

---

# Fully-Decentralized MADDPG with Networked Agents

---

**Diego Bolliger**  
diegobo@student.ethz.ch

**Lorenz Zauter**  
lzauter@student.ethz.ch

**Robert Ziegler**  
roziegler@student.ethz.ch

## Abstract

In this paper, we devise three actor-critic algorithms with decentralized training for multi-agent reinforcement learning in cooperative, adversarial, and mixed settings with continuous action spaces. To this goal, we adapt the MADDPG algorithm by applying a networked communication approach between agents. We introduce surrogate policies in order to decentralize the training while allowing for local communication during training. The decentralized algorithms achieve comparable results to the original MADDPG in empirical tests, while reducing computational cost. This is more pronounced with larger numbers of agents.

## 1 Introduction

Multi-agent reinforcement learning (MARL) is of significant practical value and has become a focus of research in the past years. MARL extends the ideas of Reinforcement Learning (RL) to settings where multiple agents interact in a common environment. These interactions can be competitive, where all agents work towards their own goal, cooperative, where all agents work towards the same goal, or a mixture of the two, e.g. agents working in teams against each other. This includes many practical examples, such as multi robot control, coordinated flight or control in power distribution systems.

We are interested in partially observable stochastic games (POSG), where each agent decides independently about its actions, has its own reward function and the goal is to maximize the individual reward as well as the overall reward. The reward functions are only accessible by the agents themselves and in the execution phase, there is only local information available. Recall that a classical single agent RL system such as a Markov decision process assumes that the agent controls all actions and has access to the entire state. Hence, the POSG setting poses two major challenges: how do agents account for missing information and how can they learn to cooperate despite being partially unaware of what the other agents are doing or aiming to do?

One actor-critic approach presented in [2] makes use of a centralized controller which has access to all information during training, whereas the actions are deterministic and independent during execution as described above. This is achieved by adapting the deep deterministic policy gradient (DDPG) algorithm to the multi-agent setting (MADDPG). A different approach is taken in [6] where both learning and execution is decentralized and information is shared through a communication network. However, the authors assume that all agents observe the full state.

Centralized training can be unfeasible in settings of large number of agents or when communication is limited. In this paper, we investigate possibilities of combining the algorithms from [2] and [6] with the goal of developing *fully decentralized* MARL algorithms. That is, both training and execution are decentralized. An additional goal of this paper is to enable the algorithms to work only with local observations for each agent in training and execution, thus lifting further restrictions imposed by the previous authors.

To this end, we first devise a fully decentralized version of the MADDPG algorithm and in a second step introduce a communication network through which information can be exchanged during training. We develop these algorithms by first considering a cooperative setting, where all agents work towards a common goal. Then afterwards, we adapt them to mixed and adversarial settings. For the evaluation of all our algorithms, we conduct experiments using the multi-particle environment (MPE). [4]

In concrete terms, our contributions are as follows.

- We develop and evaluate a fully decentralized version of MADDPG.
- We extend the fully decentralized MADDP to share critic parameters through a communication network during training. For this, we develop two different algorithms and evaluate both.
- Lastly, we modify all algorithms for the adversarial setting.

## 2 Related Work

In their overview paper [7] from 2021 Zhang et. al. describe a variety of settings. Our work focuses on decentralized algorithms in the setting of stochastic games, what the authors call ‘networked algorithms’, because communication during training is decentralized through a (possibly) time varying network.

The MADDPG algorithm from Lowe et. al. [2] is applicable for both cooperative and competitive settings. By deploying centralized  $Q$ -learning during the training phase, the algorithm is able to handle non-differentiable environment dynamics while yielding decentralized policies during execution time. This is achieved by approximating the shared  $Q$  function by a neural network that has access to the policies of all agent and coordinates accordingly. While this algorithm leads to good performance, it uses a centralized controller with its aforementioned drawbacks. Acknowledging the centralized nature of their algorithm, the authors propose to mitigate the required access to the policies of all  $N$  agents by complementing each agent with  $N - 1$  approximator networks which learn the policies of the other agents. This however, although being more decentralized, increases the scalability problem even further as the number of parameters scales quadratically in  $N$ .

In order to obtain scalability to large amounts of agents, Cui et. al. [1] maintain a distribution to model the mean field of the model. With mean field control (MFC), the agents are able to correlate their actions during training, being effectively guided by the mean field. This requires rigid assumptions on the environment. Through a radial basis function kernel, the training is then further centralized, leading effectively again to an approach that relies on centralized training to achieve decentralized execution.

Zhang et. al [6] take a different approach to obtain an algorithm that trains in a decentralized manner. The authors model the agents on a time varying communication network where the agents perform the actor step individually. In contrast to the other approaches presented so far, their algorithm maintains the decentralization aspect also during the training process, allowing for communication only through connected agents. More precisely, agents average their critic parameters with connected agents (i.e. locally) in a consensus step. While this approach performs well and is one of the few fully decentralized approaches to MARL, we found in our experiments that the consensus update – effectively overriding the critics of the agents with the average of their neighbors’ critics – can lead to instability. Furthermore, Zhang et. al. require the agents to observe the full state space, thus not allowing for partially observable settings.

### 3 Main Results

#### 3.1 Problem Setting

For our work we are interested in the general setting of partially observable stochastic games.

**Definition 1 (POSG)** A partially observable stochastic game (POSG) is defined as a tuple,

$$(\mathcal{D}, \mathcal{S}, \mathcal{A}, P, \mathcal{R}, \mathcal{O}, p_{\mathcal{O}}, \delta_0),$$

where

- a)  $\mathcal{D} = \{1, \dots, N\} = [N]$  is the set of  $N \in \mathbb{N}$  agents.
- b)  $\mathcal{S}$  is a finite set of states,
- c)  $\mathcal{A} = \prod_{n \in \mathcal{D}} A^n$  is the finite set of joint actions and  $A^i$  the set of actions of agent  $i$ ,
- d)  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability function,
- e)  $\mathcal{R} = \{R^i\}_{i \in \mathcal{D}}$  is the set of individual reward functions  $R^i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  of the agents  $i \in [N]$ ,
- f)  $\mathcal{O} = \prod_{i \in \mathcal{D}} O^i$  is the finite set of joint observations, where  $O^i$  is the set of observations available to agent  $i$ ,
- g)  $p_{\mathcal{O}}$  is the observation probability function, and
- h)  $\delta_0$  is the initial state distribution at stage  $t = 0$ .

During execution, as for the Dec-POMDP setting (see [3]) the agents of our POSG are assumed to only act based on their individual observations and to not communicate any further. We assume throughout that action spaces are *continuous*, i.e. open subsets of some Euclidean space.

To facilitate cooperation between the agents in the case of our networked algorithms (algorithm 3 and 4 in A.1) we allow communication between agents along a communication network  $G_t$ . Similar to [6], we define the networked setting below:

**Definition 2 (POSG with networked agents)** A partially observable stochastic game with networked agents is a tuple

$$(\mathcal{D}, \mathcal{S}, \mathcal{A}, P, \mathcal{R}, \mathcal{O}, p_{\mathcal{O}}, \delta_0, \{G_t\}_{t \geq 0}).$$

The first part of the tuple is a POSG, amended by a time varying undirected graph  $G_t := (\mathcal{D}, E_t)$  with the vertex set  $\mathcal{D}$  and the set of edges at time  $t$ ,  $E_t \subset \{(i, j) \in \mathcal{D} \times \mathcal{D} : i \neq j\}$ . Communication between two agents  $i, j \in \mathcal{D}$  at time  $t$  is only possible if  $(i, j) \in E_t$ .

We note that in principle, agents could communicate through the network during action execution *and* training, but we restrict ourselves to communication during training. This is without loss of generality, since the former can be included as actions and observations if needed.

#### 3.2 Fully-Decentralized MADDPG

To decentralize training of the MADDPG algorithm, see section A.1, we propose that each agent  $i$  keeps its own local replay buffer

$$\mathcal{D}^i = (o_i^j, a^j, r_i^j, o_i^{\prime j})_{j=1}^B,$$

of size  $B$ , containing its own observations  $o_i$ , the true joint actions taken  $a = (a_1, \dots, a_N)$ , its own rewards  $r_i$  and next observations  $o_i'$ . The local critic  $Q_i$  is then updated using solely the local replay buffer of agent  $i$ .

A key problem in decentralizing training lies in how to obtain the critic targets  $y_i$  with only local observations. MADDPG needs to access the policies and observations of all agents, which is undesirable in our decentralized setting. The proposed policy approximation in [2] may not need

access to policies of other agents, but still requires access to all local observations while also increasing computational cost.

We propose that instead of approximating the actual policies of the other agents in the optimization target

$$\hat{y}^j = r^j + \gamma Q_i^{\mu^i}(\mathbf{x}'^j, \hat{a}'_1, \dots, a_i, \dots, \hat{a}'_N) \Big|_{\hat{a}'_k = \hat{\mu}_i^k(o_i^j)},$$

which requires knowledge of all local observation, we let each agent learn *surrogate policies* which mimic the behavior of the other agents by only using local observations of the agent in question. Concretely, we let each agent learn a joint policy

$$\boldsymbol{\mu}_{\theta_i}(o_i) = [\boldsymbol{\mu}_{\theta_i}(o_i)]_{k=1}^N,$$

using the sampled policy gradient with respect to the surrogate policies

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu^i}(o_i^j, a^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k} \right). \quad (3.1)$$

Note that the surrogate policies model the other agents as one centrally controlled actor. When interacting with the environment, each agent then uses the joint surrogate policy to select its action. Practically speaking, we increase the output dimension of the actor network to the dimension  $d_A$  of the joint action space and use this output as a plug-in estimator for the joint action to create approximations for the target  $y_i^j$  for the critic update, i.e.

$$\hat{y}_i^j = r_i^j + \gamma Q_i^{\mu^i}(o_i^j, a^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k}.$$

We think of the surrogate policies as a form of imagined, intrinsic representation of the other agents which help to learn a good policy as imaginary friends. This reflects intuitively human behaviour in learning a team task, where it is beneficial to think about the actions of the team members to anticipate their actions.

Note that the critic still learns with transitions  $(o_i, a, r_i, o'_i)$  from the local replay buffer, where the true joint actions are used, meaning that we assume that each agent's critic has access to the true joint actions. We remark that our algorithm, as is, can be used in the cooperative setting only, as each agent  $i$  imagines the other agents' policies to maximize their own (local) reward  $r_i$ .

One drawback of this approach is that all agents train independently which could lead to optimizing for non-compatible policies. We conclude that there is a trade-off between decentralization of training and cooperation.

### 3.3 Fully-Decentralized MADDPG with Networked Agents

The algorithm we propose in section 3.2, in theory, may fail to lead to true cooperative behavior since each agent might optimize for a different joint policy which might contradict each other. This is problematic if the agents are supposed to cooperate. In the following, we propose a compromise between cooperation and decentralization in the form of networked agents, where parameters of the critic are passed to neighbors in a possibly time varying communication network (definition 2).

We represent the communication network at time  $t$  by a right stochastic matrix  $C_t \in \mathbb{R}^{N \times N}$ , where  $C_t(i, j) > 0$  if agent  $i$  receives information from agent  $j$  and  $C_t(i, j) = 0$  otherwise. The weights of the communication matrix  $C_t$  can be adjusted depending on the setting. If communication is not equally possible between all agents or depends on other factors (i.e. proximity), this formulation is able to model these situations as well. Note that the network is exclusively used to pass parameters of the critic and the policy execution does not depend on the network.

#### 3.3.1 Networked Agents with hard consensus update

One option to structure this communication, is to average the critics of neighboring agents after each learning step. Thus effectively overriding the local critics by the (weighted) average of itself and its neighbors. We call this *hard consensus update*, which is performed after each learning interval. The idea is inspired by [6].

The consensus step is then performed, for each agent  $i$ , as

$$\mu_i \leftarrow \sum_{j=1}^N \mathbf{C}_t(i, j) \mu_j,$$

where  $\mu_k$  are the network parameters of the critic of agent  $k$ . This means, in the case where  $\mathbf{C}_t \equiv I$ , we recover fully decentralized learning as in section 3.2. On the other hand, in the case where  $\mathbf{C}_t(i, j) = \frac{1}{N}, \forall i, j$ , the agents effectively learn fully centralized since the same critic is used by all agents. This requires the critics to be parameterized in the same way for all agents.

### 3.3.2 Networked Agents with soft consensus update

As the hard consensus update interferes quite bluntly with the learning of the individual critics, we propose a more subtle approach that leaves the integration of the communication requirement to the agents themselves. We call this *soft consensus update*.

Instead of updating the critics after every learning interval, we change the critic loss to include a penalty term based on the squared relative error of critic parameters of connected agents:

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^\mu(\sigma_i^j, a_i^j) \right)^2 + \underbrace{\zeta \cdot \sum_{j=1}^N \mathbf{C}_t(i, j) \frac{\|\mu_i - \mu_j\|^2}{\|\mu_j\|^2}}_{\text{consensus penalty}},$$

We include a hyperparameter  $\zeta$  to balance individual and collective learning. For the norm on the parameter space, we propose the standard Euclidean/Frobenius norm but other choices are possible too. For numerical stability, a small constant should be added to the the denominator.

### 3.4 Extension to Adversarial and Mixed Settings

As described above, the fully-decentralized algorithms can not be applied as is to adversarial and mixed settings, in contrast to the original MADDPG. This is because, when training the actor determining the surrogate actions, we maximize for a joint action by the sampled policy gradient (3.1) that maximizes the critic value which corresponds only to the own reward of the agent. Since we do not assume access the reward of other agents, another approach must be taken. For this, we let each agent take the assumption that adversarial agents have as goal to diminish their own reward. The actor is then updated in two steps. For  $N$  agents forming  $K$  teams, where team  $k$  has size  $l$ , consisting of agents  $k_1, \dots, k_l$  and having adversaries  $\bar{k}_1, \dots, \bar{k}_{N-l}$ , the policy gradient is first sampled as

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(\sigma_i^j, a^j) \Bigg|_{\substack{a_k^j = [\mu_{\theta_i}(\sigma_i^j)]_k \\ a_{\bar{k}}^j = [a_{\text{obs}}^j]_{\bar{k}}} \right), \quad (3.2)$$

meaning that we plug in the observed adversarial actions  $[a_{\text{obs}}^j]_{\bar{k}}$  (which we assumed to be able to observe) and only optimize according to the surrogate joint team action by taking a gradient *ascent* step.

Then, directly following, we optimize for the surrogate adversarial actions, where the policy gradient is sampled according to

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(\sigma_i^j, a^j) \Bigg|_{\substack{a_k^j = [\mu_{\theta_i}(\sigma_i^j)]_{\bar{k}} \\ a_{\bar{k}}^j = [a_{\text{obs}}^j]_k} \right), \quad (3.3)$$

and a gradient *descent* step is taken.

## 4 Empirical Results

### 4.1 Network architecture

In order to keep our results comparable to [2], we implemented the algorithm with simple MLPs with five hidden layers with 256 hidden nodes. Contrary to earlier suspicions, bigger networks or networks with a more specialized architecture such as GRUs and LSTMs did not improve results.

### 4.2 The communication matrix

One of the most versatile parts of our networked algorithms is their adaptability to various communication situations. As comparability is essential, we decided to construct the communication matrix  $C_t$  for  $N$  agents as follows. For cooperative settings, we choose

$$C_t = \begin{pmatrix} 1 - \eta & \frac{\eta}{N-1} & \cdots & \frac{\eta}{N-1} \\ \frac{\eta}{N-1} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{\eta}{N-1} \\ \frac{\eta}{N-1} & \cdots & \frac{\eta}{N-1} & 1 - \eta \end{pmatrix}$$

where  $\eta$  is a communication hyperparameter.

For mixed settings of type 1 vs.  $N$ , we use the communication matrix of the form

$$C_t = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 - \eta & \frac{\eta}{N-1} & \cdots & \frac{\eta}{N-1} \\ 0 & \frac{\eta}{N-1} & \ddots & \cdots & \frac{\eta}{N-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \frac{\eta}{N-1} & \cdots & \cdots & 1 - \eta \end{pmatrix}.$$

### 4.3 Cooperative setting

#### 4.3.1 Simple Spread Environment

For training and testing of our algorithms in a cooperative setting, we use the ‘simple spread’ environment of the pettingzoo library [4]. The environment consists of  $N$  agents and  $N$  landmarks on a 2D plane. The agents are controlled by the algorithms we are testing and the landmarks are stationary. The goal of the agents is to cover the landmarks while avoiding collisions. They receive a global reward based on the minimum distance of the closest agent to each landmark and are penalized for each collision. Each agent receives the same reward

$$R = \sum_{i=1}^N \min_{j=1, \dots, N} \{\|\ell_i - p_j\|\} - \sum_{i=1}^N \sum_{j=i+1}^N \mathbb{1}_{\{\|p_i - p_j\| < \epsilon\}},$$

where  $\ell_i$  is the position of landmark  $i$ , and  $p_j$  is the position of agent  $j$ , and  $\epsilon$  is a collision threshold. The observations are their relative positions to the landmarks and to the other agents, as well as their own position and velocity. The actions for each agent are to accelerate continuously in some direction and thus changing their trajectory.

#### 4.3.2 Simple Spread Experiment

In the first set of experiments, we are interested in the obtained scores of our new algorithms in comparison with the standard MADDPG algorithm. Figure 1 shows how the evaluation score changes for all four algorithms when training on the experimental setup described in section 4.3.1. We conduct five experiments, for  $N = 2, 3, 4, 5, 10$  agents, training them over  $s = 10^5, 3 \cdot 10^5, 4 \cdot 10^4, 5 \cdot 10^5$  and  $3 \cdot 10^5$  steps respectively. The last experiment was limited due to the computational cost.

As we can see in Figure 1 (also see Figure 4 in the appendix), the standard MADDPG outperforms our algorithms in terms of final score for  $N = 2, 3$  and about the same score is reached for  $N = 4, 5$ . However for  $N = 4, 5$ , MADDPG takes longer to converge.

For  $N = 10$  agents, MADDPG and the decentralized algorithm with hard consensus update are unstable, whereas our other algorithms remain stable. However, MADDPG seems to improve towards the end and a longer training run would have been desirable (it was not possible to repeat the experiment for us due to the computational cost involved).

Since MADDPG has the most information during training and does not have to rely on guesses about the other agents behaviors, it is to be expected that the performance of MADDPG is best in the long run. However, as we see with the training runs with more than three agents, the decentralized algorithms converge much faster. We believe that this is due to the decentralization which reduces complexity of the critic by considering less information. This effect is most pronounced in the case  $N = 10$ , where scaling becomes an issue.

The experiments with  $N = 10$  agents were also performed with a sparse matrix which gave similar results, see B.3.

The fully decentralized algorithm (Algorithm 2 or Section 3.2) performs promisingly. While converging significantly slower for fewer agents, its score tends towards the score of MADDPG. However, the networked algorithms are not able to outperform the fully decentralized algorithm. These train equally fast as the fully decentralized algorithm and especially for  $N = 5$  agents considerably faster than the (centralized) MADDPG algorithm.

While the algorithms with hard and the soft consensus update seem to perform quite similar, the soft consensus update is more stable in case of more aggressive updating (compare B.2).

However, we clearly notice diminishing returns when the number of agents increase compared to the starting scores of effectively random agents. We conclude that the MADDPG algorithm performs poorly in the simple spread already for a low number of agents. This is also confirmed by qualitative examination of the evaluation episodes, which indicate that effective policies are only learned for 2 and 3 agents.

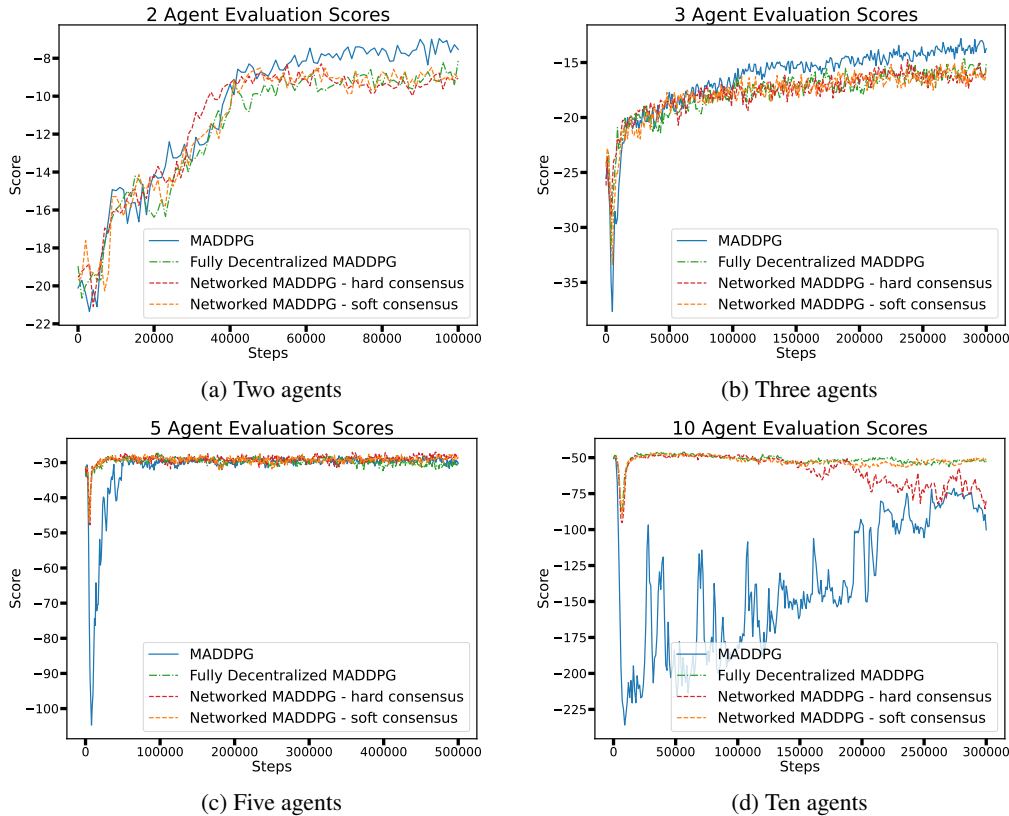


Figure 1: Comparison of evaluation score averaged over 100 test episodes for the simple spread environment between standard MADDPG, the fully decentralized MADDPG, and the fully decentralized MADDPG with hard and soft consensus update.

## 4.4 Adversarial and Mixed Setting

### 4.4.1 Simple Adversary Environment

For testing adversarial and mixed settings, we use the ‘simple adversary’ environment of the pettingzoo library [4]. There are  $N$  good agents, one adversary, and  $N$  landmarks from which one is a target landmark. All agents observe the relative position of landmarks and other agents. In addition, the good agents (but not the adversary) observe the relative position to the target. The adversary is rewarded by its distance to the target, and the good agents by their minimum distance. To avoid the adversary following the others to the target, the good agents need to split up to deceive the adversary. The rewards are

$$R_j = - \min_{j=1, \dots, N} \{ \|\ell_{\text{target}} - p_j\|_2 \} - R_{\text{adv}}$$
$$R_{\text{adv}} = - \|\ell_{\text{target}} - p_{\text{adv}}\|_2,$$

where  $R_j$  are the rewards of the good agents,  $p_j$  are the positions of the good agents and  $\ell_{\text{target}}$  is the position of the target landmark, whereas  $R_{\text{adv}}$  and  $p_{\text{adv}}$  are the reward and the position of the adversarial.

### 4.4.2 Simple Adversary Experiment

In this set of experiments we trained the agents in a 1 vs. 1 and a 1 vs. 2 setting for  $10^5$  and  $3 \cdot 10^5$  steps respectively. The plots are included in the appendix B.4, as they show little information. Even after extensive hyper parameter exploration, the agents do not learn effective policies at all. Since this is a shortcoming of the underlying MADDPG algorithm, we cannot expect our algorithms to improve on this.

## 5 Discussion and Outlook

We developed and evaluated three fully decentralized versions of the MADDPG algorithm, which can be applied to cooperative, adversary, and mixed settings. The first algorithm, 2, uses an approach of learning surrogate policies to approximate the behavior of the other agents while not relying on further information than the own observations. The second and third algorithm, 3 and 4 additionally use a communication network to exchange critic parameters. The hard consensus update averages the critics of neighboring agents after each learning step, while the soft consensus update includes a consensus penalty term.

Overall, all three algorithms show promising results that are similar to the performance of the original MADDPG algorithm. We noticed that fully-decentralized MADDPG with networked agents and hard consensus update can be unstable if hyperparameters are not carefully chosen. This problem did not occur with the soft consensus step. In our experiments, all algorithms are able to learn effective policies whenever the original MADDPG algorithm is able to do so. However, we observed shortcomings of the original MADDPG algorithm in the cooperative setting for more than three agents and in the adversarial setting in general.

Even with the original MADDPG algorithm, learning a good policy is difficult and the results often lack real cooperation. By extensive experimentation, we concluded that neither the network architecture nor the hyperparameters were at fault. Therefore, in future work, we would like to apply the decentralization techniques we developed, i.e. learning joint surrogate policies and networked agents approaches, to other algorithms such as MAPPO, which perform considerably better in MPE environments [5].



## References

- [1] Kai Cui, Sascha H. Hauck, Christian Fabian, and Heinz Koepl. Learning decentralized partially observable mean field control for artificial collective behavior. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=99tKiMVJhY>.
- [2] Ryan Lowe, Yi I 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*, 30, 2017.
- [3] Frans A. Oliehoek. *Decentralized POMDPs*, pages 471–503. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012. ISBN 978-3-642-27645-3. doi: 10.1007/978-3-642-27645-3\_15. URL [https://doi.org/10.1007/978-3-642-27645-3\\_15](https://doi.org/10.1007/978-3-642-27645-3_15).
- [4] J Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:15032–15043, 2021.
- [5] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative, multi-agent games, 2022. URL <https://arxiv.org/abs/2103.01955>.
- [6] Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pages 5872–5881. PMLR, 2018.
- [7] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms, 2021. URL <https://arxiv.org/abs/1911.10635>.

## A Appendix

### A.1 MADDPG

MADDPG as presented in [2] is a version of an actor-critic algorithm for MARL with deterministic continuous policies for each agent.

**Theorem 1 (Policy gradient theorem)** *Let  $J(\theta_i) = \mathbb{E}_{s \sim P, a \sim \pi_{\theta_i}} [R^i]$  the expected reward of each agent  $i$ . Then we can write the gradient of the policy as:*

$$\nabla_{\theta_i} J(\theta_i) = \mathbb{E}_{s \sim \mu, a_i \sim \pi_i} [\nabla_{\theta_i} \log \pi_{\theta_i}(a_i | o_i) \nabla_{a_i} Q_i^{\pi}(\mathbf{x}, a_1, \dots, a_N)], \quad (\text{A.1})$$

where  $\mathbf{x} = (o_1, \dots, o_N)$  are the observations of all agents,  $Q_i^{\pi}(\mathbf{x}, a_1, \dots, a_N)$  is a centralized action-value function,  $\pi_{\theta_i}(a_i | o_i)$  is the parameterized policy of agent  $i$  utilizing only local observations and  $\theta_i$  are the policy parameters for the agents.

The authors of [2] extend this framework to deterministic policies  $\mu_{\theta} : \mathcal{S} \rightarrow \mathcal{A}$ . By applying Theorem 1 to deterministic policies Lowe et al. find the actor update to be

$$\nabla_{\theta_i} J(\boldsymbol{\mu}_{\theta_i}) = \mathbb{E}_{\mathbf{x}, a_i \sim \mathcal{D}} \left[ \nabla_{\theta_i} \boldsymbol{\mu}_{\theta_i}(a_i | o_i) \nabla_{a_i} Q_i^{\boldsymbol{\mu}}(\mathbf{x}, a_1, \dots, a_N) \Big|_{a_i = \boldsymbol{\mu}_{\theta_i}(o_i)} \right]. \quad (\text{A.2})$$

Further, they update the critics  $Q_i^{\boldsymbol{\mu}}$  by minimizing the losses

$$\mathcal{L}(\theta_i) = \mathbb{E}_{\mathbf{x}, a_i \sim \mathcal{D}} \left[ (y - Q_i^{\boldsymbol{\mu}}(\mathbf{x}, a_1, \dots, a_N))^2 \right], \quad (\text{A.3})$$

where

$$y = r_i + \gamma Q_i^{\boldsymbol{\mu}'}(\mathbf{x}', a'_1, \dots, a'_N) \Big|_{a'_k = \boldsymbol{\mu}'_k(o_k)} \quad (\text{A.4})$$

where  $\boldsymbol{\mu}' = \{\boldsymbol{\mu}'_{\theta_1}, \dots, \boldsymbol{\mu}'_{\theta_N}\}$  is the set of target policies with delayed parameters  $\theta'_i$ . Actor, critic and target functions are all approximated with neural networks and the gradients are estimated using samples from a replay buffer  $(\mathbf{x}, \mathbf{a}, \mathbf{x}', \mathbf{a}')$  which stores past interactions between agents and environment.

## A.2 Algorithms: Variants of Multi-Agent Deep Deterministic Policy Gradient

Algorithm 1 shows the original algorithm for Multi-Agent Deep Deterministic Policy Gradient (MADDPG) from [2].

Algorithm 2 is our fully-decentralized variant of the MADDPG algorithm, while algorithms 3 and 4 are the networked approaches with hard and soft consensus update respectively. Algorithms 5 to 7 show our adaptations of the previously discussed algorithms to the mixed and adversary setting.

---

### Algorithm 1 Multi-Agent Deep Deterministic Policy Gradient

---

```

for episode = 1 to  $M$  do
  Initialize a random process  $\mathcal{N}$  for action exploration
  Receive initial state  $\mathbf{x}$ 
  for  $t = 1$  to max-episode-length do
    for each agent  $i$ , select action  $a_i = \mu_{\theta_i}(o_i) + \mathcal{N}_t$  w.r.t. the current policy and exploration
    Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r$  and next state  $\mathbf{x}'$ 
    Store transition  $(\mathbf{x}, a, r, \mathbf{x}')$  in replay buffer  $\mathcal{D}$ 
     $\mathbf{x} \leftarrow \mathbf{x}'$ 
    for agent  $i = 1$  to  $N$  do
      Sample a random minibatch of  $S$  samples  $\{(\mathbf{x}^j, a^j, r^j, \mathbf{x}'^j)\}$  from  $\mathcal{D}$ 
      Set  $y^j = r^j + \gamma Q_i^{\mu'}(\mathbf{x}'^j, a'_1, \dots, a'_N)|_{a'_k = \mu'_k(o'_k)}$ 
      Update critic by minimizing the loss:  $\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j (y^j - Q_i^\mu(\mathbf{x}^j, a^j))^2$ 
      Update the actor policy using the sampled policy gradient:
      
$$\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} Q_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_N^j)|_{a_i = \mu_i(o_i^j)}$$

    end for
    Update actor and critic target network parameters for each agent  $i$ :
    
$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$$

  end for
end for

```

---

---

**Algorithm 2** Fully-Decentralized MADDPG

---

**for** episode = 1 to  $M$  **do**  
  Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
  Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
  **for**  $t = 1$  to max-episode-length **do**  
    **for** agent  $i = 1$  to  $N$  **do**  
      Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
      Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
       $\mathbf{x}' := (o'_1, \dots, o'_N)$   
      Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$   
    **end for**  
     $\mathbf{x} \leftarrow \mathbf{x}'$   
    **if**  $t \bmod \text{learning-interval} = 0$  **then**  
      **for** agent  $i = 1$  to  $N$  **do**  
        Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$   
        Set  $y_i^j = r_i^j + \gamma Q_i^{\mu_i}(o_i'^j, a_i^j)|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i'^j)]_k}$ .  
        Update critic by minimizing the loss:  
          
$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\mu_i}(o_i^j, a_i^j) \right)^2$$
  
        Update the actor policy using the sampled policy gradient:  
          
$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k} \right)$$
  
      **end for**  
      Update actor and critic target network parameters for each agent  $i$ :  
          
$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$
  
    **end if**  
  **end for**  
**end for**

---

---

**Algorithm 3** Fully-Decentralized MADDPG with Hard Consensus Update
 

---

**for** episode = 1 to  $M$  **do**  
 Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
 Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
**for**  $t = 1$  to max-episode-length **do**  
   **for** agent  $i = 1$  to  $N$  **do**  
     Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
     Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
        $\mathbf{x}' := (o'_1, \dots, o'_N)$   
     Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$   
   **end for**  
    $\mathbf{x} \leftarrow \mathbf{x}'$   
   **if**  $t \bmod \text{learning-interval} = 0$  **then**  
     **for** agent  $i = 1$  to  $N$  **do**  
       Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$   
       Set  $y_i^j = r_i^j + \gamma Q_i^{\mu_i}(o_i'^j, a_i^j)|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k}$ .  
       Update critic by minimizing the loss:  
       
$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\mu_i}(o_i^j, a_i^j) \right)^2$$
  
       Update the actor policy using the sampled policy gradient:  
       
$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k} \right)$$
  
     **end for**  
     Update actor and critic target network parameters for each agent  $i$ :  
       
$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$
  
     *Consensus update:*  
     **for** agent  $i = 1$  to  $N$  **do**  
       Update critic parameters:  
       
$$\mu_i \leftarrow \sum_{j=1}^N C_t(i, j) \mu_j$$
  
     **end for**  
   **end if**  
   **end for**  
**end for**

---

---

**Algorithm 4** Fully-Decentralized Networked MADDPG with Soft Consensus Update
 

---

**for** episode = 1 to  $M$  **do**  
 Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
 Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
**for**  $t = 1$  to max-episode-length **do**  
   **for** agent  $i = 1$  to  $N$  **do**  
     Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
     Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
        $\mathbf{x}' := (o'_1, \dots, o'_N)$   
     Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$   
   **end for**  
    $\mathbf{x} \leftarrow \mathbf{x}'$   
   **if**  $t \bmod \text{learning-interval} = 0$  **then**  
     **for** agent  $i = 1$  to  $N$  **do**  
       Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$   
       Set  $y_i^j = r_i^j + \gamma Q_i^{\boldsymbol{\mu}^i}(o_i^j, a_i^j)|_{a_k^j = [\boldsymbol{\mu}'_{\theta_i}(o_i'^j)]_k}$ .  
       Update critic by minimizing the loss:  
       
$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\boldsymbol{\mu}^i}(o_i^j, a_i^j) \right)^2 + \sum_{j=1}^N \mathcal{C}_t(i, j) \|\mu_i - \mu_j\|$$
  
       Update the actor policy using the sampled policy gradient:  
       
$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\boldsymbol{\mu}^i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k} \right)$$
  
     **end for**  
     Update actor and critic target network parameters for each agent  $i$ :  
       
$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$$
  
   **end if**  
**end for**  
**end for**

---

---

**Algorithm 5** Fully-Decentralized MADDPG for Mixed Settings
 

---

**for** episode = 1 to  $M$  **do**  
 Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
 Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
**for**  $t = 1$  to max-episode-length **do**  
   **for** agent  $i = 1$  to  $N$  **do**  
     Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
     Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
        $\mathbf{x}' := (o'_1, \dots, o'_N)$   
     Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$

**end for**

$\mathbf{x} \leftarrow \mathbf{x}'$

**if**  $t \bmod \text{learning-interval} = 0$  **then**

**for** agent  $i = 1$  to  $N$  **do**

    Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$

    Set  $y_i^j = r_i^j + \gamma Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k}$ .

    Update critic by minimizing the loss:

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\mu_i}(o_i^j, a_i^j) \right)^2$$

    Update the actor policy for team actions (gradient ascent), using the sampled policy gradient with adversarial actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

    Update the actor policy for adversarial actions (gradient descent), using the sampled policy gradient with team actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

**end for**

  Update actor and critic target network parameters for each agent  $i$ :

$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$

**end if**

**end for**

**end for**

---

---

**Algorithm 6** Fully-Decentralized Networked MADDPG with Hard Consensus Update, Mixed
 

---

**for** episode = 1 to  $M$  **do**  
 Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
 Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
**for**  $t = 1$  to max-episode-length **do**  
   **for** agent  $i = 1$  to  $N$  **do**  
     Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
     Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
        $\mathbf{x}' := (o'_1, \dots, o'_N)$   
     Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$

**end for**

$\mathbf{x} \leftarrow \mathbf{x}'$

**if**  $t \bmod \text{learning-interval} = 0$  **then**

**for** agent  $i = 1$  to  $N$  **do**

    Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$

    Set  $y_i^j = r_i^j + \gamma Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}'_{\theta_i}(o_i^j)]_k}$ .

    Update critic by minimizing the loss:

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\mu_i}(o_i^j, a_i^j) \right)^2$$

    Update the actor policy for team actions (gradient ascent), using the sampled policy gradient with adversarial actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

    Update the actor policy for adversarial actions (gradient descent), using the sampled policy gradient with team actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

**end for**

  Update actor and critic target network parameters for each agent  $i$ :

$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$

*Consensus update:*

**for** agent  $i = 1$  to  $N$  **do**

    Update critic parameters:

$$\mu_i \leftarrow \sum_{j=1}^N \mathbf{C}_t(i, j) \mu_j$$

**end for**

**end if**

**end for**

**end for**

---



---

**Algorithm 7** Fully-Decentralized Networked MADDPG with Soft Consensus Update, Mixed
 

---

**for** episode = 1 to  $M$  **do**  
 Initialize random process  $\mathcal{N} \in \mathbb{R}^N$   
 Receive initial observation  $\mathbf{x} := (o_1, \dots, o_N)$   
**for**  $t = 1$  to max-episode-length **do**  
   **for** agent  $i = 1$  to  $N$  **do**  
     Select action  $a_i = [\boldsymbol{\mu}_{\theta_i}(o_i)]_i + \mathcal{N}_t$  w.r.t. the current policy and exploration  
     Execute actions  $a = (a_1, \dots, a_N)$  and observe reward  $r_i$  and next observations  
        $\mathbf{x}' := (o'_1, \dots, o'_N)$   
     Store  $(o_i, a_i, r_i, o'_i)$  in replay buffer  $\mathcal{D}^i$

**end for**

$\mathbf{x} \leftarrow \mathbf{x}'$

**if**  $t \bmod \text{learning-interval} = 0$  **then**

**for** agent  $i = 1$  to  $N$  **do**

    Sample a random minibatch of  $S$  samples  $\{(o_i^j, a_i^j, r^j, o_i'^j)\}$  from  $\mathcal{D}^i$

    Set  $y_i^j = r_i^j + \gamma Q_i^{\mu_i}(o_i^j, a_i^j) \Big|_{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k}$ .

    Update critic by minimizing the loss:

$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_{j=1}^S \left( y_i^j - Q_i^{\mu_i}(o_i^j, a_i^j) \right)^2 + \sum_{j=1}^N \mathcal{C}_t(i, j) \|\mu_i - \mu_j\|$$

    Update the actor policy for team actions (gradient ascent), using the sampled policy gradient with adversarial actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a^j) \Bigg|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

    Update the actor policy for adversarial actions (gradient descent), using the sampled policy gradient with team actions fixed:

$$\nabla_{\theta_i} J \approx \nabla_{\theta_i} \left( \frac{1}{S} \sum_{j=1}^S Q_i^{\mu_i}(o_i^j, a^j) \Bigg|_{\substack{a_k^j = [\boldsymbol{\mu}_{\theta_i}(o_i^j)]_k \\ a_k^j = [a_{\text{obs}}^j]_k}} \right),$$

**end for**

  Update actor and critic target network parameters for each agent  $i$ :

$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$

**end if**

**end for**

**end for**

---

## B Further experimentation

### B.1 Simple spread experiment for four agents

As mentioned in section 4.3.2 the evaluation scores for four agents look similar to those of five.

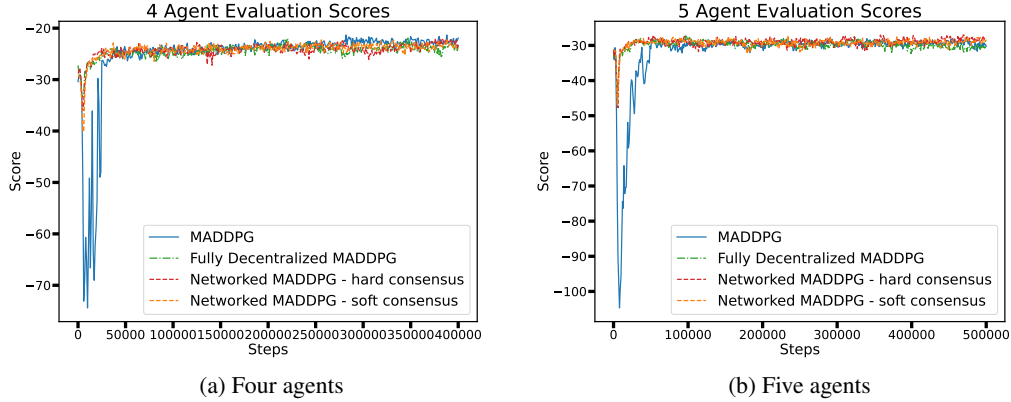


Figure 2: Comparison of evaluation scores with four and five agents for the simple spread environment. Training was done for two agents (a) over 400000 steps and for five agents (b) over 500000.

### B.2 Changing the communication hyperparameter

Altering the communication hyperparameter (see section 4.2) from  $\eta = 0.001$  to  $\eta = 0.05$  results in slower learning of the networked algorithms and instability of the hard consensus update at  $N = 3$ . For larger  $N$  the hard consensus algorithm did not converge anymore while the lower rate of communication did not affect the performance of the networked algorithm with soft consensus update.

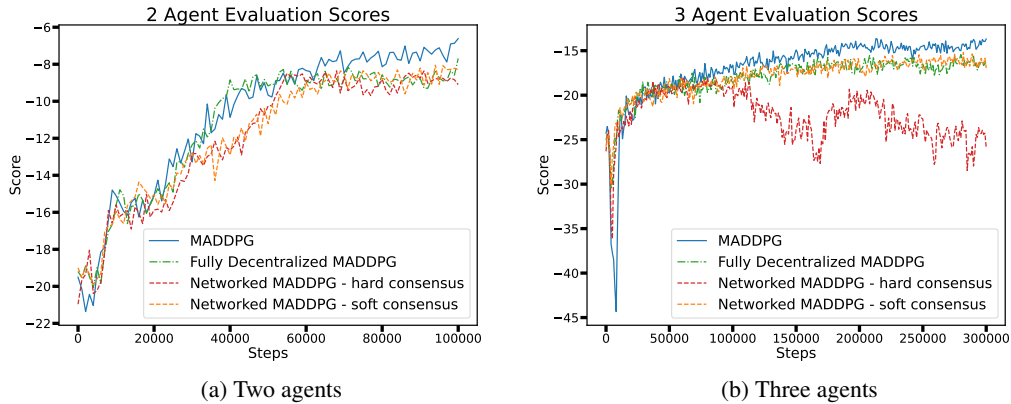


Figure 3: Comparison of evaluation scores with less connectivity over training steps for the simple spread environment between standard MADDPG (blue), the fully decentralized MADDPG (green) and the fully decentralized MADDPG with hard (red) and soft (orange) consensus update. Training was done for two agents (a) over 100000 steps and for three agents (b) over 300000.

### B.3 Sparsity of the communication matrix

A fundamental feature of the networked algorithms is the ability to control how much agents communicate with each other during training. To test whether less communication has an influence on the training result, we compared the communication matrix described in section 4.2 for the collaborative setting with a sparser matrix of the following form:

$$C_t = \begin{pmatrix} 1 - \eta & \eta & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & & \ddots & 0 \\ \vdots & & & \ddots & \eta \\ 0 & \dots & \dots & 0 & 1 - \eta \end{pmatrix}.$$

However, contrary to our expectations, the effect is negligible. Both for the fully connected and weighted communication network as well as for the circular and weighted network, the soft consensus update outperformed the hard consensus update. The computational cost however is lowered considerably by using the sparse communication network.

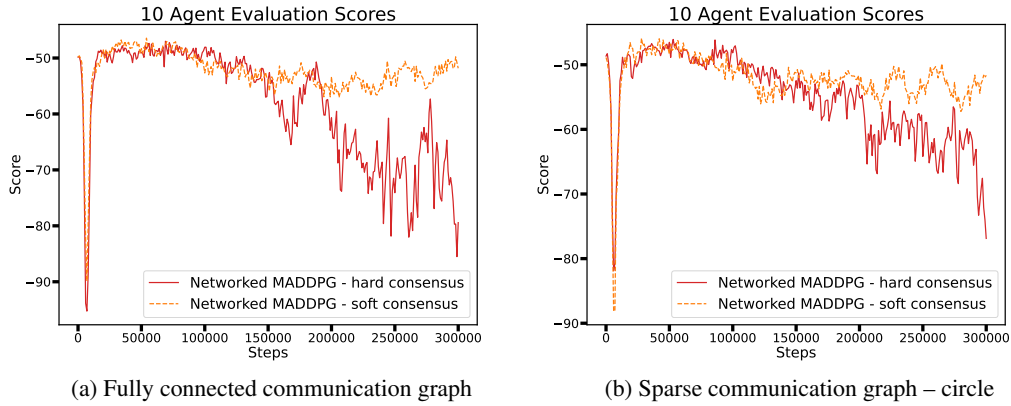


Figure 4: Comparison of evaluation scores with ten agents for the simple spread environment with a fully connected communication graph (4a) and a circle (4b).

## B.4 Adversarial and mixed settings

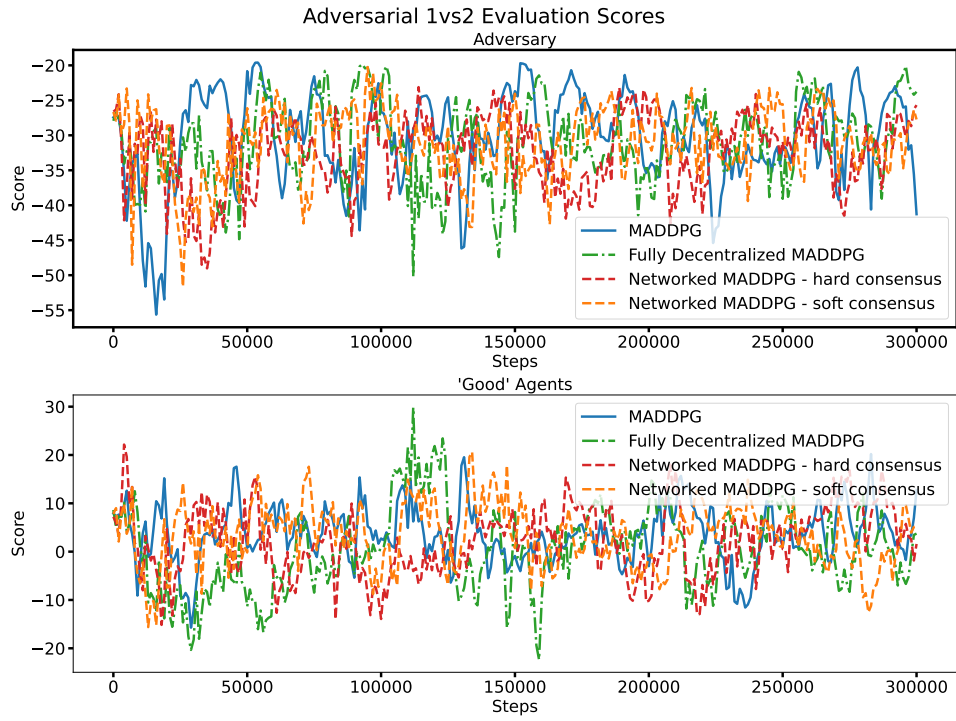


Figure 6: Performance of the algorithms in a mixed setting of one agent against two. The evaluation scores of the single agent is plotted above while the teams scores are shown below.

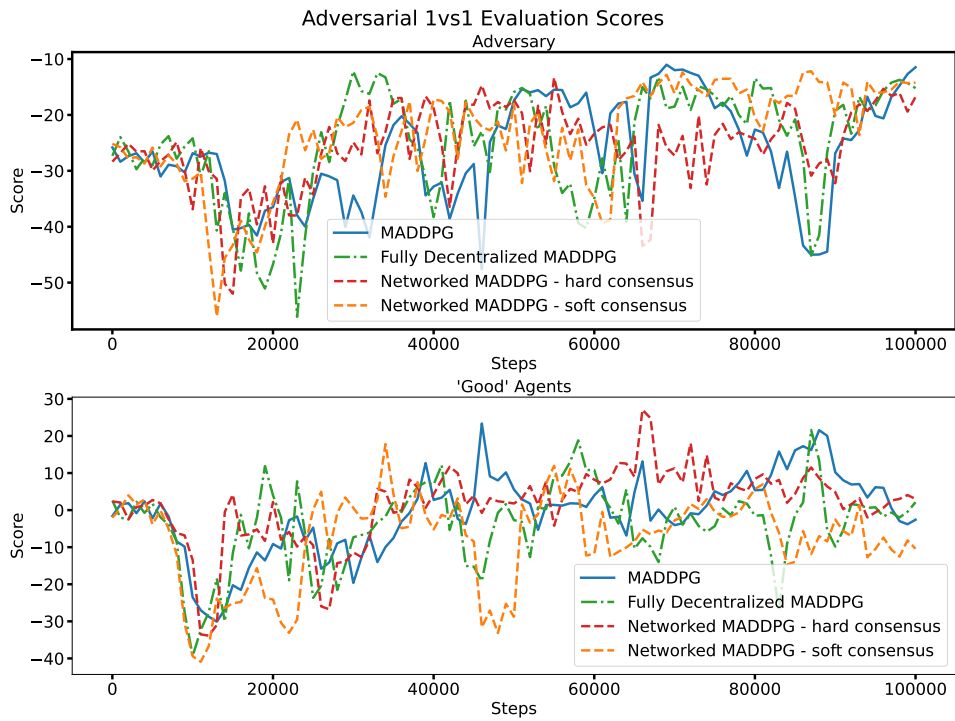


Figure 5: Performance of the algorithms in an adversarial setting of one agent against one. The two agents are shown in separate plots.