2.82 ± 3.56	19.14 ± 7.50 *	26.50 ± 1.42 *	17.96 ± 8.82 *	18.44 ± 2.51	26.88 ± 5.61
SMAC	2s_vs_1sc	15.75 ± 1.08	20.23 ± 0.01	20.15 ± 0.10 *	10.35 ± 2.20	18.48 ± 3.28 *	19.88 ± 0.38 *	20.25 ± 0.01	17.22 ± 0.90	18.23 ± 0.47
	3s5z	16.85 ± 1.43	14.44 ± 2.08	14.77 ± 2.51	15.05 ± 0.81	19.55 ± 0.45 *	19.38 ± 0.30	19.77 ± 0.11	19.08 ± 0.29	18.40 ± 0.70
	corridor	10.86 ± 1.04	8.38 ± 2.90	7.35 ± 0.48	4.92 ± 0.10	4.98 ± 0.42	10.79 ± 0.34	10.02 ± 0.19	16.20 ± 0.44	12.27 ± 2.28
	MMM2	11.06 ± 1.36	13.11 ± 4.27 *	10.29 ± 4.60	6.57 ± 0.53	8.34 ± 0.81	9.26 ± 3.08	8.51 ± 0.76	11.42 ± 2.18 *	15.12 ± 3.42
3s_vs_5z	11.80 ± 1.05 *	4.41 ± 0.02	4.26 ± 0.24	6.21 ± 1.97	4.13 ± 0.05	5.39 ± 0.74	4.62 ± 0.11	14.42 ± 2.22 *	15.13 ± 3.86	
LBF	8x8-2p-2f-c	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	0.43 ± 0.02	0.95 ± 0.07 *	1.00 ± 0.00	1.00 ± 0.00	0.99 ± 0.01 *	0.57 ± 0.15
	8x8-2p-2f-2s-c	0.99 ± 0.01 *	1.00 ± 0.00	0.83 ± 0.03	0.66 ± 0.04	0.92 ± 0.02	1.00 ± 0.00	0.92 ± 0.05	0.99 ± 0.01 *	0.98 ± 0.01
	10x10-3p-3f	0.54 ± 0.12	1.00 ± 0.01	0.96 ± 0.01	0.20 ± 0.03	0.31 ± 0.08	0.99 ± 0.01 *	0.98 ± 0.01	0.39 ± 0.06	0.22 ± 0.03
	10x10-3p-3f-2s	0.77 ± 0.06	0.59 ± 0.02	0.72 ± 0.02	0.29 ± 0.02	0.27 ± 0.09	0.96 ± 0.01	0.70 ± 0.04	0.76 ± 0.04	0.78 ± 0.05
	15x15-3p-5f	0.11 ± 0.02	0.74 ± 0.12	0.62 ± 0.12 *	0.10 ± 0.01	0.09 ± 0.02	0.72 ± 0.05 *	0.49 ± 0.19 *	0.11 ± 0.01	0.06 ± 0.01
	15x15-4p-3f	0.16 ± 0.03	0.99 ± 0.00	0.90 ± 0.01	0.17 ± 0.01	0.21 ± 0.07	0.94 ± 0.05 *	0.89 ± 0.02	0.13 ± 0.02	0.10 ± 0.01
15x15-4p-5f	0.17 ± 0.00	0.77 ± 0.10	0.69 ± 0.05 *	0.15 ± 0.01	0.14 ± 0.03	0.76 ± 0.04 *	0.45 ± 0.17	0.14 ± 0.01	0.08 ± 0.01	
RWARE	Tiny 2p	0.06 ± 0.08	5.56 ± 1.68	8.70 ± 4.04	0.24 ± 0.21	1.46 ± 0.34	6.16 ± 1.94	17.48 ± 4.52	0.06 ± 0.08	0.06 ± 0.08
	Tiny 4p	0.44 ± 0.10	18.02 ± 5.87	14.10 ± 5.20	0.44 ± 0.24	1.48 ± 0.48	31.56 ± 4.29	38.74 ± 2.99	0.34 ± 0.17	0.16 ± 0.15
	Small 4p	0.04 ± 0.05	3.10 ± 0.68	5.78 ± 0.42	0.16 ± 0.12	0.16 ± 0.10	5.00 ± 0.68	13.78 ± 1.63	0.04 ± 0.05	0.06 ± 0.08

Table 8: Average returns and 95% confidence interval over five seeds for all nine algorithms without parameter sharing in all 25 tasks.

Tasks \ Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX
Matrix Games									
Climbing	130.20 ± 4.37	164.72 ± 0.69	171.68 ± 0.41	150.05 ± 3.75	187.50 ± 41.19	169.43 ± 1.01	171.62 ± 0.39	132.52 ± 3.20	132.43 ± 3.37
Penalty k=0	246.73 ± 1.39	243.65 ± 0.93	247.72 ± 0.04	247.88 ± 0.05	247.10 ± 0.38	246.61 ± 0.52	247.73 ± 0.03	247.03 ± 1.85	247.03 ± 0.99
Penalty k=25	49.70 ± 0.24	46.95 ± 0.16	88.05 ± 79.01	86.14 ± 75.19	48.32 ± 0.10	125.90 ± 95.14	88.06 ± 79.03	50.00 ± 0.99	50.00 ± 0.99
Penalty k=50	49.70 ± 0.24	44.31 ± 0.34	46.99 ± 0.23	47.12 ± 0.07	46.54 ± 0.20	46.56 ± 0.33	46.99 ± 0.24	50.00 ± 0.99	50.00 ± 0.99
Penalty k=76	49.70 ± 0.24	41.64 ± 0.50	45.44 ± 0.36	45.74 ± 0.12	44.77 ± 0.31	44.88 ± 0.47	45.44 ± 0.37	50.00 ± 0.99	50.00 ± 0.99
Penalty k=100	49.70 ± 0.24	38.94 ± 0.66	43.89 ± 0.48	44.33 ± 0.16	42.99 ± 0.41	43.21 ± 0.62	43.90 ± 0.52	50.00 ± 0.99	50.00 ± 0.99
MPE									
Speaker-Listener	-30.49 ± 4.67	-23.33 ± 2.67	-22.78 ± 3.10	-17.79 ± 0.97	-33.88 ± 3.38	-19.48 ± 2.92	-20.51 ± 2.85	-29.49 ± 2.36	-20.31 ± 1.84
Spread	-160.10 ± 1.59	-141.31 ± 3.86	-142.86 ± 4.49	-149.53 ± 2.41	-184.72 ± 2.99	-139.54 ± 4.78	-139.20 ± 4.58	-158.60 ± 2.06	-157.04 ± 1.38
Adversary	7.82 ± 0.19	9.18 ± 1.52	9.46 ± 1.02	7.24 ± 2.13	7.28 ± 0.76	9.43 ± 1.93	10.20 ± 0.24	8.06 ± 0.27	8.81 ± 0.52
Tag	12.59 ± 1.60	9.59 ± 4.30	11.90 ± 3.31	1.91 ± 2.09	14.28 ± 5.43	11.79 ± 5.40	10.90 ± 5.47	10.71 ± 0.69	14.27 ± 2.81
SMAC									
2s_vs_1sc	13.37 ± 0.35	19.69 ± 0.05	8.59 ± 1.90	7.91 ± 0.16	15.32 ± 3.21	16.22 ± 3.75	19.68 ± 0.08	15.12 ± 0.46	15.66 ± 0.82
3s5z	14.23 ± 0.84	12.65 ± 1.00	12.26 ± 0.98	7.65 ± 0.54	17.13 ± 1.11	17.94 ± 0.28	19.03 ± 0.19	17.06 ± 0.28	15.46 ± 0.59
MMM2	8.27 ± 0.64	6.98 ± 1.72	6.23 ± 0.88	3.11 ± 0.12	3.82 ± 0.36	9.85 ± 0.19	9.41 ± 0.19	12.72 ± 0.56	8.05 ± 2.05
corridor	8.38 ± 0.82	9.81 ± 1.08	7.80 ± 0.68	4.99 ± 0.17	7.85 ± 0.46	7.78 ± 0.73	8.18 ± 0.43	9.25 ± 0.85	10.76 ± 1.79
3s_vs_5z	9.15 ± 0.46	4.26 ± 0.02	4.20 ± 0.22	3.53 ± 0.39	3.20 ± 0.89	4.91 ± 0.41	4.55 ± 0.02	10.85 ± 1.61	10.09 ± 1.10
LBF									
8x8-2p-2f-c	0.95 ± 0.02	0.96 ± 0.01	0.91 ± 0.02	0.40 ± 0.05	0.63 ± 0.13	0.97 ± 0.00	0.93 ± 0.02	0.93 ± 0.00	0.90 ± 0.00
8x8-2p-2f-2s-c	0.97 ± 0.00	0.96 ± 0.00	0.50 ± 0.00	0.52 ± 0.02	0.63 ± 0.08	0.97 ± 0.00	0.79 ± 0.03	0.96 ± 0.00	0.96 ± 0.00
10x10-3p-3f	0.77 ± 0.05	0.92 ± 0.01	0.85 ± 0.01	0.15 ± 0.00	0.24 ± 0.03	0.92 ± 0.01	0.85 ± 0.02	0.80 ± 0.03	0.81 ± 0.04
10x10-3p-3f-2s	0.78 ± 0.02	0.76 ± 0.01	0.61 ± 0.01	0.22 ± 0.03	0.23 ± 0.05	0.80 ± 0.00	0.62 ± 0.00	0.86 ± 0.02	0.86 ± 0.02
15x15-3p-5f	0.11 ± 0.02	0.39 ± 0.09	0.33 ± 0.13	0.07 ± 0.00	0.06 ± 0.00	0.40 ± 0.03	0.24 ± 0.09	0.18 ± 0.01	0.08 ± 0.03
15x15-4p-3f	0.29 ± 0.04	0.81 ± 0.01	0.60 ± 0.01	0.10 ± 0.01	0.16 ± 0.03	0.74 ± 0.03	0.65 ± 0.06	0.55 ± 0.06	0.12 ± 0.04
15x15-4p-5f	0.17 ± 0.03	0.49 ± 0.10	0.40 ± 0.11	0.11 ± 0.00	0.09 ± 0.01	0.44 ± 0.03	0.25 ± 0.04	0.25 ± 0.01	0.11 ± 0.02
RWARE									
Tiny 2p	0.01 ± 0.01	2.02 ± 0.68	5.36 ± 3.05	0.07 ± 0.05	0.37 ± 0.04	2.23 ± 0.45	10.39 ± 3.04	0.01 ± 0.01	0.01 ± 0.01
Tiny 4p	0.13 ± 0.04	6.10 ± 1.38	9.24 ± 4.20	0.25 ± 0.04	0.50 ± 0.12	10.73 ± 1.31	25.14 ± 1.44	0.10 ± 0.02	0.05 ± 0.02
Small 4p	0.01 ± 0.00	1.02 ± 0.10	3.64 ± 0.38	0.03 ± 0.02	0.07 ± 0.02	1.27 ± 0.09	7.06 ± 0.63	0.01 ± 0.01	0.01 ± 0.01

F Visualisation of the Evaluation Returns During Training

Figure 11 presents the evaluation returns that are achieved during training by the nine algorithms with parameter sharing in the 25 different tasks.

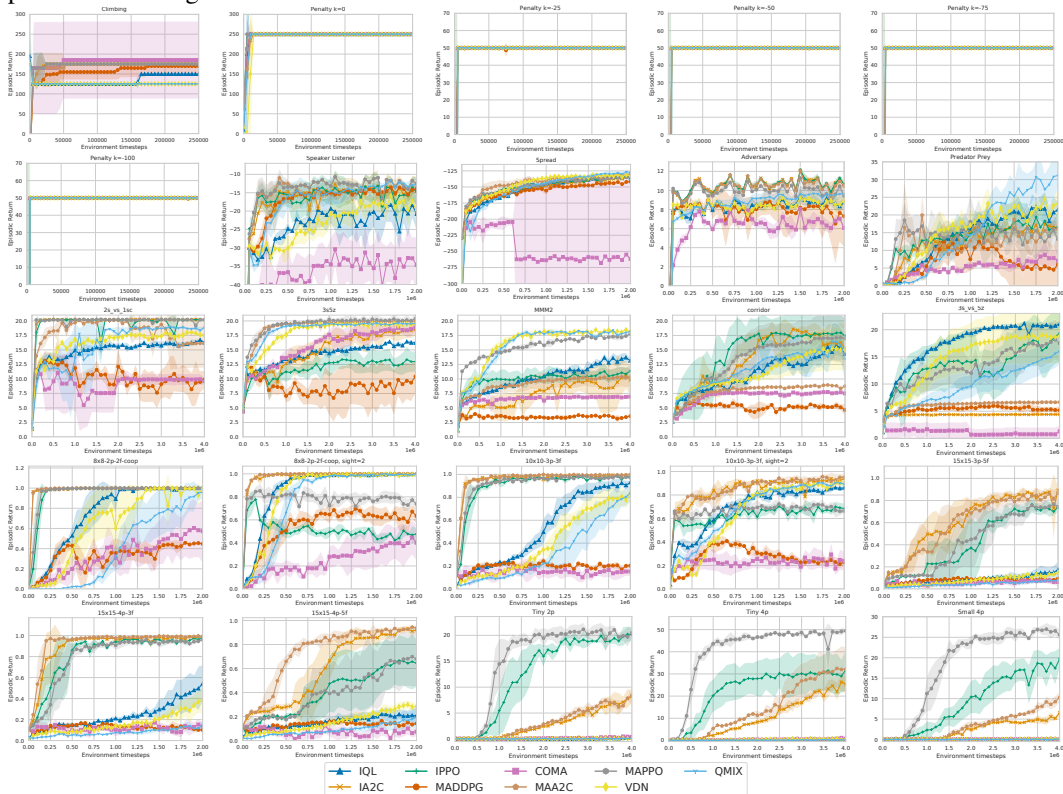


Figure 11: Episodic returns of all algorithms with parameter sharing in all environments showing the mean and the 95% confidence interval over five different seeds.

G SMAC Win-Rates

Table 9: Maximum win-rate and 95% confidence interval over five seeds for all nine algorithms with parameter sharing in all SMAC tasks.

Tasks \Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX
2s_vs_1sc	0.61 ± 0.04	1.00 ± 0.00	1.00 ± 0.00	0.21 ± 0.20	0.34 ± 0.41	1.00 ± 0.00	1.00 ± 0.00	0.74 ± 0.04	0.85 ± 0.03
3s5z	0.39 ± 0.04	0.72 ± 0.23	0.17 ± 0.13	0.15 ± 0.30	0.81 ± 0.19	0.99 ± 0.01	0.96 ± 0.01	0.92 ± 0.05	0.94 ± 0.01
corridor	0.44 ± 0.20	0.80 ± 0.08	0.82 ± 0.25	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.68 ± 0.34	0.44 ± 0.37	0.53 ± 0.29
MMM2	0.27 ± 0.08	0.14 ± 0.17	0.15 ± 0.15	0.00 ± 0.00	0.00 ± 0.00	0.01 ± 0.01	0.73 ± 0.07	0.89 ± 0.04	0.89 ± 0.04
3s_vs_5z	0.67 ± 0.17	0.00 ± 0.00	0.72 ± 0.43	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.41 ± 0.43	0.62 ± 0.31	0.43 ± 0.37

Table 10: Maximum win-rate and 95% confidence interval over five seeds for all nine algorithms without parameter sharing in all SMAC tasks.

Tasks \Algs.	IQL	IA2C	IPPO	MADDPG	COMA	MAA2C	MAPPO	VDN	QMIX
2s_vs_1sc	0.49 ± 0.12	1.00 ± 0.00	0.99 ± 0.01	0.06 ± 0.10	0.79 ± 0.40	0.96 ± 0.04	1.00 ± 0.00	0.64 ± 0.11	0.84 ± 0.02
3s5z	0.48 ± 0.23	0.27 ± 0.25	0.27 ± 0.25	0.13 ± 0.09	0.91 ± 0.08	0.87 ± 0.07	0.94 ± 0.03	0.81 ± 0.06	0.70 ± 0.11
corridor	0.05 ± 0.06	0.31 ± 0.39	0.17 ± 0.30	0.00 ± 0.00	0.01 ± 0.02	0.06 ± 0.12	0.01 ± 0.01	0.08 ± 0.11	0.40 ± 0.32
MMM2	0.04 ± 0.08	0.06 ± 0.13	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.58 ± 0.04	0.23 ± 0.16
3s_vs_5z	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	0.11 ± 0.22	0.23 ± 0.24

H Hyperparameter Optimisation

Table 11: Range of hyperparameters that was evaluated in each environment. N/A means that this hyperparameter was not optimised, and that we used one that was either proposed in the original paper or was found to be the best in the rest of the environments. If only one value is presented it means that this hyperparameter was used for all algorithms in this task.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64/128	64/128	64/128	64/128	64/128
learning rate	0.0001/0.0003/0.0005	0.0001/0.0003/0.0005	0.0005	0.0001/0.0003/0.0005	0.0001/0.0003/0.0005
reward standardisation	True/False	True/False	True/False	True/False	True/False
network type	FC	FC/GRU	FC/GRU	FC/GRU	FC/GRU
evaluation epsilon	0.0/0.05	0.0/0.05	0.0/0.05	0.0/0.05	0.0/0.05
epsilon anneal	50,000/200,000	50,000/200,000	50,000	50,000/200,000	50,000/200,000
target update	200(hard)/0.01(soft)	200(hard)/0.01(soft)	N/A	200(hard)/0.01(soft)	200(hard)/0.01(soft)
entropy coefficient	0.01/0.001	0.01/0.001	N/A	0.01/0.001	0.01/0.001
n-step	5/10	5/10	N/A	5/10	5/10

The parameters of each algorithms are optimised for each environment in one of its tasks and are kept constant for the rest of the tasks within the same environment. Each combination of hyperparameters is evaluated for three different seeds. The combination of hyperparameters that achieved the maximum evaluation, averaged over the three seeds, is used for producing the results presented in this work. Table 11 presents the range of hyperparameters we evaluated in each environment, on the respective applicable algorithms. In general, all algorithms were evaluated in approximately the same number of hyperparameter combination for each environment to ensure consistency. To reduce the computational cost, the hyperparameter search was limited in SMAC compared to the other environments. However, several of the evaluated algorithms have been previously evaluated in SMAC and their best hyperparameters are publicly available in their respective papers.

I Selected Hyperparameters

In this section we present the hyperparameters used in each task. In the off-policy algorithms we use an experience replay to break the correlation between consecutive samples [Lin, 1992, Mnih et al., 2015]. In the on-policy algorithms we use parallel synchronous workers to break the correlation between consecutive samples [Mnih et al., 2015]. The size of the experience replay is either 5K episodes or 1M samples, depending on which is smaller in terms of used memory. Exploration in Q-based algorithms is done with epsilon-greedy, starting with $\epsilon = 1$ and linearly reducing it to 0.05. Additionally, in Q-based algorithms we select action with epsilon-greedy (with a small epsilon value) to ensure that the agents are not stuck. The evaluation epsilon is the hyperparameter that is optimised during the hyperparameter optimisation, with possible values between 0 and 0.05. In the stochastic policy algorithms, we perform exploration by sampling their categorical policy. During execution, in the stochastic policy algorithms, we sample their policy instead of computing the action that maximises the policy. The computation of the temporal difference targets is done using the Double Q-learning [Hasselt, 2010] update rule. In IPPO and MAPPO the number of update epochs per training batch is 4 and the clipping value of the surrogate objective is 0.2.

Tables 12 and 13 present the hyperparameters in all environments for the IQL algorithm with and without parameter sharing respectively.

Table 12: Hyperparameters for IQL with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	128	128	64
learning rate	0.0003	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	False	True	True
network type	FC	FC	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.05	0.05
epsilon anneal	50,000	200,000	50,000	200,000	50,000
target update	200 (hard)	0.01 (soft)	200 (hard)	200 (hard)	0.01 (soft)

Table 13: Hyperparameters for IQL without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	128	64	64	64
learning rate	0.0001	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	True	True	True
network type	FC	FC	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.05	0.05
epsilon anneal	50,000	200,000	50,000	50,000	50,000
target update	0.01 (soft)	0.01 (soft)	200 (hard)	200 (hard)	0.01 (soft)

Tables 14 and 15 present the hyperparameters in all environments for the IA2C algorithm with and without parameter sharing respectively.

Table 14: Hyperparameters for IA2C with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	64	128	128	64
learning rate	0.0005	0.0005	0.0005	0.0005	0.0005
reward standardisation	True	True	True	True	True
network type	FC	GRU	FC	GRU	FC
entropy coefficient	0.01	0.01	0.01	0.001	0.01
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	5	5	5	5	5

Table 15: Hyperparameters for IA2C without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	64	64	64
learning rate	0.0001	0.0005	0.0005	0.0005	0.0005
reward standardisation	True	True	True	True	True
network type	FC	FC	FC	GRU	FC
entropy coefficient	0.01	0.01	0.01	0.01	0.01
target update	200 (hard)	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	5	10	5	5	5

Tables 16 and 17 present the hyperparameters in all environments for the IPPO algorithm with and without parameter sharing respectively.

Table 16: Hyperparameters for IPPO with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	64	128	128	128
learning rate	0.0005	0.0003	0.0005	0.0003	0.0005
reward standardisation	True	True	False	False	False
network type	FC	GRU	GRU	FC	GRU
entropy coefficient	0.001	0.01	0.001	0.001	0.001
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)
n-step	5	5	10	5	10

Table 17: Hyperparameters for IPPO without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	128	64	128	128
learning rate	0.0005	0.0001	0.0005	0.0001	0.0005
reward standardisation	True	True	True	False	False
network type	FC	FC	FC	GRU	FC
entropy coefficient	0.001	0.01	0.001	0.001	0.001
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)
n-step	5	10	10	5	10

Tables 18 and 19 present the hyperparameters in all environments for the MADDPG algorithm with and without parameter sharing respectively.

Table 18: Hyperparameters for MADDPG with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	128	64	64
learning rate	0.0003	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	False	True	False
network type	FC	GRU	GRU	FC	FC
actor regularisation	0.001	0.001	0.01	0.001	0.001
target update	200 (hard)	200 (hard)	0.01 (soft)	200 (hard)	0.01 (soft)

Table 19: Hyperparameters for MADDPG without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	128	64	64
learning rate	0.0005	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	True	True	False
network type	FC	GRU	FC	FC	FC
actor regularisation	0.001	0.01	0.001	0.001	0.001
target update	200 (hard)	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)

Tables 20 and 21 present the hyperparameters in all environments for the COMA algorithm with and without parameter sharing respectively.

Table 20: Hyperparameters for COMA with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	64	128	128	64
learning rate	0.0005	0.0003	0.0005	0.0001	0.0005
reward standardisation	True	True	True	True	True
network type	FC	GRU	FC	GRU	FC
entropy coefficient	0.01	0.001	0.01	0.001	0.01
target update	0.01 (soft)	200 (hard)	0.01 (soft)	200 (hard)	0.01 (soft)
n-step	5	10	5	10	5

Table 21: Hyperparameters for COMA without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	128	128	64
learning rate	0.0003	0.0005	0.0005	0.0001	0.0005
reward standardisation	True	True	True	True	False
network type	FC	GRU	GRU	GRU	FC
entropy coefficient	0.01	0.01	0.01	0.001	0.01
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	10	10	5	5	5

Tables 22 and 23 present the hyperparameters in all environments for the MAA2C algorithm with and without parameter sharing respectively.

Table 22: Hyperparameters for MAA2C with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	128	128	64
learning rate	0.003	0.0005	0.0005	0.0005	0.0005
reward standardisation	True	True	True	True	True
network type	FC	GRU	FC	GRU	FC
entropy coefficient	0.001	0.01	0.01	0.01	0.01
target update	0.01	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	10	5	5	10	5

Table 23: Hyperparameters for MAA2C without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	128	128	128	64
learning rate	0.0005	0.0003	0.0005	0.0005	0.0005
reward standardisation	True	True	True	True	True
network type	FC	GRU	FC	GRU	FC
entropy coefficient	0.001	0.01	0.01	0.01	0.01
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	10	5	5	5	5

Tables 24 and 25 present the hyperparameters in all environments for the MAPPO algorithm with and without parameter sharing respectively.

Table 24: Hyperparameters for MAPPO with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	64	64	128	128
learning rate	0.0005	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	False	False	False
network type	FC	FC	GRU	FC	FC
entropy coefficient	0.001	0.01	0.001	0.001	0.001
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)	0.01 (soft)
n-step	5	5	10	5	10

Table 25: Hyperparameters for MAPPO without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	128	64	128	128
learning rate	0.0005	0.0001	0.0005	0.0001	0.0005
reward standardisation	True	True	True	False	False
network type	FC	FC	GRU	FC	FC
entropy coefficient	0.001	0.01	0.001	0.001	0.001
target update	0.01 (soft)	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)
n-step	5	5	10	10	10

Tables 26 and 27 present the hyperparameters in all environments for the VDN algorithm with and without parameter sharing respectively.

Table 26: Hyperparameters for VDN with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	128	128	128	64
learning rate	0.0001	0.0005	0.0005	0.0003	0.0005
reward standardisation	True	True	True	True	True
network type	FC	FC	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.0	0.05
epsilon anneal	200,000	50,000	50,000	200,000	50,000
target update	0.01 (soft)	200 (hard)	200 (hard)	0.01 (soft)	0.01 (soft)

Table 27: Hyperparameters for VDN without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	64	64	64
learning rate	0.0005	0.0005	0.0005	0.0001	0.0005
reward standardisation	True	True	True	True	True
network type	FC	FC	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.05	0.05
epsilon anneal	50,000	50,000	50,000	50,000	50,000
target update	0.01 (soft)	200 (hard)	200 (hard)	200 (hard)	0.01 (soft)

Tables 28 and 29 present the hyperparameters in all environments for the QMIX algorithm with and without parameter sharing respectively.

Table 28: Hyperparameters for QMIX with parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	64	64	128	64	64
learning rate	0.0003	0.0005	0.005	0.0003	0.0005
reward standardisation	True	True	True	True	True
network type	FC	GRU	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.05	0.05
epsilon anneal	200,000	200,000	50,000	200,000	50,000
target update	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)	0.01 (soft)

Table 29: Hyperparameters for QMIX without parameter sharing.

	Matrix Games	MPE	SMAC	LBF	RWARE
hidden dimension	128	128	64	64	64
learning rate	0.0005	0.0003	0.0005	0.0001	0.0003
reward standardisation	True	True	True	True	True
network type	FC	GRU	GRU	GRU	FC
evaluation epsilon	0.0	0.0	0.05	0.05	0.05
epsilon anneal	50,000	200,000	50,000	50,000	50,000
target update	0.01 (soft)	0.01 (soft)	200 (hard)	0.01 (soft)	0.01 (soft)