# Deep Multi-Agent Reinforcement Learning for Decentralized Continuous Cooperative Control

**Christian Schroeder de Witt*** [1]  **Bei Peng*** [2]  **Pierre-Alexandre Kamienny** [3]  **Philip Torr** [1]  **Wendelin Böhmer** [2]
**Shimon Whiteson** [2]

## Abstract

Deep multi-agent reinforcement learning (MARL) holds the promise of automating many real-world cooperative robotic manipulation and transportation tasks. Nevertheless, decentralised cooperative robotic control has received less attention from the deep reinforcement learning community, as compared to single-agent robotics and multi-agent games with discrete actions. To address this gap, this paper introduces Multi-Agent Mujoco, an easily extensible multi-agent benchmark suite for robotic control in continuous action spaces. The benchmark tasks are diverse and admit easily configurable partially observable settings. Inspired by the success of single-agent continuous value-based algorithms in robotic control, we also introduce COMIX, a novel extension to a common discrete action multi-agent $Q$-learning algorithm. We show that COMIX significantly outperforms state-of-the-art MADDPG on a partially observable variant of a popular particle environment and matches or surpasses it on Multi-Agent Mujoco. Thanks to this new benchmark suite and method, we can now pose an interesting question: what is the key to performance in such settings, the use of value-based methods instead of policy gradients, or the factorisation of the joint $Q$-function? To answer this question, we propose a second new method, FacMADDPG, which factors MADDPG's critic. Experimental results on Multi-Agent Mujoco suggest that factorisation is the key to performance.

[1]Dept. of Engineering Science, University of Oxford, United Kingdom [2]Dept. of Computer Science, University of Oxford, United Kingdom [3]Facebook AI Research, Paris, France. *: Equal Contribution. Correspondence to: Christian Schroeder de Witt <christian.schroeder@stcatz.ox.ac.uk>.

## 1. Introduction

While reinforcement learning (RL) has shown promise in learning optimal control policies for a variety of single-agent robot control problems, ranging from idealised multi-joint simulations (Todorov et al., 2012; Gu et al., 2016; Haarnoja et al., 2018) to complex grasping control problems (Kalashnikov et al., 2018; Andrychowicz et al., 2020), many real-world robot control tasks can be naturally framed as multiple decentralised collaborating agents. Cooperative manipulation tasks arise in autonomous aerial construction (Augugliaro et al., 2013; 2014), industrial manufacturing (Caccavale & Uchiyama, 2008), and agricultural robotics (Shamshiri et al., 2018) and have so far received comparatively little attention from the deep RL community.

Cooperative robotics present a number of challenges when compared to many conventional multi-agent tasks. For example, unlike in established multi-agent benchmarks, such as StarCraft II (Vinyals et al., 2017; Samvelyan et al., 2019), robotic actuators are usually continuous, so learning algorithms must scale to large continuous joint action spaces. Furthermore, such tasks are often partially observable, which arises from varying sensory equipment, included limited fields of view, together with communication constraints due to latency, power or environmental limitations (Ong et al., 2009). In fact, many such applications require fully decentralised policies for safety reasons, as communication cannot be guaranteed under all circumstances (Takadama et al., 2003).

Even when execution must be decentralised, deep reinforcement learning policies are typically trained in a centralised fashion in a simulator or laboratory. The framework of *Centralised Training with Decentralised Execution* (CTDE) (Oliehoek et al., 2008; Kraemer & Banerjee, 2016) allows policy training to exploit extra information that is not available during execution in order to accelerate learning (Lowe et al., 2017; Foerster et al., 2018; Rashid et al., 2018).

Although some multi-agent benchmark environments for continuous control exist (Leibo et al., 2017; Liu et al., 2019), few environments specialise in cooperative control and even fewer model partial observability. Moreover, existing bench-

marks, like the popular Multi-agent Particle Environment (Leibo et al., 2017), are not complex enough to meaningfully compare methods intended for robotic control. After performing a comprehensive search, we decided to introduce Multi-Agent Mujoco, a novel benchmark that fits our requirements of being a diverse, publicly available, partially observable cooperative robotics simulation that effectively captures the nature of cooperative robotic manipulation tasks.

Starting from the popular fully observable single-agent robotic control suite Mujoco (Todorov et al., 2012) included with OpenAI Gym (Brockman et al., 2016), we decompose single robots into individual segments controlled by different agents. We introduce partial observability during execution by allowing the user to limit the observation distance within the model graph formed by joints and limbs, as well as allowing fine-grained control over which segment attributes may be observed at different distance levels. Multi-Agent Mujoco is available as open-source software.

While there is a diverse portfolio of multi-agent algorithms for cooperative tasks with discrete action spaces (Foerster et al., 2018; Rashid et al., 2018; Schroeder de Witt et al., 2019), decentralised continuous control algorithms have been largely limited to deep deterministic policy gradient approaches (Lowe et al., 2017). Single-agent $Q$-learning approaches to continuous control exist, but the involved greedy action maximisation usually requires strong constraints on the functional form of the $Q$-values (Gu et al., 2016; Amos et al., 2017) or an approximate maximisation procedure based on search heuristics (Kalashnikov et al., 2018). Neither approach scales well when joint action spaces are large and values poorly approximated by constrained functions, as might be expected in cooperative multi-agent robotic tasks.

To this end, we introduce a novel $Q$-learning algorithm, COMIX, that employs a decentralisable joint action-value function with a per-agent factorisation (Rashid et al., 2018). This allows the application of the cross-entropy method (Kalashnikov et al., 2018) for greedy action maximisation on a per-agent basis, circumventing scaling issues related to large joint action spaces. Importantly, we find that COMIX significantly outperforms the state-of-the-art MADDPG in a partially observable variant of a Continuous Predator-Prey toy environment (Lowe et al., 2017). We then benchmark both on Multi-Agent Mujoco, which is a more realistic decentralisable cooperative robotic control setting. We find that COMIX outperforms MADDPG in two of the three scenarios tested and performs similarly to MADDPG in the third one. These results suggest that continuous $Q$-learning is a compelling alternative to deterministic policy gradients for decentralised cooperative multi-agent tasks.

Thanks to this new benchmark suite and method, we can now pose an interesting question: what is the key to per-

formance in such settings, the use of value-based methods instead of policy gradients, or the factorisation of the joint $Q$-function? To answer this question, we introduce Factored MADDPG (FacMADDPG), a novel variant of MADDPG where the centralised critic is factored into individual critic networks similarly to COMIX. We find that, interestingly, FacMADDPG performs similarly to COMIX in the Predator-Prey toy environment, as well as on a single Multi-Agent Mujoco environment. This suggests that it is indeed the value-function factorisation that is key to performance in tasks such as these.

## 2. Background

We consider a *fully cooperative multi-agent task* in which a team of cooperative agents choose sequential actions in a stochastic, partially observable environment. It can be modeled as a *decentralised partially observable Markov decision process* (Dec-POMDP Oliehoek et al., 2016), defined by a tuple $\langle \mathcal{N}, \mathcal{S}, \mathcal{U}, P, r, \mathcal{Z}, O, \rho, \gamma \rangle$. Here $\mathcal{N} := \{1, \ldots, N\}$ denotes the set of $N$ agents and $s \in \mathcal{S}$ describes the discrete or continuous state of the environment. The initial state $s_0 \sim \rho$ is drawn from distribution $\rho$, and at each time step $t$, all agents $a \in \mathcal{N}$ choose simultaneously discrete or continuous actions $u_t^a \in \mathcal{U}$, yielding the joint action $\boldsymbol{u}_t := \{u_t^a\}_{a=1}^N \in \mathcal{U}^N$. After executing the joint action $\boldsymbol{u}_t$ in state $s_t$, the next state $s_{t+1} \sim P(s_t, \boldsymbol{u}_t)$ is drawn from transition kernel $P$ and the collaborative reward $r_t = r(s_t, \boldsymbol{u}_t)$ is returned to the team.

In a Dec-POMDP, the true state of the environment cannot be directly observed by the agents. Each agent $a \in \mathcal{N}$ draws an individual observation $z_t^a \in \mathcal{Z}, \boldsymbol{z}_t := \{z_t^a\}_{a=1}^N$, from the observation kernel $O(s_t, a)$. The history of an agent's observations and actions is denoted $\tau_t^a \in \mathcal{T}_t := (\mathcal{Z} \times \mathcal{U})^t \times \mathcal{Z}$, and the set of all agents' histories $\boldsymbol{\tau}_t := \{\tau_t^a\}_{a=1}^N$. Agent $a$ chooses its actions with a decentralised policy $u_t^a \sim \pi(\cdot | \tau_t^a)$ based only on its individual history. The collaborating team of agents aim to learn a *joint policy* $\pi(\boldsymbol{u} | \boldsymbol{\tau}_t) := \prod_{a=1}^N \pi^a(u^a | \tau_t^a)$ that maximises their expected discounted return, $\mathbb{E}[\sum_{t=0}^\infty \gamma^t r_t]$. This joint policy induces a joint action-value function $Q^\pi$ which estimates the expected discounted return when the agents take joint action $\boldsymbol{u}_t$ with histories $\boldsymbol{\tau}_t$ in state $s_t$ and then follow some joint policy $\pi$:

$$Q^\pi(s_t, \boldsymbol{\tau}_t, \boldsymbol{u}_t) := \mathbb{E}\Big[ \sum_{i=0}^\infty \gamma^i r_{t+i} \; \Big|\; {}^{s_{k+1}\sim P(s_k, \boldsymbol{u}_k)}_{{}^{\boldsymbol{z}_{k+1}\sim O(s_{k+1},\cdot)}_{\boldsymbol{u}_{k+1}\sim \pi(\cdot|\boldsymbol{\tau}_{k+1})}}, \forall k \geq t \Big]$$
$$= r_t + \gamma \mathbb{E}\big[ Q^\pi(s_{t+1}, \boldsymbol{\tau}_{t+1}, \boldsymbol{u}_{t+1}) \; \big|\; {}^{s_{t+1}\sim P(s_t, \boldsymbol{u}_t)}_{\boldsymbol{u}_{t+1}\sim \pi(\cdot|\boldsymbol{\tau}_{t+1})} \big], \quad (1)$$

where $\gamma \in [0, 1)$ is a discount factor.

### 2.1. Deep $Q$-Learning

*Deep $Q$-Network* (DQN Mnih et al., 2015) uses a deep neural network to estimate the action-value function,

$Q(s, \boldsymbol{z}_t, \boldsymbol{u}; \theta) \approx \max_\pi Q^\pi(s, \boldsymbol{\tau}, \boldsymbol{u})$, where $\theta$ are the parameters of the network. For the sake of simplicity, we restrict us here to feed-forward networks, which condition on the last observations $\boldsymbol{z}_t^a$, rather than the entire agent histories $\boldsymbol{\tau}_t$. The network parameters $\theta$ are trained by gradient descent on the mean squared regression loss:

$$\mathcal{L}_{[\theta]}^{\text{DQN}} := \mathbb{E}_\mathcal{D}\left[\left(y_t^{\text{DQN}} - Q(s_t, \boldsymbol{z}_t, \boldsymbol{u}_t; \theta)\right)^2\right], \quad (2)$$
$$y_t^{\text{DQN}} := r_t + \gamma \max_{\boldsymbol{u}'} Q(s_{t+1}, \boldsymbol{z}_{t+1}, \boldsymbol{u}'; \theta^-),$$

where the expectation is estimated with transitions $(s_t, \boldsymbol{z}_t, \boldsymbol{u}_t, r_t, s_{t+1}, \boldsymbol{z}_{t+1}) \sim \mathcal{D}$ sampled from an *experience replay buffer* $\mathcal{D}$ (Lin, 1992b). The use of replay buffer reduces correlations in the observation sequence. To further stabilise learning, $\theta^-$ denotes parameters of a *target network* that are only periodically copied from the most recent $\theta$.

## 2.2. Centralised Training with Decentralised Execution

A simple and natural approach to solving Dec-POMDPs is to let each agent $a$ learn an individual action-value function $Q_a$ independently, as in *independent Q-learning* (IQL Tan, 1993). IQL serves as a surprisingly strong benchmark in both cooperative and competitive MARL tasks with discrete actions (Tampuu et al., 2017). However, IQL has no convergence guarantees since, as agents independently explore and update their policies, the environment becomes nonstationary from each agent's viewpoint. An alternative solution is to employ *centralised training with decentralised execution* (CTDE Kraemer & Banerjee, 2016). CTDE allows the learning algorithm to access all local action-observation histories $\mathcal{T}$ and global state $s$, as well as share gradients and parameters, but each agent's executed policy can condition only on its own action-observation history $\tau^a \in \mathcal{T}_t$.

## 2.3. VDN and QMIX

Value decomposition networks (VDN Sunehag et al., 2018) and QMIX (Rashid et al., 2018) are two representative examples of value function factorisation (Koller & Parr, 1999) that aim to efficiently learn a centralised but factored action-value function. They both work in cooperative MARL tasks with discrete actions, using CTDE. To ensure consistency between the centralised and decentralised polices, VDN factors the joint action-value function $Q_{tot}(s, \boldsymbol{\tau}, \boldsymbol{u}; \boldsymbol{\theta})$ into a sum of individual action-value functions $Q_a(\tau^a, u^a; \theta^a)$,[1] one for each agent $a$, that condition only on individual action-observation histories:

$$Q_{tot}(s, \boldsymbol{\tau}, \boldsymbol{u}; \boldsymbol{\theta}) := \sum_{a=1}^N Q_a(\tau^a, u^a; \theta^a). \quad (3)$$

By contrast, QMIX represents the joint action-value function as a monotonic function of individual action-value functions. The main insight is that, to extract decentralised policies that are consistent with their centralised counterparts, it suffices to constrain $Q_{tot}$ to be monotonic in each $Q_a$: $\frac{\partial Q_{tot}}{\partial Q_a} \geq 0, \forall a \in \mathcal{N}$. Thus, in QMIX, the joint action-value function $Q_{tot}$ is represented as:

$$Q_{tot}(s, \boldsymbol{\tau}, \boldsymbol{u}; \boldsymbol{\theta}, \phi) := f_\phi\left(s, \{Q_a(\tau^a, u^a; \theta^a)\}_{a=1}^N\right), \quad (4)$$

where $f_\phi$ is a *mixing network* that takes as input the *agent network* outputs $Q_a$ and mixes them monotonically, producing the values of $Q_{tot}$. Monotonicity can be guaranteed by non-negative mixing weights. These weights are generated by separate *hypernetworks* (Ha et al., 2016), parameterised by $\phi$, which condition on the full state $s$. This allows it to learn different monotonic mixing weights in each state.

In both methods, the loss function is analogous to the standard DQN loss of (2), where $Q$ is replaced by $Q_{tot}$. During execution, each agent selects actions greedily with respect to its own $Q_a$.

## 2.4. MADDPG

*Multi-agent deterministic policy gradient* (MADDPG Lowe et al., 2017) is an actor-critic method that works in both cooperative and competitive MARL tasks with discrete or continuous action spaces. MADDPG was originally designed for the more general case of *partially observable stochastic games* (Kuhn, 1953). Here we discuss a version specific to Dec-POMDPs and consider continuous actions. We assume each agent $a$ has a deterministic policy $\mu^a$, parameterised by $\theta^a$, with $\boldsymbol{\mu}(\boldsymbol{\tau}; \boldsymbol{\theta}) := \{\mu^a(\tau^a; \theta^a)\}_{a=1}^N$. MADDPG learns a centralised critic $Q_a^{\boldsymbol{\mu}}(s, \boldsymbol{u}; \phi^a)$ for each agent $a$ that conditions on the full state $s$ and the joint actions $\boldsymbol{u}$ of all agents. The policy gradient for $\theta^a$ is:

$$\nabla_{\theta^a} \mathcal{L}_{[\theta^a]}^{\boldsymbol{\mu}} := -\mathbb{E}_\mathcal{D}\left[\nabla_{\theta^a}\mu^a(\tau_t^a; \theta^a)\nabla_{u^a}Q_a^{\boldsymbol{\mu}}(s_t, \hat{\boldsymbol{u}}_t^a; \phi^a)\big|_{u^a=\mu^a(\tau_t^a)}\right],$$
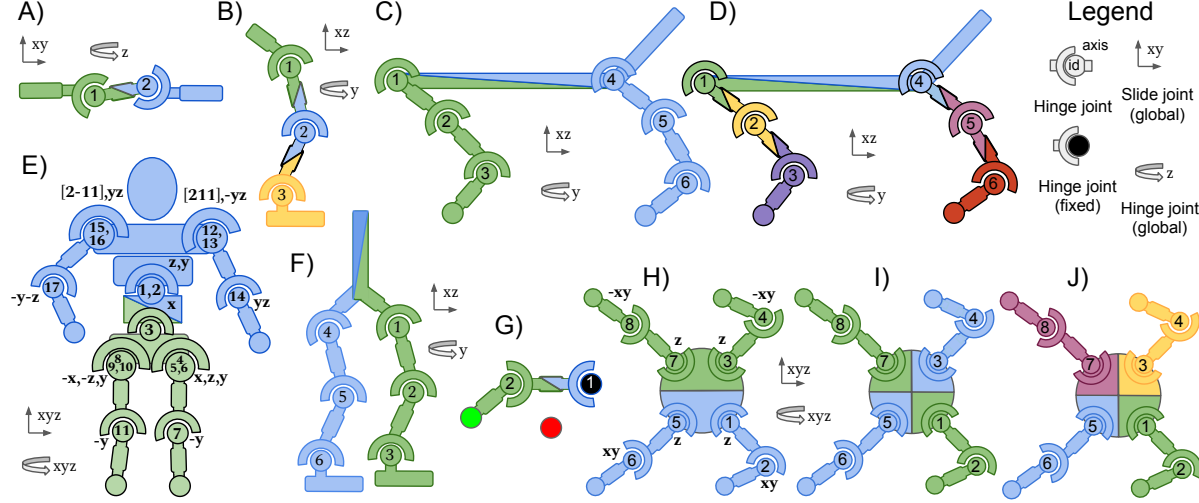
where $\hat{\boldsymbol{u}}_t^a := \{u_t^1, \dots, u_t^{a-1}, u^a, u_t^{a+1}, \dots, u_t^N\}$ and $s_t, \boldsymbol{u}_t, \boldsymbol{\tau}_t$ are sampled from a replay buffer $\mathcal{D}$. The centralised action-value function $Q_a^{\boldsymbol{\mu}}$ of each agent $a$ is trained by minimising the following loss

$$\mathcal{L}_{[\phi^a]}^{\text{DPG}} := \mathbb{E}_\mathcal{D}\left[\left(y_t^a - Q_a^{\boldsymbol{\mu}}(s_t, \boldsymbol{u}_t; \phi^a)\right)^2\right], \quad (5)$$
$$y_t^a := r_t + \gamma Q_a^{\boldsymbol{\mu}}\left(s_{t+1}, \boldsymbol{\mu}(\boldsymbol{\tau}_{t+1}; \boldsymbol{\theta}'); \phi'^a\right).$$

where transitions are sampled from a replay buffer $\mathcal{D}$ (Lin, 1992a) and $\boldsymbol{\theta}'$ and $\phi'^a$ are target-network parameters.
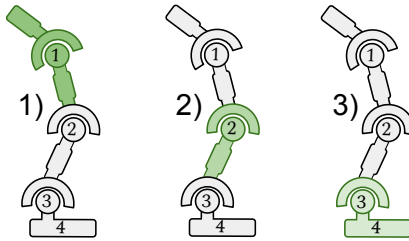
## 3. Multi-Agent Mujoco

Multi-Agent Mujoco is a novel benchmark for decentralised cooperative continuous multi-agent robotic control. Starting from the popular fully observable single-agent robotic

---

[1]Strictly speaking, each $Q_a$ is a *utility function* since by itself it does not estimate an expected return. We refer to $Q_a$ as action-value function for simplicity.

*Figure 1.* **Agent partitionings for Multi-Agent Mujoco environments:** A) 2-Agent Swimmer [2x1], B) 3-Agent Hopper [3x1], C) 2-Agent HalfCheetah [2x3], D) 6-agent HalfCheetah [6x1], E) 2-Agent Humanoid and 2-Agent HumanoidStandup (each [1x9,1x8]), F) 2-Agent Walker G) 2-Agent Reacher [2x1], H) 2-Agent Ant [2x4], I) 2-Agent Ant Diag [2x4], J) 4-Agent Ant [4x2]. Colours indicate agent partitionings. Each joint corresponds to a single controllable motor. Split partitions indicate shared body segments. Square brackets indicate [(number of agents) x (joints per agent)]. Joint IDs are in order of definition in the corresponding OpenAI Gym XML asset files (Brockman et al., 2016). Global joints indicate degrees of freedom of the center of mass of the composite robotic agent.

Mujoco (Todorov et al., 2012) control suite included with OpenAI Gym (Brockman et al., 2016), we create novel scenarios in which multiple agents have to solve a task cooperatively. This is achieved by first representing a given single robotic agent as a *body graph*, where vertices (joints) are connected by adjacent edges (body segments), as shown in Figure 1. We then partition the body graph into disjunct sub-graphs, one for each agent, each of which contains one or more joints that can be controlled.



*Figure 2.* **Observations by distance for 3-Agent Hopper (as seen from agent 1).** 1) corresponds to joints and body parts at zero graph distance from agent 1, 2) corresponds to joints and body parts observable at unit graph distance and 3) at two unit graph distances.

Hence, each agent's action space in Multi-Agent Mujoco is given by the joint action space over all motors controllable by that agent. For example, for the agent corresponding to the green partition in 2-Agent HalfCheetah (Figure 1, C) consists of three joints (joint ids 1,2 and 3) and four adjacent

body segments. Each joint has action space $[-1, 1]$, hence the joint action space of this agent $\mathcal{A}_{joint} = [-1, 1]^3$.

For each agent $a$, observations are constructed in a two-stage process. First, we infer which body segments and joints are observable by agent $a$. Each agent can always observe all joints within its own sub-graph. A configurable parameter $k \geq 0$ determines the maximum graph distance to the agent's subgraph at which joints are observable. Body segments directly attached to observable joints are themselves observable. The agent observation is then given by a fixed order concatenation of the representation vector of each observable graph element. Depending on configuration, representation vectors may include attributes such as position, velocity and external body forces.

Restricting both the observation distance $k$, as well as limiting the set of observable element categories imposes partial observability. However, task goals remain unchanged from the single-agent variants (see Table 1 in the Appendix), except that the goals must be reached collaboratively by multiple agents: we simply repurpose the original single-agent reward signal as a team reward signal.

## 4. COMIX and FacMADDPG

**COMIX.** $Q$-learning has shown considerable success in multi-agent settings with discrete action spaces (Rashid et al., 2018). However, performing greedy action selection in $Q$-learning requires evaluating $\arg\max_{\mathbf{u}} Q_{tot}(s, \boldsymbol{\tau}, \boldsymbol{u})$,

where $Q_{tot}$ is the joint state-action value function. In discrete action spaces, this operation can be performed efficiently through enumeration unless the action space is extremely large. In continuous action spaces, however, enumeration is impossible. Hence, existing continuous $Q$-learning approaches in single-agent settings either impose constraints on the form of $Q$-value to make maximisation easy (Gu et al., 2016; Amos et al., 2017), at the expense of estimation bias, or perform greedy action selection only in approximation (Kalashnikov et al., 2018). Neither approach scales easily to the large joint action spaces inherent to multi-agent settings, as 1) the joint action space grows exponentially in the number of agents, and 2) training $Q_{tot}$ to select maximal actions becomes impractical when there are many agents.

This highlights the importance of learning a centralised but factored $Q_{tot}$. To factor large joint action spaces efficiently in a decentralisable fashion, COMIX models a joint state-action value function $Q_{MIX} = f\big(s, Q_1(\tau^1, u^1; \theta^1), \dots, Q_N(\tau^N, u^N; \theta^N)\big)$, where $Q_a$ are per-agent utility functions used for greedy action selection. Similarly to QMIX (Rashid et al., 2018), COMIX imposes a monotonicity constraint on $f$ to keep joint action selection compatible with action selection from individual utility functions.

COMIX performs greedy selection of actions $u^a$ with respect to utility functions $Q_a(\tau^a, u^a; \theta^a)$ for each agent $a$ using the *cross-entropy method* (CEM De Boer et al., 2005), a sampling-based derivative-free heuristic search method that has been successfully used to find approximate maxima of nonconvex $Q$-networks in a number of single-agent robotic control tasks (Kalashnikov et al., 2018). The centralised but factored $Q_{MIX}$ allows us to use CEM to sample actions for each agent independently and to use the individual utility function $Q_a$ to guide the selection of maximal actions. We rely on CEM instead of other continuous $Q$-learning approaches (Gu et al., 2016; Amos et al., 2017) because of its empirical success (see Section 5).

CEM iteratively draws a batch of $N$ random samples from a candidate distribution $\mathcal{D}_k$, e.g., a Gaussian, at each iteration $k$. The best $M < N$ samples are then used to fit a new Gaussian distribution $\mathcal{D}_{k+1}$, and this process repeats $K$ times. For COMIX, we use a CEM hyperparameter configuration similar to Qt-Opt (Kalashnikov et al., 2018), where $N = 64$, $M = 6$, and $K = 2$.[2] Gaussian distributions are initialised with mean $\mu = 0$ and standard deviation $\sigma = 1$. Algorithm 2 outlines the full CEM process for COMIX.

**FacMADDPG.** Learning a centralised critic conditioning on a large joint agent observation space can be difficult (Iqbal & Sha, 2019). We introduce FacMADDPG, a novel

---

[2]We empirically find 2 iterations to suffice.

variant of MADDPG with an agent-specific factorisation that facilitates the learning of a centralised critic in Dec-POMDPs. In FacMADDPG, all agents share a single centralised critic that is factored as $Q^{\boldsymbol{\mu}_T} =$

$$
\begin{aligned}
g_\phi\big(s, Q_1(\tau^1, u^1, \dots u^N; \theta^1), \dots, \\
Q_N(\tau^N, u^1, \dots, u^N; \theta^N)\big)
\end{aligned}
\tag{6}
$$

where $g$ is a function represented by a monotonic mixing network. Although the monotonicity requirement on $g_\phi$ is no longer required as the critic is not used for greedy action selection, FacMADDPG does impose monotonicity on $g_\phi$ in order to keep the factorisation comparable to the one employed by COMIX. We find that FacMADDPG significantly outperforms an ablation without monotonicity constraints (see Appendix 10).

## 5. Experimental Setup

**Partially Observable Continuous Predator-Prey.** The mixed *simple tag* environment (Figure 3) introduced by Leibo et al. (2017) is a variant of the classic predator-prey game. Three slower cooperating circular agents (red), each with continuous movement action spaces $u^a \in \mathbb{R}^2$, must catch a faster circular prey (green) on a randomly generated two-dimensional toroidal plane with two large landmarks blocking the way.

To obtain a purely cooperative environment, we replace the prey's policy by a hard-coded heuristic, that, at any time step, moves the prey to the sampled position with the largest distance to the closest predator. If one of the cooperative agents collides with the prey, a team reward

---

**Algorithm 1** Algorithmic description of COMIX. The function CEM is defined in Appendix 10.

**function** COMIX
  Initialise ReplayBuffer, $\theta, \theta^-, \phi, \phi^-$
  **for** each training episode $e$ **do**
    $s_0, \mathbf{z}_0 \leftarrow \text{EnvInit}()$
    **for** $t := 0$ until $t = T$ step 1 **do**
      $\mathbf{u}_t \leftarrow \text{CEM}(Q_1, \dots, Q_N, \tau_t^1, \dots, \tau_t^N)$
      $\langle s_{t+1}, \mathbf{z}_{t+1}, r_t \rangle \leftarrow \text{EnvStep}(\mathbf{u}_t)$
      $\text{ReplayBuffer} \leftarrow \langle s_t, \mathbf{u}_t, \mathbf{z}_t, r_t, s_{t+1}, \mathbf{z}_{t+1} \rangle$
    **end for**
    $\{\langle s_i, \boldsymbol{u}_i, \boldsymbol{z}_i, r_i, s_i', \boldsymbol{z}_i' \rangle\}_{i=1}^b \sim \text{ReplayBuffer}$
    $y_i \leftarrow r_i + \gamma \max_{\boldsymbol{u}_i'} Q_{\text{tot}}(s_i', \boldsymbol{z}_i', \boldsymbol{u}_i'; \boldsymbol{\theta}^-, \phi^-), \ \forall i$
    $\mathcal{L} \leftarrow \sum_{i=1}^b \Big(y_i - Q_{\text{tot}}(s_i, \boldsymbol{z}_i, \boldsymbol{u}_i; \boldsymbol{\theta}, \phi)\Big)^2$
    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}$
    $\phi \leftarrow \phi - \alpha \nabla_\phi \mathcal{L}$
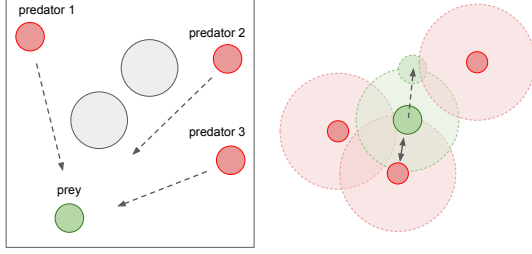  **end for**
**end function**

*Figure 3.* Continuous Predator-Prey. Left: Top-down view of toroidal plane, with predators (red), prey (green) and obstacles (grey). Right: Illustration of the prey's avoidance heuristic. Observation radii of both agents and prey are indicated.

of $+10$ is emitted; otherwise, no reward is given. In the original simple tag environment, each agent can observe the relative positions of other two agents, the relative position and velocity of the prey, and the relative positions of the landmarks. This means each agent's private observation provides an almost complete representation of the true state of the environment.

To introduce partial observability to the environment, we add an agent *view radius*, which restricts the agents from receiving information about other entities (including all landmarks, the other two agents, and the prey) that are out of range. Specifically, we set the value of view radius such that the agents can only observe other agents roughly 60% of the time. We open-source the full set of multi-agent particle environments with added partial observability.[3]

**Multi-Agent Mujoco.** All Multi-Agent Mujoco environments are configured according to the default configuration of Multi-Agent Mujoco, where each agent can observe both velocity and position of its own body parts and positions only at graph distances greater than zero. We set maximum observation distances to $k = 2$ for 2-Agent HalfCheetah and 3-Agent Hopper, $k = 1$ for 2-Agent Walker. Default team reward signals are used (see Table 1 in the Appendix).

**Ablations.** We also introduce a number of novel ablations in order to study diverse aspects of our method in isolation: 1) *COVDN:* we factor the joint action-value function $Q_{tot}$ into a sum of individual action-value functions $Q_a$ as in VDN, and use CEM to learn $Q_a$ for each agent $a$, 2) *COMIX-NAF:* we factor $Q_{tot}$ assuming mixing monotonicity as in QMIX, and add quadratic function constraints on each $Q_a$ based on Normalized Advantage Functions (NAF Gu et al., 2016), and 3) *COVDN-NAF:* we represent $Q_{tot}$ assuming additive mixing as in VDN, and add quadratic function constraints on $Q_a$ based on NAF.

---

[3] https://github.com/schroederdewitt/
multiagent_mujocohttps://github.com/
schroederdewitt/multiagent-particle-envs/

**Evaluation Procedure.** We evaluate each method's performance using the following procedure: for each run of a method, we pause training every fixed number of timesteps (2000 timesteps for Continuous Predator-Prey and 4000 timesteps for Multi-Agent Mujoco) and run 10 independent episodes with each agent performing greedy decentralised action selection. The mean value of these 10 episode returns are then used to evaluate the performance of learned policies. See Appendix 9 for further experimental details.

# 6. Empirical Results

Figure 4 shows that COMIX significantly outperforms MADDPG on Continuous Predator-Prey (Figure 4a), both in terms of absolute performance and learning speed. On Multi-Agent Mujoco, COMIX outperforms MADDPG in absolute terms on 2-Agent Walker scenario (Figure 4c), while MADDPG cannot learn it effectively. On 2-Agent HalfCheetah (Figure 4b), COMIX outperforms MADDPG in absolute terms and has lower limit variance. On 3-Agent Hopper (Figure 4d), COMIX performs similarly to MAD-DPG but has significantly lower variance across seeds.

Despite the ability to represent a richer form of coordination with its functionally unconstrained critic (Son et al., 2019), in our experiments MADDPG is not able to outperform COMIX, which uses a monotonically constrained mixing network.

We hypothesise that this is because MADDPG's critic directly conditions on the joint observations and actions of all agents. COMIX, by contrast, represents the optimal joint action-value function using a monotonic mixing function of per-agent utilities. Early in training, MADDPG's critic estimator may thus be more prone to picking up non-trivial suboptimal coordination patterns than COMIX. Such local minima might be hard to subsequently escape.

By contrast, the monotonicity constraint on COMIX's mixing network may smooth the optimisation landscape, allowing COMIX to avoid suboptimal local minima more efficiently than MADDPG. In other words, COMIX's network architecture imposes a suitable prior that captures the forms of additive-monotonic coordination required to solve these tasks.

These results an interesting question: What is the key to performance in such settings, the use of value-based methods instead of policy gradients, or the choice of factorisation of the joint $Q$-value function? Previous work on this question (Bescuca, 2019) is inconclusive due to the confounder that the policy gradient methods studied (COMA and Central-V (Foerster et al., 2018)) are on-policy, while the respective $Q$-learning method, QMIX (Rashid et al., 2018), is off-policy with experience replay (Lin, 1992a).
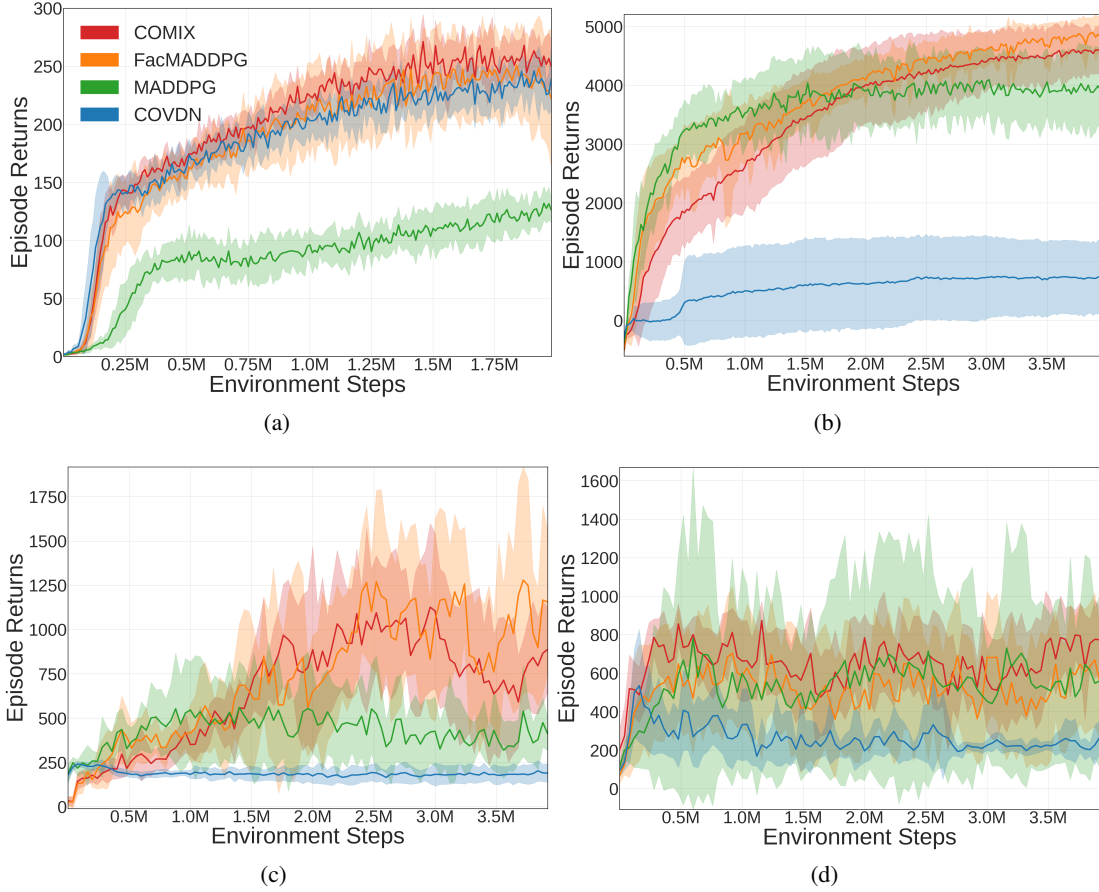
*Figure 4.* Mean episode return on (a) Continuous Predator-Prey, (b) 2-Agent HalfCheetah [2x3], (c) 2-Agent Walker [2x3], and (d) 3-Agent Hopper [3x1]. The mean across 10 seeds is plotted and the 95% confidence interval is shown shaded.

To address this question, we evaluate the performance of FacMADDPG, which factors MADDPG's critic in the same manner as COMIX's joint $Q$-value function. We find that FacMADDPG performs similarly to COMIX on both Continuous Predator-Prey and all three Multi-Agent Mujoco tasks (see Figure 4b-4d). As both COMIX and FacMAD-DPG are off-policy algorithms, this shows that the factorisation of the joint $Q$-value function is the key to performance in these decentralised continuous cooperative multi-agent tasks.

**Ablations.** We find that COVDN performs drastically worse than both COMIX and MADDPG across all Multi-Agent Mujoco tasks (shown in Figure 4a-4d), demonstrating the necessity of the non-linear mixing of agent utilities and conditioning on the state information in order to achieve competitive performance in such tasks.

Figure 5 shows that, compared to its ablations COVDN-NAF and COMIX-NAF, COMIX is noticeably more stable on Continuous Predator-Prey. COVDN-NAF has sharp drops in performance at the late stage of training, while COMIX-
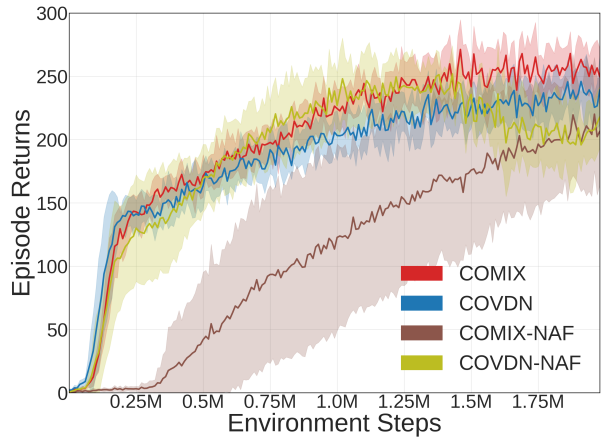


*Figure 5.* Mean episode return on Continuous Predator-Prey comparing COMIX and ablations. The mean across 10 seeds is plotted and the 95% confidence interval is shown shaded.

NAF converges significantly slower than COMIX and is much more varied across seeds. This demonstrates that greedy action selection based on CEM heuristic search is both more stable and performant than simple exact methods in practice.

Finally, we evaluate an ablation of FacMADDPG without the monotonicity constraint. We find that this method does not learn at all in Continuous Predator-Prey, and performs significantly worse than both COMIX and MADDPG on 2-Agent HalfCheetah (see Figure 7 in Appendix 10). This suggests monotonicity matters.

## 7. Related Work

While several MARL benchmarks with continuous action spaces have been released, few are simultaneously diverse, fully cooperative, decentralisable and admit partial observability. The Multi-Agent Particle suite (Lowe et al., 2017) features a few decentralisable tasks in a fully observable planar point mass toy environment. Presumably due to its focus on real-world robotic control, RoboCup Soccer Simulation (Kitano et al., 1997; Stone & Sutton, 2001; Riedmiller et al., 2009) does not currently feature an easily configurable software interface for MARL, nor suitable AI-controlled benchmark opponents. Liu et al. (2019) introduce MuJoCo Soccer Environment, a multi-agent soccer environment with continuous simulated physics that cannot be used in a purely cooperative setting and does not admit partial observability.

Of the existing environments most similar but not as diverse as Multi-Agent Mujoco, Wang et al. (2018) introduce two decomposed Mujoco environments, Centipede and Snakes, the latter of which being similar to Multi-Agent Mujoco's 2-Agent Swimmer. Ackermann et al. (2019) evaluate on a single environment that is similar to a particular configuration of 2-Agent Ant, but do not consider tasks across different numbers of agents and Mujoco scenarios.

A number of multi-agent variants of deep deterministic policy gradients (Lillicrap et al., 2015; Lowe et al., 2017) have been proposed for MARL in continuous action spaces: MADDPG-M (Kilinc & Montana, 2018) uses communication channels in order to overcome observation noise in partially observable settings. By contrast, we consider fully decentralised settings without communication. R-MADDPG (Wang et al., 2019) equips MADDPG with recurrent policies and critics in a partially observable setting with communication. As we are primarily interested in the relative performance between policy gradients and continuous $Q$-learning approaches, we employ feed-forward network architectures to avoid the complexities of recurrent network training.

NerveNet (Wang et al., 2018) achieves policy transfer across robotic agents with different numbers of repeated units. Unlike COMIX, NerveNet is not fully decentralisable as it requires explicit communication channels. Iqbal & Sha (2019) introduce MAAC, a variant of MADDPG for stochastic games, in which the centralised critics employ an attention mechanism on top of agent-specific observation embeddings. Unlike FacMADDPG, MAAC explicitly addresses settings where agents receive both individual and team rewards.

Besides VDN and QMIX, QTRAN (Son et al., 2019) allows for arbitrary utility function mixings by introducing auxiliary losses that align utility function maxima with maxima of the joint $Q$-function. Despite being more expressive, QTRAN does not scale well to complex environments, such as StarCraft II (Samvelyan et al., 2019; Bohmer et al., 2019) and may not generalise well to continuous action spaces due to the point-wise nature of its auxiliary losses.

Continuous $Q$-learning has so far been studied almost exclusively in the fully observable single-agent setting. Two distinct approaches to making greedy action selection tractable have emerged: Both Normalized Advantage Functions (NAF Gu et al., 2016) and Partially Input-Convex Neural Networks (PICNN Amos et al., 2017) constrain the functional form of the state-action value function approximator such as to guarantee an easily identifiable global maximum. In contrast, heuristic search approaches, such as Cross-Entropy Maximisation (CEM Mannor et al., 2003), forfeit global guarantees but allow for unconstrained $Q$-learning approximators. CEM (Mannor et al., 2003) has been used successfully in single-agent robotic simulations (Kalashnikov et al., 2018). As for COMIX, we find that ablations using NAF perform poorly (see section 5).

## 8. Conclusion

In order to stimulate research into decentralised cooperative robotic control, we introduce a novel benchmark suite, Multi-Agent Mujoco. Multi-Agent Mujoco consists of a diverse set of multi-agent tasks with continuous action spaces and is easily extensible. We also introduce COMIX, a novel $Q$-learning algorithm that factors the joint action space into per-agent action spaces to overcome scalability problems in continuous greedy action selection.

Our results show that COMIX performs competitively with, or even outperforms, MADDPG both on a traditional benchmark environment, as well as on Multi-Agent Mujoco. Futhermore, we introduce a second method FacMADDPG, a novel variant of MADDPG where the centralised critic is factored into individual critic networks similarly to COMIX.

We find that, interestingly, FacMADDPG performs similarly to COMIX in Predator-Prey toy environment, as well as on a single Multi-Agent Mujoco environment. This shows that the factoration of the joint $Q$-value is the key to performance in decentralised continuous cooperative multi-agent tasks.

Future work will investigate the utility of recent amortisation techniques for cross-entropy maximisation (Van de Wiele et al., 2020) and the relationship between exploration strategies in continuous $Q$-learning and deterministic policy gradient approaches. We also plan to extend Multi-Agent Mujoco to contain more challenging multi-agent environments composed of multiple robotic agents rather than decompositions of single robotic agents.

## Acknowledgements

## References

Ackermann, J., Gabler, V., Osa, T., and Sugiyama, M. Reducing overestimation bias in multi-agent domains using double centralized critics. *arXiv preprint arXiv:1910.01465*, 2019.

Amos, B., Xu, L., and Kolter, J. Z. Input convex neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 146–155. JMLR. org, 2017.

Andrychowicz, O. M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.

Augugliaro, F., Mirjan, A., Gramazio, F., Kohler, M., and D'Andrea, R. Building tensile structures with flying machines. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 3487–3492. IEEE, 2013.

Augugliaro, F., Lupashin, S., Hamer, M., Male, C., Hehn, M., Mueller, M. W., Willmann, J. S., Gramazio, F., Kohler, M., and D'Andrea, R. The flight assembled architecture installation: Cooperative construction with flying machines. *IEEE Control Systems Magazine*, 34(4): 46–64, 2014.

Bescuca, M. Factorised critics in deep multi-agent reinforcement learning. In *Master Thesis, University of Oxford*, 2019.

Bohmer, W., Kurin, V., and Whiteson, S. Deep coordination graphs. *arXiv preprint arXiv:1910.00091*, 2019.

Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

Caccavale, F. and Uchiyama, M. Cooperative Manipulators. In Siciliano, B. and Khatib, O. (eds.), *Springer Handbook of Robotics*, pp. 701–718. Springer Berlin Heidelberg, 2008.

Ciosek, K. and Whiteson, S. Expected policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

De Boer, P.-T., Kroese, D. P., Mannor, S., and Rubinstein, R. Y. A tutorial on the cross-entropy method. *Annals of operations research*, 134(1):19–67, 2005.

Foerster, J. N., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. Counterfactual multi-agent policy gradients. In *Thirty-second AAAI conference on artificial intelligence*, 2018.

Gu, S., Lillicrap, T., Sutskever, I., and Levine, S. Continuous deep q-learning with model-based acceleration. In *International Conference on Machine Learning*, pp. 2829–2838, 2016.

Ha, D., Dai, A., and Le, Q. V. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.

Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pp. 1856–1865, 2018.

Iqbal, S. and Sha, F. Actor-attention-critic for multi-agent reinforcement learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2961–2970, Long Beach, California, USA, 09–15 Jun 2019. PMLR.

Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.

Kilinc, O. and Montana, G. Multi-agent deep reinforcement learning with extremely noisy observations. *arXiv preprint arXiv:1812.00922*, 2018.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Kitano, H., Asada, M., Kuniyoshi, Y., Noda, I., Osawa, E., and Matsubara, H. Robocup: A challenge problem for ai. *AI magazine*, 18(1):73–73, 1997.

Koller, D. and Parr, R. Computing factored value functions for policies in structured mdps. In *Proceedings of IJCAI*, pp. 1332–1339, 1999.

Kraemer, L. and Banerjee, B. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.

Kuhn, H. Extensive games and the problem of information. *Annals of Mathematics Studies*, 28, 1953.

Leibo, J. Z., Zambaldi, V., Lanctot, M., Marecki, J., and Graepel, T. Multi-agent Reinforcement Learning in Sequential Social Dilemmas. *arXiv preprint arXiv:1702.03037*, 2017.

Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

Lin, L.-J. Reinforcement learning for robots using neural networks. In *Dissertation, Carnegie Mellon University*, 1992a.

Lin, L.-J. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992b.

Liu, S., Lever, G., Merel, J., Tunyasuvunakool, S., Heess, N., and Graepel, T. Emergent coordination through competition. *arXiv preprint arXiv:1902.07151*, 2019.

Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, O. P., and Mordatch, I. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in neural information processing systems*, pp. 6379–6390, 2017.

Mannor, S., Rubinstein, R. Y., and Gat, Y. The cross entropy method for fast policy search. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 512–519, 2003.

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.

Oliehoek, F. A., Spaan, M. T. J., and Nikos Vlassis. Optimal and Approximate Q-value Functions for Decentralized POMDPs. *JAIR*, 32:289–353, 2008.

Oliehoek, F. A., Amato, C., et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.

Ong, S. C., Png, S. W., Hsu, D., and Lee, W. S. Pomdps for robotic tasks with mixed observability. In *Robotics: Science and systems*, volume 5, pp. 4, 2009.

Plappert, M., Houthooft, R., Dhariwal, P., Sidor, S., Chen, R. Y., Chen, X., Asfour, T., Abbeel, P., and Andrychowicz, M. Parameter space noise for exploration. *arXiv preprint arXiv:1706.01905*, 2018.

Rashid, T., Samvelyan, M., Witt, C. S., Farquhar, G., Foerster, J., and Whiteson, S. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 4292–4301, 2018.

Riedmiller, M., Gabel, T., Hafner, R., and Lange, S. Reinforcement learning for robot soccer. *Autonomous Robots*, 27(1):55–73, 2009.

Samvelyan, M., Rashid, T., de Witt, C. S., Farquhar, G., Nardelli, N., Rudner, T. G. J., Hung, C.-M., Torr, P. H. S., Foerster, J., and Whiteson, S. The StarCraft Multi-Agent Challenge. *CoRR*, abs/1902.04043, 2019.

Schroeder de Witt, C., Foerster, J., Farquhar, G., Torr, P., Boehmer, W., and Whiteson, S. Multi-agent common knowledge reinforcement learning. In *Advances in Neural Information Processing Systems 32*, pp. 9924–9935. Curran Associates, Inc., 2019.

Shamshiri, R., Weltzien, C., Hameed, I., Yule, I., Grift, T., Balasundram, S., Pitonakova, L., Ahmad, D., and Chowdhary, G. Research and development in agricultural robotics: A perspective of digital farming. *International Journal of Agricultural and Biological Engineering*, 11: 1–14, 2018.

Son, K., Kim, D., Kang, W. J., Hostallero, D. E., and Yi, Y. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 5887–5896, 2019.

Stone, P. and Sutton, R. S. Scaling reinforcement learning toward RoboCup soccer. In *Icml*, volume 1, pp. 537–544. Citeseer, 2001.

Sunehag, P., Lever, G., Gruslys, A., Czarnecki, W. M., Zambaldi, V., Jaderberg, M., Lanctot, M., Sonnerat, N., Leibo, J. Z., Tuyls, K., and et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2085–2087. International Foundation for Autonomous Agents and Multiagent Systems, 2018.

Takadama, K., Matsumoto, S., Nakasuka, S., and Shimo-hara, K. A reinforcement learning approach to fail-safe design for multiple space robots?cooperation mechanism without communication and negotiation schemes. *Advanced Robotics*, 17(1):21–39, 2003.

Tampuu, A., Matiisen, T., Kodelja, D., Kuzovkin, I., Korjus, K., Aru, J., Aru, J., and Vicente, R. Multiagent cooperation and competition with deep reinforcement learning. *PloS one*, 12(4), 2017.

Tan, M. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.

Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012.

Van de Wiele, T., Warde-Farley, D., Mnih, A., and Mnih, V. Q-learning in enormous action spaces via amortized approximate maximization. *arXiv preprint arXiv:2001.08116*, 2020.

Vinyals, O., Ewalds, T., Bartunov, S., Georgiev, P., Vezhnevets, A. S., Yeo, M., Makhzani, A., Kttler, H., Agapiou, J., Schrittwieser, J., Quan, J., Gaffney, S., Petersen, S., Simonyan, K., Schaul, T., van Hasselt, H., Silver, D., Lillicrap, T., Calderone, K., Keet, P., Brunasso, A., Lawrence, D., Ekermo, A., Repp, J., and Tsing, R. StarCraft II: A New Challenge for Reinforcement Learning. *arXiv preprint arXiv:1708.04782*, 2017.

Wang, R. E., Everett, M. D., and How, J. P. R-MADDPG for partially observable environments and limited communication. In *Proceedings of the Reinforcement Learning for Real Life workshop (at ICML)*, 2019.

Wang, T., Liao, R., Ba, J., and Fidler, S. NerveNet: Learning structured policy with graph neural networks. In *6th International Conference on Learning Representations, ICLR*, 2018.

# 9. Experimental Details

In all experiments, we use a replay buffer of size $10^6$, the target networks are updated via soft target updates with $\tau = 0.001$, and we scale the gradient norms during training to be at most 0.5. The mixing network used in COMIX, COMIX-NAF, and FacMADDPG consists of a single hidden layer of 64 units, utilising an ELU non-linearity. The hypernetworks are then sized to produce weights of appropriate size. The hypernetworks producing the first layer weights and final layer weights and bias of the mixing network all consist of a single hidden layer of 64 units with a ReLU non-linearity. For COMIX and its ablations, to speed up the learning, we share the parameters of the agent networks across all agents. Similarly, in FacMADDPG, a single actor and critic network is shared among all agents, while in MADDPG, there is a separate actor and critic network for each agent as in the original algorithm.

## 9.1. Continuous Predator-Prey

The architecture of all agent networks is a MLP with 2 hidden layers of 64 units and ReLU non-linearities, except for COVDN-NAF and COMIX-NAF where we replace ReLU units with tanh units as it leads to better performance. In both MADDPG and FacMADDPG, the architecture of all agent networks and critic networks is also a MLP with 2 hidden layers of 64 units and ReLU non-linearities, while the final output layer of the actor was a tanh layer, to bound the actions. The global state consists of the observations of all agents. COMIX and MADDPG can both take advantage of the extra state information available during training, while COVDN ignores it. During training and testing, we restrict each episode to have a length of 25 time steps. Training lasts for 2 million timesteps. To encourage exploration, we use uncorrelated, mean-zero Gaussian noise during training (for all 2 million timesteps). We set $\gamma = 0.85$ for all experiments. We train on a batch size of 1024 after every timestep. All neural networks are trained using Adam (Kingma & Ba, 2014) optimiser with a learning rate of 0.01. To evaluate the learning performance, the training is paused after every 2000 timesteps during which 10 test episodes are run with agents performing action selection greedily in a decentralised fashion.

## 9.2. Multi-Agent Mujoco

The architecture of all agent networks is a MLP with 2 hidden layers with 400 and 300 units respectively, similar to the setting used in OpenAI Spinning Up.[4] All neural networks use ReLU non-linearities for all hidden layers, except for COVDN-NAF and COMIX-NAF where we found tanh units lead to better performance. In both MADDPG and FacMADDPG, the architecture of all agent networks and critic networks is also a MLP with 2 hidden layers with 400 and 300 units respectively, while the final output layer of the actor was a tanh layer, to bound the actions. The global state consists of the full state information returned by the original OpenAI Gym (Brockman et al., 2016). COMIX and MADDPG can both take advantage of the extra state information available during training, while COVDN ignores it. During training and testing, we restrict each episode to have a length of 1000 time steps. Training lasts for 4 million timesteps. To encourage exploration, we use uncorrelated, mean-zero Gaussian noise during training (for all 4 million timesteps). We also use the same trick as in OpenAI Spinning Up to improve exploration at the start of training. For a fixed number of steps at the beginning (we set it to be 10000), the agent takes actions which are sampled from a uniform random distribution over valid actions. After that, it returns to normal Gaussian exploration. We set $\gamma = 0.99$ for all experiments. We train on a batch size of 100 after every timestep. All neural networks are trained using Adam optimiser with a learning rate of 0.001. To evaluate the learning performance, the training is paused after every 4000 timesteps during which 10 test episodes are run with agents performing action selection greedily in a decentralised fashion.

## 9.3. Exploration

The choice of exploration strategy plays a substantial role in the performance of deep deterministic policy gradient algorithms (Ciosek & Whiteson, 2018). To keep exploration strategies comparable across MADDPG and $Q$-learning based COMIX, we restrict ourselves to noising in action spaces rather than parameter space (Plappert et al., 2018). MADDPG's official codebase [5] uses additive Gaussian noise with a standard deviation that is itself given by an additional policy output that is learnt end-to-end within the policy gradient loss. As $Q$-learning does not allow explicitly predict a policy output, we cannot apply a comparable strategy for COMIX. However, we find empirically that constant i.i.d. noising in the action spaces exhibits similar performance at lower variance than learnt noise on 2-Agent HalfCheetah (see Figure 6: Right). Even on Continuous Predator-Prey, a significantly less complex environment on which MADDPG's official codebase was tuned on,

---

[4]https://spinningup.openai.com/en/latest/.
[5]https://github.com/openai/maddpg.git

learnt exploration does not result in better limit performance than i.i.d. Gaussian noise (see Figure 6: Left).
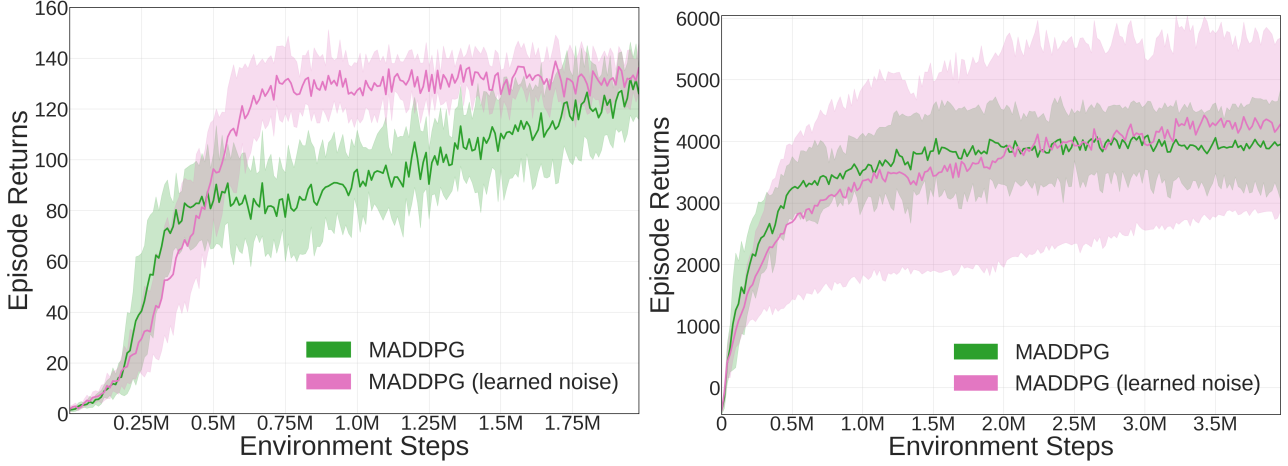


*Figure 6.* Mean episode return on **Left:** Continuous Predator-Prey and **Right:** 2-Agent HalfCheetah comparing MADDPG with constant i.i.d. Gaussian noise and MADDPG with learned Gaussian noise. The mean across 10 seeds is plotted and the 95% confidence interval is shown shaded.

## 10. Critic Mixing Network Constraints in FacMADDPG

As in an actor-critic setting, the critic is not used for greedy action selection, FacMADDPG does not strictly require a monotonocity constraint on its critic mixing network. However, we find empirically that introducing the monotonicity requirement significantly increases performance (see Figure 7). This supports the hypothesis that introducing monotonicity constraints strikes a reasonable trade-off between having independent critics with limited coordinative ability and the case where excess coordinative expressivity in the unconstrained critic leads to an increase in learning difficulty.
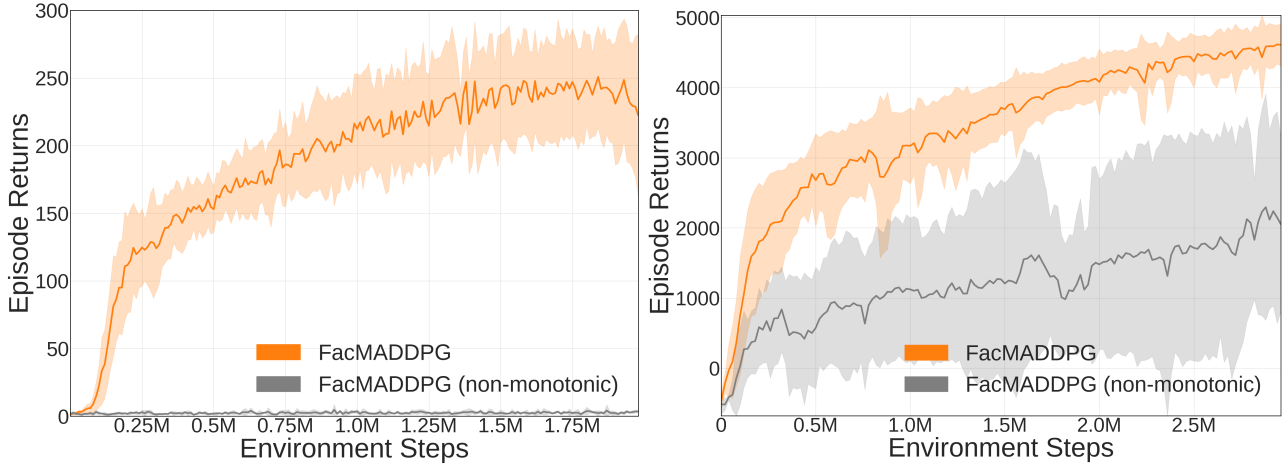


*Figure 7.* Mean episode return on **Left:** Continuous Predator-Prey and **Right:** 2-Agent HalfCheetah comparing FacMADDPG and FacMADDPG without monotonicity constraints on the mixing network of the critic. The mean across 10 seeds is plotted and the 95% confidence interval is shown shaded.

| Task | Goal | Special observations | Reward function |
|------|------|----------------------|-----------------|
| 2-Agent Swimmer | Maximise +ve $x$-speed. | - | $\frac{\Delta x}{\Delta t} + 0.0001\alpha$ |
| 2-Agent Reacher | Fingertip (green) needs to reach target at random location (red). | Target is only visible to green agent. | $-\|\text{distance from fingertip to target}\|_2^2$ $+\alpha$ |
| 2-Agent Ant | Maximise +ve $x$-speed. | All agents can observe the central torso. | $\frac{\Delta x}{\Delta t} + 5 \cdot 10^{-4}\|\text{external contact forces}\|_2^2$ $+0.5\alpha + 1$ |
| 2-Agent Ant Diag | Maximise +ve $x$-speed. | All agents can observe the central torso. | $\frac{\Delta x}{\Delta t} + 5 \cdot 10^{-4}\|\text{external contact forces}\|_2^2$ $+0.5\alpha + 1$ |
| 2-Agent HalfCheetah | Maximise +ve $x$-speed. | - | $\frac{\Delta x}{\Delta t} + 0.1\alpha$ |
| 2-Agent Humanoid | Maximise +ve $x$-speed. | - | $0.25\frac{\Delta x}{\Delta t} + \min(10,$ $5 \cdot 10^{-6}\|\text{external contact forces}\|_2^2)$ |
| 2-Agent HumanoidStandup | Maximise +ve $x$-speed. | - | $\frac{y}{\Delta t} + \min(10,$ $5 \cdot 10^{-6}\|\text{external contact forces}\|_2^2)$ |
| 3-Agent Hopper | Maximise +ve $x$-speed. | - | $\frac{\Delta x}{\Delta t} + 0.001\alpha + 1.0$ |
| 4-Agent Ant | Maximise +ve $x$-speed. | All agents can observe the central torso. | $\frac{\Delta x}{\Delta t} + 5 \cdot 10^{-4}\|\text{external contact forces}\|_2^2$ $+0.5\alpha + 1$ |
| 6-Agent HalfCheetah | Maximise +ve $x$-speed. | - | $0.25\frac{\Delta x}{\Delta t} + \min(10,$ $5 \cdot 10^{-6}\|\text{external contact forces}\|_2^2)$ |

*Table 1.* Overview of tasks contained in Multi-Agent Mujoco. We define $\alpha$ as an action regularisation term $-\|\mathbf{u}\|_2^2$.

---

**Algorithm 2** For each agent $a$, we perform $n_c$ CEM iterations. Hyper-parameters $d_i \in \mathbb{N}$ control how many actions are sampled at the $i$th iteration.

---

**function** CEM $(Q_1, \ldots, Q_N, \tau_1, \ldots, \tau_N)$
  **for** $a := 1,\ a \le N$ **do**
    $\boldsymbol{\mu}_a \leftarrow \mathbf{0} \in \mathbb{R}^{|\mathcal{A}_a|}$
    $\boldsymbol{\sigma}_a \leftarrow \mathbf{1} \in \mathbb{R}^{|\mathcal{A}_a|}$
    **for** $i := 1,\ i \le n_c$ **do**
      **for** $j := 1,\ j \le d_i$ **do**
        $\mathbf{v}'_{aj} \sim \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\sigma}_a)$
        $\mathbf{v}_{aj} \leftarrow \tanh(\mathbf{v}'_{aj})$
        $q_{aj} \leftarrow Q_a(\tau_a, \mathbf{v}_{aj})$
        $j \leftarrow j + 1$
      **end for**
      **if** $i < n_c$ **then**
        $U \leftarrow \{\mathbf{v}'_{al} \mid q_{al} \in \text{top}k_i(q_{a1}, \ldots, q_{ad_i}), \forall l \in \{1 \ldots N\}\}$
        $\boldsymbol{\mu}_a \leftarrow \text{sample\_mean}(U)$
        $\boldsymbol{\sigma}_a \leftarrow \text{sample\_std}(U)$
      **else** {Right}
        $m \leftarrow \arg\max_j q_{aj}$
        $\mathbf{u}_a \leftarrow \mathbf{v}_{am}$
      **end if**
      $i \leftarrow i + 1$
    **end for**
    $a \leftarrow a + 1$
  **end for**
  **return** $\langle \mathbf{u}_1, \ldots, \mathbf{u}_n \rangle$
**end function**