

---

# FACMAC: Factored Multi-Agent Centralised Policy Gradients

---

**Bei Peng**\*<sup>†</sup>  
University of Liverpool

**Tabish Rashid**\*  
University of Oxford

**Christian A. Schroeder de Witt**\*  
University of Oxford

**Pierre-Alexandre Kamienny**<sup>†</sup>  
Facebook AI Research

**Philip H. S. Torr**  
University of Oxford

**Wendelin Böhmer**  
Delft University of Technology

**Shimon Whiteson**  
University of Oxford

## Abstract

We propose **FAC**tored **MU**lti-**AG**ent **C**entralised policy gradients (FACMAC), a new method for cooperative multi-agent reinforcement learning in both discrete and continuous action spaces. Like MADDPG, a popular multi-agent actor-critic method, our approach uses deep deterministic policy gradients to learn policies. However, FACMAC learns a centralised but factored critic, which combines per-agent utilities into the joint action-value function via a non-linear monotonic function, as in QMIX, a popular multi-agent  $Q$ -learning algorithm. However, unlike QMIX, there are no inherent constraints on factoring the critic. We thus also employ a nonmonotonic factorisation and empirically demonstrate that its increased representational capacity allows it to solve some tasks that cannot be solved with monolithic, or monotonically factored critics. In addition, FACMAC uses a centralised policy gradient estimator that optimises over the entire joint action space, rather than optimising over each agent’s action space separately as in MADDPG. This allows for more coordinated policy changes and fully reaps the benefits of a centralised critic. We evaluate FACMAC on variants of the multi-agent particle environments, a novel multi-agent MuJoCo benchmark, and a challenging set of StarCraft II micromanagement tasks. Empirical results demonstrate FACMAC’s superior performance over MADDPG and other baselines on all three domains.

## 1 Introduction

Significant progress has been made in cooperative multi-agent reinforcement learning (MARL) under the paradigm of *centralised training with decentralised execution* (CTDE) [27, 18] in recent years, both in value-based [41, 31, 37, 32, 44, 30] and actor-critic [23, 9, 36, 13, 8] approaches. Most popular multi-agent actor-critic methods such as COMA [9] and MADDPG [23] learn a *centralised critic* with decentralised actors. The critic is centralised to make use of all available information (i.e., it can condition on the global state and the joint action) to estimate the joint action-value function  $Q_{tot}$ , unlike a *decentralised critic* that estimates the local action-value function  $Q_a$  based only on individual observations and actions for each agent  $a$ .<sup>3</sup> Even though the joint action-value function

---

\*Equal contribution. Correspondence to: Bei Peng <bei.peng@liverpool.ac.uk>

<sup>†</sup>Work done while the authors were at the University of Oxford.

<sup>3</sup>COMA learns a single centralised critic for all cooperative agents due to parameter sharing. For each agent the critic has different inputs and can thus output different values for the same state and joint action.

these actor-critic methods can represent is not restricted, in practice they significantly underperform value-based methods like QMIX [31] on the challenging StarCraft Multi-Agent Challenge (SMAC) [35] benchmark [33, 32].

In this paper, we propose a novel approach called *FACTored Multi-Agent Centralised policy gradients* (FACMAC), which works for both discrete and continuous cooperative multi-agent tasks. Like MADDPG, our approach uses deep deterministic policy gradients [20] to learn decentralised policies. However, FACMAC learns a single *centralised but factored critic*, which factors the joint action-value function  $Q_{tot}$  into per-agent utilities  $Q_a$  that are combined via a non-linear monotonic function, as in the popular  $Q$ -learning algorithm QMIX [31]. While the critic used in COMA and MADDPG is also centralised, it is monolithic rather than factored.<sup>4</sup> Compared to learning a monolithic critic, our factored critic can potentially scale better to tasks with a larger number of agents and/or actions. In addition, in contrast to other value-based approaches such as QMIX, there are no inherent constraints on factoring the critic. This allows us to employ rich value factorisations, including *nonmonotonic* ones, that value-based methods cannot directly use without forfeiting decentralisability or introducing other significant algorithmic changes. We thus also employ a nonmonotonic factorisation and empirically demonstrate that its increased representational capacity allows it to solve some tasks that cannot be solved with monolithic, or monotonically factored critics.

In MADDPG, a separate policy gradient is derived for each agent individually, which optimises its policy assuming all other agents’ actions are fixed. This could cause the agents to converge to sub-optimal policies in which no single agent wishes to change its action unilaterally. In FACMAC, we use a new *centralised* gradient estimator that optimises over the entire joint action space, rather than optimising over each agent’s action space separately as in MADDPG. The agents’ policies are thus trained as a single joint-action policy, which can enable learning of more coordinated behaviour, as well as the ability to escape sub-optimal solutions. The centralised gradient estimator fully reaps the benefits of learning a centralised critic, by not implicitly marginalising over the actions of the other agents in the policy-gradient update. The gradient estimator used in MADDPG is also known to be vulnerable to relative overgeneralisation [47]. To overcome this issue, in our centralised gradient estimator, we sample all actions from all agents’ current policies when evaluating the joint action-value function. We empirically show that MADDPG can quickly get stuck in local optima in a simple continuous matrix game, whereas our centralised gradient estimator finds the optimal policy. While Lyu et al. [24] recently show that merely using a centralised critic (with per-agent gradients that optimise over each agent’s actions separately) does not necessarily lead to better coordination between agents, our centralised gradient estimator re-establishes the value of using centralised critics.

Most recent works on continuous MARL focus on evaluating their algorithms on the multi-agent particle environments [23], which feature a simple two-dimensional world with some basic simulated physics. To demonstrate FACMAC’s scalability to more complex continuous domains and to stimulate more progress in continuous MARL, we introduce *Multi-Agent MuJoCo* (MAMuJoCo), a new, comprehensive benchmark suite that allows the study of decentralised continuous control. Based on the popular single-agent MuJoCo benchmark [5], MAMuJoCo features a wide variety of novel robotic control tasks in which multiple agents within a single robot have to solve a task cooperatively.

We evaluate FACMAC on variants of the multi-agent particle environments [23] and our novel MAMuJoCo benchmark, which both feature continuous action spaces, and the challenging SMAC benchmark [35], which features discrete action spaces. Empirical results demonstrate FACMAC’s superior performance over MADDPG and other baselines on all three domains. In particular, FACMAC scales better when the number of agents (and/or actions) and the complexity of the task increases. Results on SMAC show that FACMAC significantly outperforms stochastic DOP [46], which recently claimed to be the first multi-agent actor-critic method to outperform state-of-the-art value-based methods on SMAC, in all scenarios we tested. Moreover, our ablations and additional experiments demonstrate the advantages of both factoring the critic and using our centralised gradient estimator. We show that, compared to learning a monolithic critic, learning a factored critic can: 1) better take advantage of the centralised gradient estimator to optimise the agent policies when the

---

In MADDPG, each agent learns its own centralised critic, as it is designed for general multi-agent learning problems, including cooperative, competitive, and mixed settings.

<sup>4</sup>We use “centralised and monolithic critic” and “monolithic critic” interchangeably to refer to the centralised critic used in COMA and MADDPG, and “centralised but factored critic” and “factored critic” interchangeably to refer to the critic used in our approach.

number of agents and/or actions is large, and 2) leverage a nonmonotonic factorisation to solve tasks that cannot be solved with monolithic or monotonically factored critics.

## 2 Background

We consider a *fully cooperative multi-agent task* in which a team of agents interacts with the same environment to achieve some common goal. It can be modeled as a *decentralised partially observable Markov decision process* (Dec-POMDP) [28] consisting of a tuple  $G = \langle \mathcal{N}, S, U, P, r, \Omega, O, \gamma \rangle$ . Here  $\mathcal{N} \equiv \{1, \dots, n\}$  denotes the finite set of agents and  $s \in S$  describes the true state of the environment. At each time step, each agent  $a \in \mathcal{N}$  selects a discrete or continuous action  $u_a \in U$ , forming a joint action  $\mathbf{u} \in \mathbf{U} \equiv U^n$ . This results in a transition to the next state  $s'$  according to the state transition function  $P(s'|s, \mathbf{u}) : S \times \mathbf{U} \times S \rightarrow [0, 1]$  and a team reward  $r(s, \mathbf{u})$ .  $\gamma \in [0, 1)$  is a discount factor. Due to the *partial observability*, each agent  $a \in \mathcal{N}$  draws an individual partial observation  $o_a \in \Omega$  from the observation kernel  $O(s, a)$ . Each agent learns a stochastic policy  $\pi_a(u_a|\tau_a)$  or a deterministic policy  $\mu_a(\tau_a)$ , conditioned only on its local action-observation history  $\tau_a \in T \equiv (\Omega \times U)^*$ . The joint stochastic policy  $\pi$  induces a joint *action-value function*:  $Q^\pi(s_t, \mathbf{u}_t) = \mathbb{E}_{s_{t+1:\infty}, \mathbf{u}_{t+1:\infty}} [R_t|s_t, \mathbf{u}_t]$ , where  $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$  is the discounted return. Similarly, the joint deterministic policy  $\mu$  induces a joint action-value function denoted  $Q^\mu(s_t, \mathbf{u}_t)$ . We adopt the *centralised training with decentralised execution* (CTDE) paradigm [27, 18], where policy training can exploit extra global information that might be available and has the freedom to share information between agents during training. However, during execution, each agent must act with only access to its own action-observation history.

**VDN and QMIX.** VDN [41] and QMIX [31] are  $Q$ -learning algorithms for cooperative MARL tasks with discrete actions. They both aim to efficiently learn a centralised but factored action-value function  $Q_{tot}^\pi$ , using CTDE. To ensure consistency between the centralised and decentralised policies, VDN and QMIX factor  $Q_{tot}^\pi$  assuming additivity and monotonicity, respectively. More specifically, VDN factors  $Q_{tot}^\pi$  into a sum of the per-agent utilities:  $Q_{tot}^\pi(\boldsymbol{\tau}, \mathbf{u}; \phi) = \sum_{a=1}^n Q_a^{\pi_a}(\tau_a, u_a; \phi_a)$ . QMIX, however, represents  $Q_{tot}^\pi$  as a continuous monotonic mixing function of each agent’s utilities:  $Q_{tot}^\pi(\boldsymbol{\tau}, \mathbf{u}, s; \phi, \psi) = f_\psi(s, Q_1^{\pi_1}(\tau_1, u_1; \phi_1), \dots, Q_n^{\pi_n}(\tau_n, u_n; \phi_n))$ , where  $\frac{\partial f_\psi}{\partial Q_a} \geq 0, \forall a \in \mathcal{N}$ . This is sufficient to ensure that the global arg max performed on  $Q_{tot}^\pi$  yields the same result as a set of individual arg max performed on each  $Q_a^{\pi_a}$ . Here  $f_\psi$  is approximated by a monotonic mixing network, parameterised by  $\psi$ . Monotonicity can be guaranteed by non-negative mixing weights. These weights are generated by separate *hypernetworks* [12], which condition on the full state  $s$ . QMIX is trained end-to-end to minimise the following loss:

$$\mathcal{L}(\phi, \psi) = \mathbb{E}_{\mathcal{D}} \left[ \left( y^{tot} - Q_{tot}^\pi(\boldsymbol{\tau}, \mathbf{u}, s; \phi, \psi) \right)^2 \right], \quad (1)$$

where the bootstrapping target  $y^{tot} = r + \gamma \max_{\mathbf{u}'} Q_{tot}^\pi(\boldsymbol{\tau}', \mathbf{u}', s'; \phi^-, \psi^-)$ . Here  $r$  is the global reward, and  $\phi^-$  and  $\psi^-$  are parameters of the target  $Q$  and mixing network, respectively, as in DQN [26]. The expectation is estimated with a minibatch of transitions sampled from an experience replay buffer  $\mathcal{D}$  [21]. During execution, each agent selects actions greedily with respect to its own  $Q_a^{\pi_a}$ .

**MADDPG.** MADDPG [23] is an extension of DDPG [20] to multi-agent settings. It is an actor-critic, off-policy method that uses the paradigm of CTDE to learn deterministic policies in continuous action spaces. In MADDPG, a separate actor and critic is learned for each agent, such that each agent can have its own arbitrary reward function. It is therefore applicable to either cooperative, competitive, or mixed settings. We assume each agent  $a$  has a deterministic policy  $\mu_a(\tau_a; \theta_a)$ , parameterised by  $\theta_a$  (abbreviated as  $\mu_a$ ), and let  $\boldsymbol{\mu} = \{\mu_a(\tau_a; \theta_a)\}_{a=1}^n$  be the set of all agent policies. In MADDPG, a *centralised and monolithic critic* that estimates the joint action-value function  $Q_a^\mu(s, u_1, \dots, u_n; \phi_a)$  is learned for each agent  $a$  separately. The critic is said to be *centralised* as it utilises information only available to it during the *centralised* training phase, the global state  $s^5$  and the actions of all agents,  $u_1, \dots, u_n$ , to estimate the joint action-value function  $Q_a^\mu$ , which is parameterised by  $\phi_a$ . This joint action-value function is trained by minimising the following loss:

$$\mathcal{L}(\phi_a) = \mathbb{E}_{\mathcal{D}} \left[ \left( y^a - Q_a^\mu(s, u_1, \dots, u_n; \phi_a) \right)^2 \right], \quad (2)$$

<sup>5</sup>If the global state  $s$  is not available, the centralised and monolithic critic can condition on the joint observations or action-observation histories.

where  $y^a = r_a + \gamma Q_a^\mu(s', u'_1, \dots, u'_n |_{u'_a = \mu_a(\tau'_a; \theta_a^-)}; \phi_a^-)$ . Here  $r_a$  is the reward received by each agent  $a$ ,  $u'_1, \dots, u'_n$  is the set of target policies with delayed parameters  $\theta_a^-$ , and  $\phi_a^-$  are the parameters of the target critic. The replay buffer  $\mathcal{D}$  contains the transition tuples  $(s, s', u_1, \dots, u_n, r_1, \dots, r_n)$ .

The following policy gradient can be calculated individually to update the policy of each agent  $a$ :

$$\nabla_{\theta_a} J(\mu_a) = \mathbb{E}_{\mathcal{D}} \left[ \nabla_{\theta_a} \mu_a(\tau_a) \nabla_{u_a} Q_a^\mu(s, u_1, \dots, u_n) \Big|_{u_a = \mu_a(\tau_a)} \right], \quad (3)$$

where the current agent  $a$ 's action  $u_a$  is sampled from its current policy  $\mu_a$  when evaluating the joint action-value function  $Q_a^\mu$ , while all other agents' actions are sampled from the replay buffer  $\mathcal{D}$ .

### 3 FACMAC

In this section, we propose a new approach called *FACTored Multi-Agent Centralised policy gradients* (FACMAC) that uses a centralised but factored critic and a centralised gradient estimator to learn continuous cooperative tasks. We start by describing the idea of learning a centralised but factored critic. We then discuss our new centralised gradient estimator and demonstrate its benefit in a simple continuous matrix game. Finally, we discuss how we adapt our method to discrete cooperative tasks.

#### 3.1 Learning a Centralised but Factored Critic

Learning a centralised and monolithic critic conditioning on the global state and the joint action can be difficult and/or impractical when the number of agents and/or actions is large [13]. We thus employ value function factorisation in the multi-agent actor-critic framework to enable scalable learning of a centralised critic in Dec-POMDPs. Another key advantage of adopting value factorisation in an actor-critic framework is that, compared to value-based methods, it allows for a more flexible factorisation as the critic's design is not constrained. One can employ any type of factorisation, including nonmonotonic factorisations that value-based methods cannot directly use without forfeiting decentralisability or introducing other significant algorithmic changes.

Specifically, in FACMAC, all agents share a centralised critic  $Q_{tot}^\mu$  that is factored as:

$$Q_{tot}^\mu(\boldsymbol{\tau}, \mathbf{u}, s; \boldsymbol{\phi}, \psi) = g_\psi(s, \{Q_a^{\mu_a}(\tau_a, u_a; \phi_a)\}_{a=1}^n), \quad (4)$$

where  $\boldsymbol{\phi}$  and  $\phi_a$  are parameters of the joint action-value function  $Q_{tot}^\mu$  and agent-wise utilities  $Q_a^{\mu_a}$ , respectively. In our canonical implementation which we refer to as FACMAC,  $g_\psi$  is a non-linear monotonic function parametrised as a mixing network with parameters  $\psi$ , as in QMIX [31]. To evaluate the policy, the centralised but factored critic is trained by minimising the following loss:

$$\mathcal{L}(\boldsymbol{\phi}, \psi) = \mathbb{E}_{\mathcal{D}} \left[ \left( y^{tot} - Q_{tot}^\mu(\boldsymbol{\tau}, \mathbf{u}, s; \boldsymbol{\phi}, \psi) \right)^2 \right], \quad (5)$$

where  $y^{tot} = r + \gamma Q_{tot}^\mu(\boldsymbol{\tau}', \boldsymbol{\mu}(\boldsymbol{\tau}'; \boldsymbol{\theta}^-), s'; \boldsymbol{\phi}^-, \psi^-)$ . Here  $\mathcal{D}$  is the replay buffer, and  $\boldsymbol{\theta}^-$ ,  $\boldsymbol{\phi}^-$ , and  $\psi^-$  are parameters of the target actors, critic, and mixing network, respectively.

Leveraging the flexibility of our approach, namely the lack of restrictions on the form of the critic, we also explore a new nonmonotonic factorisation with full representational capacity. The joint action-value function  $Q_{tot}^\mu$  is represented as a non-linear nonmonotonic mixing function of per-agent utilities  $Q_a^{\mu_a}$ . This nonmonotonic mixing function is parameterised as a mixing network, with a similar architecture to  $g_\psi$  in FACMAC, but without the constraint of monotonicity enforced by using non-negative weights. We refer to this method as FACMAC-nonmonotonic. Additionally, to better understand the advantages of factoring a centralised critic, we also explore two additional simpler factorisation schemes. These include factoring the centralised critic  $Q_{tot}^\mu$  into a sum of per-agent utilities  $Q_a^{\mu_a}$  as in VDN (FACMAC-vdn), and as a sum of  $Q_a^{\mu_a}$  and a state-dependent bias (FACMAC-vdn-s). Our value factorisation technique is general and can be readily applied to any multi-agent actor-critic algorithms that learn centralised and monolithic critics [23, 9, 8].

#### 3.2 Centralised Policy Gradients

To update the decentralised policy of each agent, a naive adaptation of the deterministic policy gradient used by MADDPG (shown in (3)) is

$$\nabla_{\theta_a} J(\mu_a) = \mathbb{E}_{\mathcal{D}} \left[ \nabla_{\theta_a} \mu_a(\tau_a) \nabla_{u_a} Q_{tot}^\mu(\boldsymbol{\tau}, u_1, \dots, u_n, s) \Big|_{u_a = \mu_a(\tau_a)} \right]. \quad (6)$$

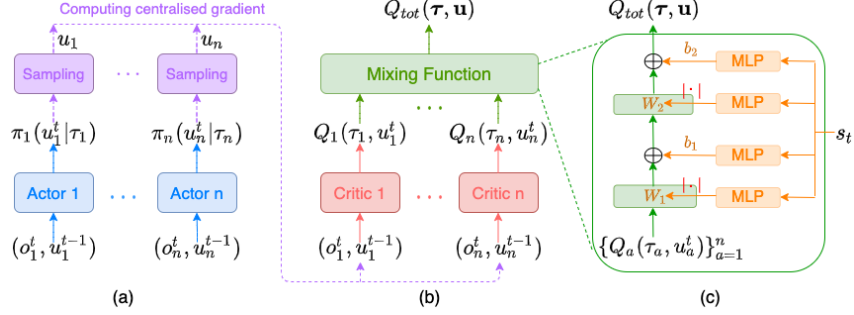


Figure 1: The overall FACMAC architecture. (a) The decentralised policy networks. There is a sampling step since we sample from the categorical distribution when using discrete actions. (b) The centralised but factored critic. (c) The non-linear monotonic mixing function.

Compared to the policy gradient used in MADDPG, the updates of all agents’ individual deterministic policies now depend on the single shared factored critic  $Q_{tot}^\mu$ , as opposed to learning and utilising a monolithic critic  $Q_a^\mu$  for each agent. However, there are two main problems in both policy gradients. First, each agent optimises its own policy assuming all other agents’ actions are fixed, which could cause the agents to converge to sub-optimal policies in which no single agent wishes to change its action unilaterally. Second, both policy gradients make the corresponding methods vulnerable to relative overgeneralisation [47] as, when agent  $a$  ascends the policy gradient based on  $Q_a^\mu$  or  $Q_{tot}^\mu$ , only its own action  $u_a$  is sampled from its current policy  $\mu_a$ , while all other agents’ actions are sampled from the replay buffer  $\mathcal{D}$ . The other agents’ actions thus might be drastically different from the actions their current policies would choose. This could cause the agents to converge to sub-optimal actions that appear to be a better choice when considering the effect of potentially arbitrary actions from the other collaborating agents.

In FACMAC, we use a new *centralised* gradient estimator that optimises over the entire joint action space, rather than optimising over each agent’s actions separately as in both (3) and (6), to achieve better coordination among agents. In addition, to overcome relative overgeneralisation, when calculating the policy gradient we sample all actions from all agents’ current policies when evaluating  $Q_{tot}^\mu$ . Our centralised policy gradient can thus be estimated as

$$\nabla_{\theta} J(\boldsymbol{\mu}) = \mathbb{E}_{\mathcal{D}} \left[ \nabla_{\theta} \boldsymbol{\mu} \nabla_{\boldsymbol{\mu}} Q_{tot}^{\mu}(\boldsymbol{\tau}, \mu_1(\tau_1), \dots, \mu_n(\tau_n), s) \right], \quad (7)$$

where  $\boldsymbol{\mu} = \{\mu_1(\tau_1; \theta_1), \dots, \mu_n(\tau_n; \theta_n)\}$  is the set of all agents’ current policies and all agents share the same actor network parameterised by  $\theta$ . However, it is not a requirement of our method for all agents to share parameters in this manner.

If the critic factorisation is linear, as in FACMAC-vdn, then the centralised gradient is equivalent to the per-agent gradients that optimise over each agent’s actions separately. This is explored in more detail by DOP [46], which restricts the factored critic to be linear to exploit this equivalence. A major benefit of our method then, is that it does not place any such restrictions on the critic. As remarked by Lyu et al. [24], merely using a centralised critic with per-agent gradients does not necessarily lead to better coordination between agents due to the two problems outlined above. Even with a factored critic, methods that use stochastic policy gradients can still suffer from the problems caused by the per-agent gradients if a fully factored stochastic policy is used. Our centralised gradient estimator, which uses deterministic policies and optimises over the entire joint action space, is required in order to fully take advantage of a centralised critic.

Figure 1 illustrates the overall FACMAC architecture. For each agent  $a$ , there is one policy network that decides which individual action (discrete or continuous) to take. There is also one critic network for each agent  $a$  that estimates the individual agent utilities  $Q_a$ , which are then combined into the joint action-value function  $Q_{tot}$  via a non-linear monotonic mixing function as in QMIX.  $Q_{tot}$  is then used by our centralised gradient estimator to help the actor update its policy parameters.

To show the benefits of our new centralised gradient estimator, we compare MADDPG with the centralised policy gradient (CPG) against the original MADDPG on a simple continuous cooperative matrix game with two agents. Figure 10 in Appendix D.1 illustrates this matrix game. There is a

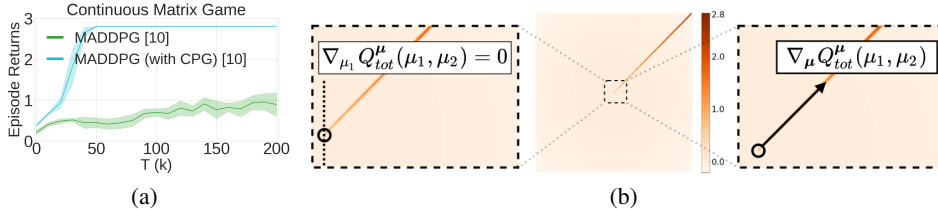


Figure 2: (a) Mean test return on Continuous Matrix Game. (b) **Left:** Per-agent policy gradient at the origin. For agent 1 (similarly for agent 2) it is 0 since the gradient term assumes the other agent’s action to be fixed and thus it only considers the relative improvements along the dotted line (agent 1’s own action space). **Right:** Our centralised policy gradient correctly determines the gradient for improving the joint action.

narrow path (shown in red) starting from the origin  $(0, 0)$  to  $(1, 1)$ , in which the reward gradually increases. Everywhere else there is a small punishment moving away from the origin, increasing in magnitude further from the origin. Experimental results are shown in Figure 2(a). MADDPG quickly gets stuck in the local optimum within  $200k$  timesteps, while MADDPG (with CPG) robustly converges to the optimal policy. Figure 2(b) visualises the differences between the per-agent and centralised policy gradients, demonstrating that the per-agent policy gradient can be wrong in our continuous matrix game and our centralised policy gradient is necessary to take advantage of the centralised critic. In Section 5, we further demonstrate the benefits of our centralised gradient estimator in more complex tasks.

### 3.3 Discrete Policy Learning

As FACMAC requires differentiable policies and the sampling process of discrete actions from a categorical distribution is not differentiable, we use the Gumbel-Softmax estimator [14] to enable efficient learning of FACMAC on cooperative tasks with discrete actions. The Gumbel-Softmax estimator is a continuous distribution that approximates discrete samples from a categorical distribution to produce differentiable samples. Moreover, we use the Straight-Through Gumbel-Softmax Estimator [14] to ensure the action dynamics during training and evaluation are the same. Specifically, during training, we sample discrete actions  $u_a$  from the original categorical distribution in the forward pass, but use the continuous Gumbel-Softmax sample  $x_a$  in the backward pass to approximate the gradients:  $\nabla_{\theta_a} u_a \approx \nabla_{\theta_a} x_a$ . We can then update the agent’s policy using our centralised policy gradient:  $\nabla_{\theta} J(\mathbf{x}) \approx \mathbb{E}_{\mathcal{D}} \left[ \nabla_{\theta} \mathbf{x} \nabla_{\mathbf{x}} Q_{tot}^{\mathbf{x}}(\tau, x_1, \dots, x_n, s) \right]$ , where  $\mathbf{x} = \{x_1, \dots, x_n\}$  is the set of continuous samples that approximates the discrete agent actions. The softmax temperature hyperparameter  $\tau$  is set to be 1 in our experiments.

## 4 Multi-Agent MuJoCo

The evaluation of continuous MARL algorithms has recently been largely limited to the simple multi-agent particle environments [23]. We believe the lack of diverse continuous benchmarks is one factor limiting progress in continuous MARL. To demonstrate FACMAC’s scalability to more complex continuous domains and to stimulate more progress in continuous MARL, we develop *Multi-Agent MuJoCo* (MAMuJoCo), a novel benchmark for continuous cooperative multi-agent robotic control. Starting from the popular fully observable single-agent robotic MuJoCo [42] control suite included with OpenAI Gym [5], we create a wide variety of novel scenarios in which multiple agents within a single robot have to solve a task cooperatively.

This design offers important benefits. It facilitates comparisons to existing literature on both the fully observable single-agent domain [29], as well as settings with low-bandwidth communication [45]. More importantly, it allows for the study of novel MARL algorithms for decentralised coordination in isolation (scenarios with multiple robots may add confounding factors such as spatial exploration), which is currently a gap in the research literature. MAMuJoCo also includes scenarios with a larger and more flexible number of agents, which takes inspiration from modular robotics [49, 19]. Compared to traditional robots, modular robots are more versatile, configurable, and scalable. We thus

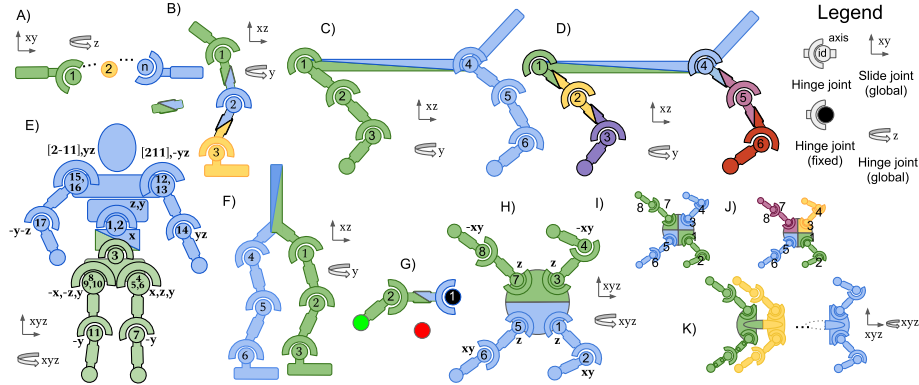


Figure 3: **Agent partitionings for MAMuJoCo environments:** A) Manyagent Swimmer, B) 3-Agent Hopper [3x1], C) 2-Agent HalfCheetah [2x3], D) 6-Agent HalfCheetah [6x1], E) 2-Agent Humanoid and 2-Agent HumanoidStandup (each [1x9,1x8]), F) 2-Agent Walker [2X3], G) 2-Agent Reacher [2x1], H) 2-Agent Ant [2x4], I) 2-Agent Ant Diag [2x4], J) 4-Agent Ant [4x2], and K) Manyagent Ant. 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 [5]. Global joints indicate degrees of freedom of the center of mass of the composite robotic agent.

develop two scenarios named ManyAgent Swimmer and ManyAgent Ant, in which one can configure an arbitrarily large number of agents (within the memory limits), each controlling a consecutive segment of arbitrary length.

Single-robot multi-agent tasks in MAMuJoCo arise 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 3. We then partition the body graph into disjoint sub-graphs, one for each agent, each of which contains one or more joints that can be controlled. Note that in ManyAgent Swimmer (see Figure 3A) and ManyAgent Ant (see Figure 3K), the number of agents are not limited by the given single robotic agent. See Appendix B for more details about MAMuJoCo.

## 5 Experimental Results

In this section we present our experimental results on our cooperative variants of the continuous *simple tag* environment introduced by Lowe et al. [23] (we refer to this environment as Continuous Predator-Prey), our novel continuous benchmark MAMuJoCo, and the challenging SMAC<sup>6</sup> [35] benchmark with discrete action spaces. In discrete cooperative tasks, we compare with state-of-the-art multi-agent actor-critic algorithms MADDPG [23], COMA [9], CentralV [9], DOP [46], VDAC-mix [40], and value-based methods QMIX [31] and QPLEX [44]. In continuous cooperative tasks, we compare with MADDPG [23] and independent DDPG (IDDPG), as well as COVDN and COMIX, two novel baselines described below. We also explore different forms of critic factorisation to better understand the advantages of factoring a centralised critic. More details about the environments, experimental setup, and training details are included in Appendix D and E. Code is available at <https://github.com/oxwhirl/facmac>.

**COVDN and COMIX.** We find that not many multi-agent value-based methods work off the shelf with continuous actions. To compare FACMAC against value-based approaches in continuous cooperative tasks, we use existing continuous *Q*-learning approaches in single-agent settings to extend VDN and QMIX to continuous action spaces. Specifically, we introduce COVDN and COMIX, which use VDN-style and QMIX-style factorisation respectively and both perform approximate greedy action selection using the cross-entropy method (CEM) [7]. CEM is a sampling-based derivative-free

<sup>6</sup>We utilise SC2.4.10., which is used by the latest PyMAREL [35] framework. The original results reported in Samvelyan et al. [35] and Rashid et al. [33] use SC2.4.6. Performance is **not** always comparable across versions.

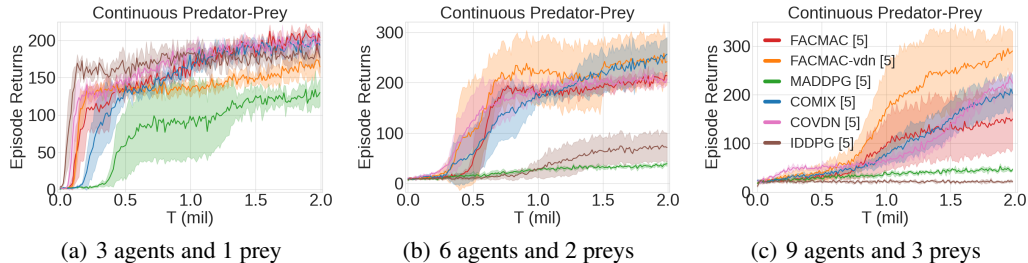


Figure 4: Mean episode return on Continuous Predator-Prey with different number of agents and preys. The mean across 5 seeds is plotted and the 95% confidence interval is shown shaded. The numbers in square brackets in the figure legend represent the number of random seeds used to run each method (similarly for all other figures).

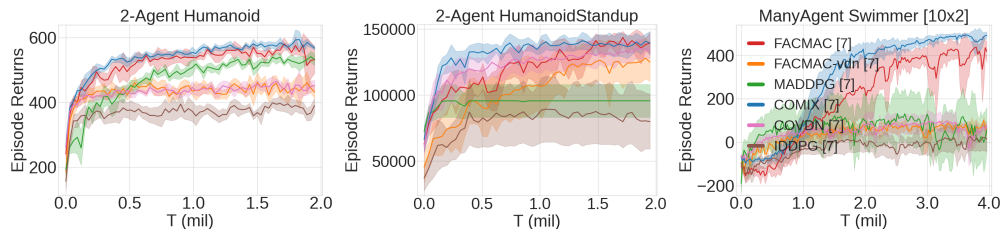


Figure 5: Mean episode return on different MAMuJoCo tasks. In ManyAgent Swimmer, we configure the number of agents to be 10, each controlling a consecutive segment of length 2. The mean across 7 seeds is plotted and the 95% confidence interval is shown shaded.

heuristic search method that has been successfully used to find approximate maxima of nonconvex  $Q$ -networks in single-agent robotic control tasks [15]. The centralised but factored  $Q_{tot}$  allows us to use CEM to sample actions for each agent independently and to use the per-agent utility  $Q_a$  to guide the selection of maximal actions. We do not consider COVDN and COMIX significant algorithmic contributions but instead merely baseline algorithms. See Appendix C for more details about them.

**FACMAC outperforms MADDPG and other baselines in both discrete and continuous action tasks.** Figure 4 and 5 illustrate the mean episode return attained by different methods on Continuous Predator-Prey with varying number of agents and different MAMuJoCo tasks, respectively. We can see that FACMAC significantly outperforms MADDPG on all these continuous cooperative tasks, both in terms of absolute performance and learning speed. On discrete SMAC tasks, Figure 6 shows that FACMAC performs significantly better than MADDPG on 4 out of 6 maps we tested, and achieves similar performance to MADDPG on the other 2 maps. Additionally, on all 6 SMAC maps, FACMAC significantly outperforms all multi-agent actor-critic baselines (COMA, CentralV, DOP, and VDAC-mix), while DOP is recently claimed to be the first multi-agent actor-critic method that outperforms state-of-the-art valued-based methods on SMAC. FACMAC is also competitive with state-of-the-art value-based methods (QMIX and QPLEX), with significantly better performance on *MMM*, *bane\_vs\_bane*, *MMM2*, and *27m\_vs\_30m*. These results demonstrate the benefits of our method for improving performance in challenging cooperative tasks with discrete and continuous action spaces.

In Continuous Predator-Prey, FACMAC-vdn scales better than FACMAC when the number of agents increases. However, on MAMuJoCo, FACMAC-vdn performs drastically worse than FACMAC in 2-Agent Humanoid and ManyAgent Swimmer (with 10 agents), demonstrating the necessity of the non-linear mixing of agent utilities and conditioning on the central state information in order to achieve competitive performance in such tasks. Furthermore, on SMAC, Figure 13 in Appendix F shows that FACMAC is noticeably more stable than FACMAC-vdn and FACMAC-vdn-s across different maps, and achieves significantly better performance on the *super hard* map *MMM2*.

Interestingly, we find that FACMAC performs similarly to COMIX on both Continuous Predator-Prey and MAMuJoCo tasks. As FACMAC and COMIX use the same value factorisation as in QMIX and are both off-policy, this suggests that, in these continuous cooperative tasks, factorisation of the



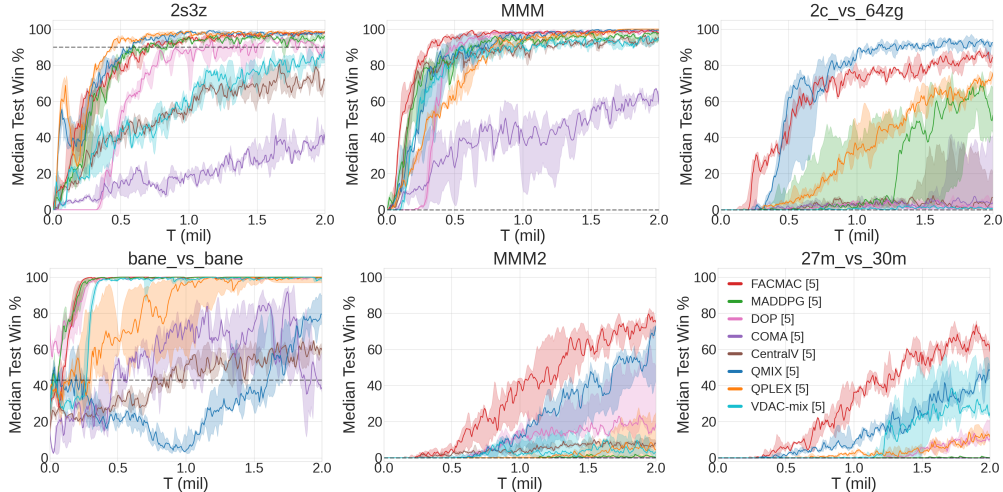


Figure 6: Median test win % on six different SMAC maps, including *2s3z* (easy), *MMM* (easy), *2c\_vs\_64zg* (hard), *bane\_vs\_bane* (hard), *MMM2* (super hard), and *27m\_vs\_30m* (super hard). The median across 5 seeds is plotted and the 25 – 75% percentiles is shown shaded. The performance of the heuristic-based algorithm is shown as a dashed line. We report the median instead of the mean as recommended by Samvelyan et al. [35] in order to avoid the effect of any outliers.

joint  $Q$ -value function plays a greater role in performance than the underlying algorithmic choices. On SMAC, however, FACMAC performs significantly better than QMIX on *MMM*, *bane\_vs\_bane*, *MMM2*, and *27m\_vs\_30m*. For instance, on *bane\_vs\_bane*, a task with 24 agents, while QMIX struggles to find the optimal policy with 2 million timesteps, FACMAC, with exactly the same value factorisation, can quickly recover the optimal policy and achieve 100% test win rate. This shows the convergence advantages of policy gradient methods in this type of multi-agent settings [43].

**FACMAC scales better as the number of agents (and/or actions) and the complexity of the task increases.** As shown in Figure 4(b) and 4(c), MADDPG performs poorly if we increase the number of agents in Continuous Predator-Prey, while both FACMAC and FACMAC-vdn achieve significantly better performance. The monolithic critic in MADDPG simply concatenates all agents’ observations into a single input vector, which can be quite large when there are many agents and/or entities and make it more difficult to learn a good critic. Factoring the critic enables scalable critic learning, by combining individual agent utilities that condition on much smaller observations into a joint action-value function. On MAMuJoCo (shown in Figure 5), similarly, the largest performance gap between FACMAC and MADDPG can be seen on ManyAgent Swimmer (with 10 agents), a task with the largest number of agents among three MAMuJoCo tasks tested.

On SMAC (shown in Figure 6), the largest performance gap between FACMAC and MADDPG can be seen on the challenging *MMM2* and *27m\_vs\_30m* with a large number of agents, which are classified as 2 *super hard* SMAC maps due to current methods’ poor performance [35]. We can see that FACMAC is able to learn to consistently defeat the enemy, whereas MADDPG fails to learn anything useful in both tasks. The second largest performance gap between FACMAC and MADDPG can be seen on the hard map *2c\_vs\_64zg*, where MADDPG not only performs significantly worse but also exhibits significantly more variance than FACMAC across seeds. While there are only 2 agents in this scenario, the number of actions each agent can choose is the largest among all 6 maps tested as there are 64 enemies. These results further demonstrate that FACMAC scales better when the number of agents (and/or actions) and the complexity of the tasks increases.

**Factoring the critic can better take advantage of our centralised gradient estimator to optimise the agent policies when the number of agents and/or actions is large.** We conduct ablation experiments to investigate the influence of factoring the critic and using the centralised gradient estimator in our method. FACMAC (without CPG) is our method without the centralised policy gradient. It uses a naive adaptation of the deterministic policy gradient used in MADDPG (shown in (6)). Thus, the only difference between FACMAC (without CPG) and MADDPG is that the previous

one learns a non-linearly factored critic while the latter one learns a monolithic critic. We also evaluate MADDPG with our centralised policy gradient and refer to it as MADDPG (with CPG).

Figure 7 shows the ablation results on SMAC and MAMuJoCo. We can see that FACMAC (without CPG) performs significantly better than MADDPG on both SMAC maps tested, demonstrating the advantages of factoring the critic in challenging coordination problems. With the centralised policy gradient, MADDPG (with CPG) performs significantly better than MADDPG on *2c\_vs\_64zg*. However, on the harder map *MMM2*, MADDPG with both policy gradients fail to learn anything useful. By contrast, FACMAC significantly outperforms FACMAC (without CPG) on *MMM2*, and has lower variance across seeds on *2c\_vs\_64zg*. Furthermore, on ManyAgent Swimmer with 2 agents, our centralised gradient estimator does not affect the performance of both MADDPG and FACMAC. However, on the same task with 10 agents, using the centralised gradient estimator significantly improves the learning performance when learning a factored critic. These results demonstrate that factoring the critic can better take advantage of our centralised gradient estimator to optimise the agent policies when the number of agents and/or actions is large.

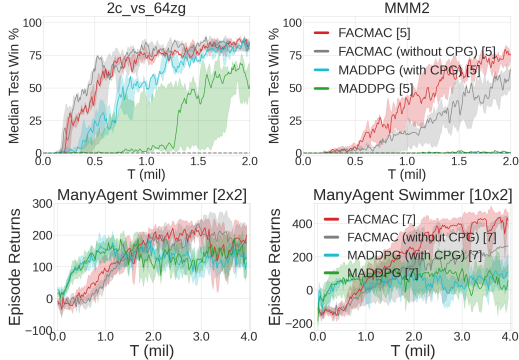


Figure 7: Ablations for different FACMAC components on SMAC and MAMuJoCo tasks.

**Nonmonotonically factored critics can solve tasks that cannot be solved with monolithic or monotonically factored critics.** In our multi-agent actor-critic framework, there are no inherent constraints on factoring the critic, we thus also employ a nonmonotonic factorisation and refer to it as FACMAC-nonmonotonic (as discussed in Section 3.1). As shown in Figure 8 (left), on our continuous matrix game (as discussed in Section 3.2), FACMAC-nonmonotonic can robustly learn the optimal policy, while both FACMAC and MADDPG converge to some sub-optimal policy within  $200k$  timesteps. On a variant of Continuous Predator-Prey task (see Appendix D.2 for more details about this task) with *nonmonotonic* value functions (i.e., an agent’s ordering over its own actions depends on other agents’ actions [25]), Figure 8 (right) shows that both FACMAC and MADDPG fail to learn anything useful, while FACMAC-nonmonotonic successfully learns to capture the prey. These results demonstrate that nonmonotonically factored critics can solve tasks that cannot be solved with monolithic or monotonically factored critics.

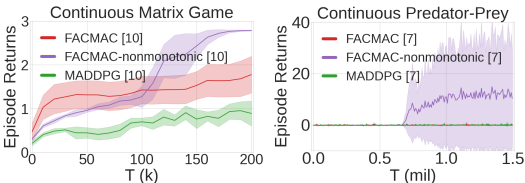


Figure 8: Mean episode return on (Left) Continuous Matrix Game and (Right) a variant of our Continuous Predator-Prey task (with 3 agents and 1 prey) with *nonmonotonic* value functions.

It is important to note that the relative performance of FACMAC and FACMAC-nonmonotonic is task dependent. On the original Continuous Predator-Prey task (with 3 agents and 1 prey), FACMAC-nonmonotonic yields similar performance to FACMAC (see Figure 12 in Appendix F). On SMAC (see Figure 13 in Appendix F), FACMAC-nonmonotonic performs similarly to FACMAC on easy maps, but exhibits significantly worse performance on harder maps. This shows that, in this type of tasks, using an unconstrained factored critic could lead to an increase in learning difficulty. We thus expect FACMAC-nonmonotonic to be more useful in tasks with nonmonotonic value functions.

## 6 Conclusion

This paper presented FACMAC, a multi-agent actor-critic method that learns decentralised policies with a centralised but factored critic, working for both discrete and continuous cooperative tasks. We showed the advantages of both factoring the critic and using the new centralised gradient estimator in our approach. We also introduced a novel benchmark suite MAMuJoCo to demonstrate FACMAC’s scalability to more complex continuous tasks. Our results on three different domains demonstrated FACMAC’s superior performance over existing MARL algorithms. Future work will explore more forms of nonmonotonic factorisation to tackle tasks with nonmonotonic value functions.

## Acknowledgements

We would like to thank Dr Martin Strohmeier and the members of the Whiteson Research Lab for their helpful feedback. We would also like to thank the anonymous reviewers for their constructive comments during the reviewing process. This project has received funding from the European Research Council (ERC), under the European Union’s Horizon 2020 research and innovation programme (grant agreement number 637713). It was also supported by an EPSRC grant (EP/M508111/1, EP/N509711/1). Christian Schroeder de Witt is generously funded by Cyber Defence Campus, Science and Technology, Armasuisse, Switzerland. The experiments were made possible by a generous equipment grant from NVIDIA.

## References

- [1] Johannes Ackermann, Volker Gabler, Takayuki Osa, and Masashi Sugiyama. Reducing overestimation bias in multi-agent domains using double centralized critics. *arXiv preprint arXiv:1910.01465*, 2019.
- [2] Brandon Amos, Lei Xu, and J Zico Kolter. Input convex neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 146–155. JMLR. org, 2017.
- [3] Ronald C Arkin. Motor schema—based mobile robot navigation. *The International journal of robotics research*, 8(4):92–112, 1989.
- [4] Marilena Bescuca. Factorised critics in deep multi-agent reinforcement learning. In *Master Thesis, University of Oxford*, 2019.
- [5] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [6] Rodney Brooks. A robust layered control system for a mobile robot. *IEEE journal on robotics and automation*, 2(1):14–23, 1986.
- [7] Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. A tutorial on the cross-entropy method. *Annals of operations research*, 134(1):19–67, 2005.
- [8] Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. LIIR: Learning individual intrinsic reward in multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems 32*, pages 4405–4416, 2019.
- [9] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [10] Shixiang Gu, Timothy Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous deep q-learning with model-based acceleration. In *International Conference on Machine Learning*, pages 2829–2838, 2016.
- [11] Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 66–83. Springer, 2017.
- [12] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- [13] Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 2961–2970, 2019.
- [14] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.

- [15] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- [16] Hiroaki Kitano, Minoru Asada, Yasuo Kuniyoshi, Itsuki Noda, Eiichi Osawa, and Hitoshi Matsubara. Robocup: A challenge problem for ai. *AI magazine*, 18(1):73–73, 1997.
- [17] Daphne Koller and Ronald Parr. Computing factored value functions for policies in structured mdps. In *Proceedings of IJCAI*, pages 1332–1339, 1999.
- [18] Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.
- [19] Haruhisa Kurokawa, Kohji Tomita, Akiya Kamimura, Shigeru Kokaji, Takashi Hasuo, and Satoshi Murata. Distributed self-reconfiguration of m-tran iii modular robotic system. *The International Journal of Robotics Research*, 27(3-4):373–386, 2008.
- [20] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *4th International Conference on Learning Representations, ICLR*, 2016.
- [21] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992.
- [22] Siqi Liu, Guy Lever, Josh Merel, Saran Tunyasuvunakool, Nicolas Heess, and Thore Graepel. Emergent coordination through competition. *arXiv preprint arXiv:1902.07151*, 2019.
- [23] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, pages 6379–6390, 2017.
- [24] Xueguang Lyu, Yuchen Xiao, Brett Daley, and Christopher Amato. Contrasting centralized and decentralized critics in multi-agent reinforcement learning. In *Proceedings of the 20th International Conference on Autonomous Agents and Multi-Agent Systems*, 2021.
- [25] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, pages 7613–7624, 2019.
- [26] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [27] Frans A. Oliehoek, Matthijs T. J. Spaan, and Nikos Vlassis. Optimal and approximate Q-value functions for decentralized pomdps. *JAIR*, 32:289–353, 2008.
- [28] Frans A Oliehoek, Christopher Amato, et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- [29] OpenAI. openai/baselines, May 2020. original-date: 2017-05-24T01:58:13Z.
- [30] Ling Pan, Tabish Rashid, Bei Peng, Longbo Huang, and Shimon Whiteson. Softmax with regularization: Better value estimation in multi-agent reinforcement learning. *arXiv preprint arXiv:2103.11883*, 2021.
- [31] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 4292–4301, 2018.
- [32] Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation. In *Advances in neural information processing systems*, 2020.

- [33] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *JMLR*, 21:178:1–178:51, 2020.
- [34] Martin Riedmiller, Thomas Gabel, Roland Hafner, and Sascha Lange. Reinforcement learning for robot soccer. *Autonomous Robots*, 27(1):55–73, 2009.
- [35] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The StarCraft Multi-Agent Challenge. *CoRR*, abs/1902.04043, 2019.
- [36] Christian Schroeder de Witt, Jakob Foerster, Gregory Farquhar, Philip Torr, Wendelin Boehmer, and Shimon Whiteson. Multi-Agent Common Knowledge Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [37] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:1905.05408*, 2019.
- [38] Kyunghwan Son, Sungsoo Ahn, Roben Delos Reyes, Jinwoo Shin, and Yung Yi. Qtran++: Improved value transformation for cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2006.12010*, 2020.
- [39] Peter Stone and Richard S. Sutton. Scaling reinforcement learning toward RoboCup soccer. In *Icml*, volume 1, pages 537–544. Citeseer, 2001.
- [40] Jianyu Su, Stephen Adams, and Peter A Beling. Value-decomposition multi-agent actor-critics. In *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence*, 2021.
- [41] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinícius Flores Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *AAMAS*, pages 2085–2087, 2018.
- [42] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.
- [43] Jianhao Wang, Zhizhou Ren, Beining Han, and Chongjie Zhang. Towards understanding linear value decomposition in cooperative multi-agent q-learning. *arXiv preprint arXiv:2006.00587*, 2020.
- [44] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling multi-agent q-learning. *arXiv preprint arXiv:2008.01062*, 2020.
- [45] Tingwu Wang, Renjie Liao, Jimmy Ba, and Sanja Fidler. NerveNet: learning structured policy with graph neural networks. In *6th International Conference on Learning Representations, ICLR*, 2018.
- [46] Yihan Wang, Beining Han, Tonghan Wang, Heng Dong, and Chongjie Zhang. Dop: Off-policy multi-agent decomposed policy gradients. In *9th International Conference on Learning Representations, ICLR*, 2021.
- [47] Ermo Wei and Sean Luke. Lenient learning in independent-learner stochastic cooperative games. *The Journal of Machine Learning Research*, 17(1):2914–2955, 2016.
- [48] Yaodong Yang, Jianye Hao, Ben Liao, Kun Shao, Guangyong Chen, Wulong Liu, and Hongyao Tang. Qatten: A general framework for cooperative multiagent reinforcement learning. *arXiv preprint arXiv:2002.03939*, 2020.
- [49] Mark Yim, Ying Zhang, and David Duff. Modular robots. *IEEE Spectrum*, 39(2):30–34, 2002.
- [50] Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit credit assignment for multi-agent actor-critic. In *Advances in neural information processing systems*, 2020.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] We mention that the relative performance of FACMAC and FACMAC-nonmonotonic is task dependent and discuss some empirical results in the last paragraph of Section 5.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Appendix B.1.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
  - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We include the code and instructions needed to reproduce the main experimental results in the supplemental material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendix E.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] For all results on Continuous Matrix Game, Continuous Predator-Prey, and MAMuJoCo, we report the mean episode return and the 95% confidence interval. For all results on SMAC, we report the median test win rate as well as the 25-75% percentiles (as in [35]).
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We describe the type of GPUs used in Appendix E.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [No] The license is already included in the original code resource we cited and could change in the future.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Related Work

Value function factorisation [17] has been widely employed in value-based MARL algorithms. VDN [41] and QMIX [31] factor the joint action-value function into per-agent utilities that are combined via a simple summation or a monotonic mixing function respectively, to ensure consistency between the  $\arg \max$  of the centralised joint-action value function and the  $\arg \max$  of the decentralised policies. This monotonicity constraint, however, prevents them from representing joint action-value functions that are characterised as *nonmonotonic* [25], i.e., an agent’s ordering over its own actions depends on other agents’ actions. A large number of recent works [37, 48, 32, 44, 38] thus focus on developing new value-based MARL algorithms that address this representational limitation, in order to learn a richer class of action-value functions.

QTRAN [37] learns an unrestricted joint action-value function and aims to solve a constrained optimisation problem in order to decentralise it, but has been shown to scale poorly to more complex tasks such as SMAC [25]. QPLEX [44] takes advantage of the dueling network architecture to factor the joint action-value function in a manner that does not restrict the representational capacity, whilst also remaining easily decentralisable, but can still fail to solve simple tasks with nonmonotonic value functions [32]. Weighted QMIX [32] introduces a weighting scheme to place more importance on the better joint actions to learn a richer class of joint action-value functions. QTRAN++ [38] addresses the gap between the empirical performance and theoretical guarantees of QTRAN. Our multi-agent actor-critic framework with decentralised actors and a centralised but factored critic, by contrast, provides a more direct and simpler way of coping with nonmonotonic tasks as one can simply factor the centralised critic in any manner without constraints. Additionally, our framework can be readily applied to tasks with continuous action spaces, whereas these value-based algorithms require additional algorithmic changes.

Most state-of-the-art multi-agent actor-critic methods [23, 9, 13, 8] learn a centralised and monolithic critic conditioning on the global state and the joint action to stabilise learning. Even though the joint action-value function they can represent is not restricted, in practice they significantly underperform value-based methods like QMIX on the challenging SMAC benchmark [32, 33]. In contrast, FACMAC utilises a centralised but factored critic to allow it to scale to the more complex tasks in SMAC, and follows the centralised policy gradient instead of per-agent policy gradients.

Lyu et al. [24] recently provide some interesting insights about the pros and cons of centralised and decentralised critics for on-policy actor-critic algorithms that use stochastic policy gradients. One important issue that they highlight is that merely utilising a centralised critic does not necessarily lead to the learning of more coordinated behaviours than using decentralised critics. This is because the use of a per-agent policy gradient can lead to the agents getting stuck in sub-optimal solutions in which no agents wishes to change their policy, as discussed in Section 3.2. Su et al. [40] recently propose the value-decomposition actor-critic (VDAC) framework that factors the centralised critic monotonically (i.e., factors the joint value function as a monotonic combination of local value functions) to learn more efficiently. However, even with a factored critic, VDAC still suffers from the problems caused by the per-agent policy gradient due to the use of a fully factored stochastic policy. Our centralised policy gradient resolves this issue by using deterministic policies and taking full advantage of the centralised training paradigm to optimise over the joint-action policy, which allows us to fully reap the benefits of a centralised critic. Furthermore, VDAC requires on-policy learning, which can be sample-efficient. In contrast, FACMAC is off-policy, we thus benefit immensely from utilising a centralised critic over a decentralised one since we avoid the issues of non-stationarity when training on older data.

Zhou et al. [50] propose to use a single centralised critic for MADDPG, whose weights are generated by hypernetworks that condition on the state, similarly to QMIX’s mixing network without the monotonicity constraints. FACMAC also uses a single centralised critic, but factorises it similarly to QMIX (not just using the mixing network) which allows for more efficient learning on more complex tasks. Of existing work, the deterministic decomposed policy gradients (DOP) algorithm proposed by Wang et al. [46] is perhaps most similar to our own approach. Deterministic DOP is off-policy and factors the centralised critic as a weighted linear sum of individual agent utilities and a state bias. It is limited to only considering linearly factored critics, which have limited representational capacity, whilst we are free to choose any method of factorisation in FACMAC to allow for the learning of a richer class of action-value functions. While they claim to be the first to introduce the idea of value function factorisation into the multi-agent actor-critic framework, it is actually first explored

by Bescuca [4], where a monotonically factored critic is learned for COMA [9]. However, their performance improvement on SMAC is limited since COMA requires on-policy learning and it is not straightforward to extend COMA to continuous action spaces. Furthermore, both works (and VDAC) only consider monotonically factored critics, whilst we employ a nonmonotonic factorisation and demonstrate its benefits. We also investigate the benefits of learning a centralised but factored critic more thoroughly, providing a better understanding about the type of tasks that can benefit more from a factored critic. Furthermore, both deterministic DOP and LICA [50] use a naive adaptation of the deterministic policy gradient used by MADDPG and suffer from the same problems as discussed in Section 3.2, while our centralised policy gradients allow for better coordination across agents in certain tasks.

## B Multi-Agent MuJoCo

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 [23] 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 [16, 39, 34] does not currently feature an easily configurable software interface for MARL, nor suitable AI-controlled benchmark opponents. Liu et al. [22] 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.

To demonstrate FACMAC’s scalability to more complex continuous domains and to stimulate more progress in continuous MARL, we develop *Multi-Agent MuJoCo* (MAMuJoCo), a novel benchmark for continuous cooperative multi-agent robotic control. Starting from the popular fully observable single-agent robotic MuJoCo [42] control suite included with OpenAI Gym [5], we create a wide variety of novel scenarios in which multiple agents within a single robot have to solve a task cooperatively. Single-robot multi-agent tasks in MAMuJoCo arise 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 3. We then partition the body graph into disjoint sub-graphs, one for each agent, each of which contains one or more joints that can be controlled. Note that in ManyAgent Swimmer (see Figure 3A) and ManyAgent Ant (see Figure 3K), the number of agents are not limited by the given single robotic agent.

Multiple agents are introduced within a single robot as partial observability arises through latency, bandwidth, and noisy sensors in a single robot. Even if communication is free and instant when it works, we want policies that keep working even when communication channels within the robot malfunction. Without access to the exact full state, local decision rules become more important and introducing autonomous agents at individual decision points (e.g., each physical component of the robot) is reasonable and beneficial. This also makes it more robust to single-point failures (e.g., broken sensors) and more adaptive and flexible as new independent decision points (thus agents) may be added easily. Similar physical decompositions can be found in early behavior-based robotic research [6, 3].

Each agent’s action space in MAMuJoCo is given by the joint action space over all motors controllable by that agent. For example, the agent corresponding to the green partition in 2-Agent HalfCheetah (Figure 3C) consists of three joints (joint ids 1, 2, and 3) and four adjacent body segments. Each joint has an action space  $[-1, 1]$ , so the action space for each agent is a 3-dimensional vector with each entry in  $[-1, 1]$ .

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 (see Figure 9 for an example). 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 the environment and configuration, representation vectors may include attributes such as position, velocity, and external body forces. In addition to joint and body-segment specific observation categories, agents can also be configured to observe the position and/or velocity attributes of the robot’s central torso.



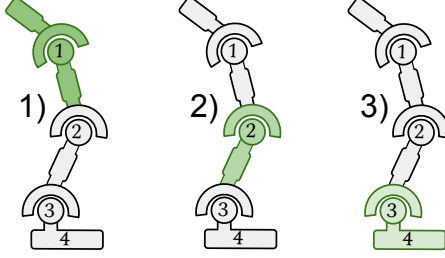


Figure 9: **Observations by distance for 3-Agent Hopper (as seen from agent 1).** Each corresponds to joints and body parts observable at 1) zero graph distance from agent 1, 2) one unit graph distance from agent 1, and 3) two unit graph distances from agent 1.

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, 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. Default team reward is summarised in Table 1.

Task	Goal	Special observations	Reward function
2-Agent Swimmer	Maximise +ve $x$ -speed.	All agents can observe velocities of the central torso.	$\frac{\Delta x}{\Delta t} + 0.0001\alpha$
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 velocities of 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 velocities of 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.	-	$\frac{\Delta x}{\Delta t} + \min(10, 5 \cdot 10^{-6} \ \text{external contact forces}\ _2^2)$
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 velocities of 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.	-	$0.25 \frac{\Delta x}{\Delta t} + \min(10, 5 \cdot 10^{-6} \ \text{external contact forces}\ _2^2)$
ManyAgent Swimmer	Maximise +ve $x$ -speed.	All agents can observe velocities of the central torso.	$\frac{\Delta x}{\Delta t} + 0.0001\alpha$
ManyAgent Ant	Maximise +ve $x$ -speed.	All agents can observe velocities of the central torso.	$\frac{\Delta x}{\Delta t} + 5 \cdot 10^{-4} \ \text{external contact forces}\ _2^2 + 0.5\alpha + 1$

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

The most similar existing environments, though not as diverse as MAMuJoCo, are the decomposed MuJoCo environments Centipede and Snakes [45]. The latter is similar to MAMuJoCo’s 2-Agent Swimmer. Ackermann et al. [1] evaluate on one environment similar to a configuration of 2-Agent Ant, but, similarly to Gupta et al. [11], do not consider tasks across different numbers of agents and MuJoCo scenarios.

### B.1 Social Impact

MAMuJoCo is a simulated environment that mimics certain fundamental aspects of robotic actuator chains. Such actuator chains are commonly found in industrial robotics, perhaps most prominently in robot arms on assembly lines, but also a variety of other civil use cases, including remote space exploration or bomb disposal. However, in order to be applicable to such use cases, policies trained in MAMuJoCo first need to undergo a transfer process to the real world, for example by using Sim2Real techniques. MAMuJoCo can not only help assess algorithmic advances, but also contribute toward understanding and mitigating risks inherent to civil autonomous robotics applications. For example, MAMuJoCo could facilitate the development of safe reinforcement learning algorithms that keep

actuator parameters of real-world robotic arms within safe ranges, thus avoiding manufacturing errors or injury.

As most developments in robotic control, algorithmic progress based on our benchmark environment, or the novel algorithms introduced in this paper, might ultimately find application in autonomous warfare. MAMuJoCo has been developed with fundamental algorithmic development for civil purposes in mind and does not provide any features that would make it particularly suitable to non-civil use. Any robotic agents trained with reinforcement learning algorithms on our environment, including our novel methods, should be certified with respect to fairness and security before real-world deployment. Like any artificial intelligence system, our proposed framework has the potential to greatly improve human productivity. However, it may also reduce the need for human workers, resulting in job losses.

## C COVDN and COMIX

$Q$ -learning has shown considerable success in multi-agent settings with discrete action spaces [41, 31, 37, 32, 44]. However, performing greedy action selection in  $Q$ -learning requires evaluating  $\arg \max_{\mathbf{u}} Q_{tot}(\boldsymbol{\tau}, \mathbf{u}, s)$ , where  $Q_{tot}$  is the joint 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 [10, 2], at the expense of estimation bias, or perform only approximate greedy action selection [15]. 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  $Q_{tot}$  required for greedy action selection 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, COVDN represents the joint action-value function  $Q_{tot}$  as a sum of the per-agent utilities  $Q_a$  as in VDN [41], while COMIX represents  $Q_{tot}$  as a non-linear monotonic combination of  $Q_a$  as in QMIX [31]. COVDN and COMIX are thus simple variants of VDN and QMIX, respectively, that scale to continuous action spaces. They both perform approximate greedy selection of actions  $u_a$  with respect to utility functions  $Q_a$  for each agent  $a$  using the cross-entropy method (CEM) [7]. CEM is 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 [15]. The centralised but factored  $Q_{tot}$  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.

In both COVDN and COMIX, CEM is used by each agent  $a$  to find an action that approximately optimises its local utility function  $Q_a$ . Specifically, 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 (with the highest utility values) are then used to fit a new Gaussian distribution  $\mathcal{D}_{k+1}$ , and this process repeats  $K$  times. We use a CEM hyperparameter configuration similar to Qt-Opt [15], where  $N = 64$ ,  $M = 6$ , and  $K = 2$ .<sup>7</sup> Gaussian distributions are initialised with mean  $\mu = 0$  and standard deviation  $\sigma = 1$ . Algorithm 1 outlines the full CEM process used in both COVDN and COMIX. Algorithm 2 outlines the full process for COMIX. Note we do not consider COVDN and COMIX significant algorithmic contributions but instead merely baseline algorithms.

## D Environment Details

### D.1 Continuous Matrix Game

Figure 10 illustrates the continuous matrix game with two agents. There is a narrow path (shown in red) starting from the origin  $(0, 0)$  to  $(1, 1)$ , in which the reward gradually increases. Everywhere else there is a small punishment moving away from the origin, increasing in magnitude further from the origin.

<sup>7</sup>We empirically find 2 iterations to suffice.

---

**Algorithm 1** 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.

---

```

for  $a := 1, a \leq N$  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  $i := 1, i \leq n_c$  do
    for  $j := 1, j \leq d_i$  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
       $U \leftarrow \{\mathbf{v}'_{al} \mid q_{al} \in \text{top}^{k_i}(q_{a1}, \dots, q_{ad_i}), \forall l \in \{1 \dots N\}\}$ 
       $\boldsymbol{\mu}_a \leftarrow \text{sample\_mean}(U)$ 
       $\boldsymbol{\sigma}_a \leftarrow \text{sample\_std}(U)$ 
    else
       $m \leftarrow \arg \max_j q_{aj}$ 
       $\mathbf{u}_a \leftarrow \mathbf{v}_{am}$ 
    end if
     $i \leftarrow i + 1$ 
  end for
   $a \leftarrow a + 1$ 
end for
return  $\langle \mathbf{u}_1, \dots, \mathbf{u}_n \rangle$ 

```

---



---

**Algorithm 2** Algorithmic description of COMIX. The function CEM is defined in Algorithm 1.

---

```

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
     $\mathbf{u}_t \leftarrow \text{CEM}(Q_1, \dots, Q_N, \tau_t^1, \dots, \tau_t^N)$ 
     $\langle s_{t+1}, \mathbf{z}_{t+1}, r_t \rangle \leftarrow \text{EnvStep}(\mathbf{u}_t)$ 
    ReplayBuffer  $\leftarrow \langle s_t, \mathbf{u}_t, \mathbf{z}_t, r_t, s_{t+1}, \mathbf{z}_{t+1} \rangle$ 
  end for
   $\{\langle s_i, \mathbf{u}_i, \mathbf{z}_i, r_i, s'_i, \mathbf{z}'_i \rangle\}_{i=1}^b \sim \text{ReplayBuffer}$ 
   $y_i \leftarrow r_i + \gamma \max_{\mathbf{u}'_i} Q_{\text{tot}}(s'_i, \mathbf{z}'_i, \mathbf{u}'_i; \boldsymbol{\theta}^-, \phi^-), \forall i$ 
   $\mathcal{L} \leftarrow \sum_{i=1}^b \left( y_i - Q_{\text{tot}}(s_i, \mathbf{z}_i, \mathbf{u}_i; \boldsymbol{\theta}, \phi) \right)^2$ 
   $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}$ 
   $\phi \leftarrow \phi - \alpha \nabla_{\phi} \mathcal{L}$ 
end for

```

---

## D.2 Continuous Predator-Prey

We consider the mixed *simple tag* environment (Figure 11) introduced by Lowe et al. [23], which 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 of +10 is emitted; otherwise, no reward is given. In the original simple tag environment, each agent can observe the relative positions of the other two agents, the relative position and velocity of the prey, and the relative

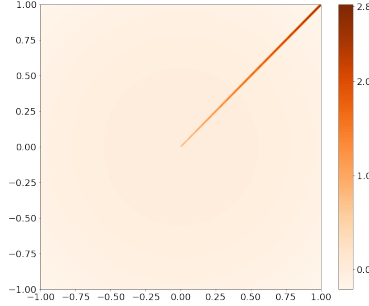


Figure 10: The continuous matrix game.

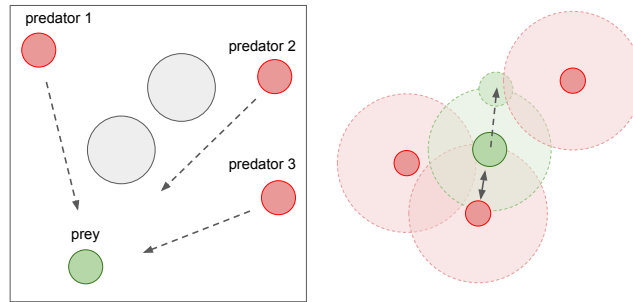


Figure 11: Continuous Predator-Prey. **Left:** Top-down view of toroidal plane, with predators (red), prey (green) and obstacles (grey). **Right:** Illustration of the prey’s avoidance heuristic. Observation radii of both agents and prey are indicated.

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 view radius such that the agents can only observe other agents roughly 60% of the time. We refer to this environment as Continuous Predator-Prey.

In addition, we implement a variant of our Continuous Predator-Prey task (with 3 agents and 1 prey), where the reward function is modified to make the task *nonmonotonic*. Specifically, if one agent collides with the prey while at least another one being close enough, a team reward of +10 is given. However, if only one agent collides with the prey without any other agent being close enough, a negative team reward of  $-1$  is given. Otherwise, no reward is provided.

### D.3 Multi-Agent MuJoCo

All MAMuJoCo environments we tested are configured according to its default configuration, where each agent can observe only positions (not velocities) of its own body parts and at graph distances greater than zero. In ManyAgent Swimmer, we configure the number of agents to be 10, each controlling a consecutive segment of length 2. We thus refer to this environment as ManyAgent Swimmer [10x2]. We set maximum observation distances to  $k = 0$  (which means each agent can observe only positions of its own body parts) for all three environments tested, including 2-Agent Humanoid, 2-Agent HumanoidStandup, and ManyAgent Swimmer [10x2]. Default team reward is used (see Table 1).

### D.4 SMAC

SMAC consists of a set of complex StarCraft II micromanagement tasks that are carefully designed to study decentralised multi-agent control. The tasks in SMAC involve combat between two armies

of units. The first army is controlled by a group of learned allied agents. The second army consists of enemy units controlled by the built-in heuristic AI. The goal of the allied agents is to defeat the enemy units in battle, to maximise the win rate. The action space consists of a set of discrete actions: move in four cardinal directions, attack any selected enemy (available if the enemy is within the agent’s shooting range), stop, and noop. Hence the number of actions increases as the number of enemies increases. All experiments on SMAC use the default reward and observation settings of the SMAC benchmark [35].

## E Experimental Details

We evaluate the performance of each method using the following procedure: for each run of a method, we pause training every fixed number of timesteps (2000 timesteps for Continuous Predator-Prey, 4000 timesteps for MAMuJoCo, and 10000 timesteps for SMAC) and run a fixed number of independent episodes (10 episodes for Continuous Predator-Prey and MAMuJoCo, and 32 episodes for SMAC) with each agent performing action selection greedily in a decentralised fashion. On both Continuous Predator-Prey and MAMuJoCo, the mean value of these episode returns are used to evaluate the performance of the learned policies. On SMAC, we use the median test win rate (i.e., the percentage of the 32 episodes where the agents defeat all enemy units within the permitted time limit) to evaluate the learned policies, as in [35]. All experiments are carried out on NVIDIA GeForce GTX 1080 GPU.

### E.1 Continuous Predator-Prey

In value-based methods COVDN and COMIX, the architecture of the shared agent network is a DRQN with a recurrent layer comprised of a GRU with a 64-dimensional hidden state, with a fully-connected layer before and after. In actor-critic methods FACMAC, FACMAC-vdn, MADDPG, and IDDPG, the architecture of the shared agent network is also a DRQN with a recurrent layer comprised of a GRU with a 64-dimensional hidden state, with a fully-connected layer before and after, while the final output layer is a tanh layer, to bound actions. The shared critic network is a MLP with 2 hidden layers of 64 units and ReLU non-linearities. All agent networks receive the current local observation and last individual action as input. In MADDPG, the centralised critic receives the global state and the joint action of all agents as input. The global state consists of the joint observations of all agents in Continuous Predator-Prey. In other actor-critic methods, there is a shared critic network that approximates per-agent utilities, which receives each agent’s local observation and individual action as input.

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 with  $\sigma = 0.1$  during training (for all 2 million timesteps). We set  $\gamma = 0.85$  for all experiments. The replay buffer contains the most recent 5000 transitions. We train on a batch size of 32 after every timestep. The hyper-parameter *batch\_size\_run* is set to be 1. For the soft target network updates we use  $\tau = 0.001$ . All neural networks (actor and critic) are trained using Adam optimiser with a learning rate of 0.01. To evaluate the learning performance, the training is paused after every 2000 timesteps during which 10 independent test episodes are run with agents performing action selection greedily in a decentralised fashion.

In FACMAC, the architecture of the mixing network consists of a single hidden layer of 64 units with an ELU non-linearity. The weights of the mixing network are produced by separate hypernetworks. The hypernetworks consist of a feedforward network without any hidden layers. The output of the hypernetwork is passed through an absolute activation function (to achieve non-negativity) and then resized into a matrix of appropriate size. In FACMAC-nonmonotonic, the architecture of the mixing network (and the hypernetworks) is similar to the one used in FACMAC, but without the constraint of monotonicity enforced by using non-negative weights.

### E.2 Multi-Agent MuJoCo

In all value-based methods, 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.<sup>8</sup> All agent

<sup>8</sup><https://spinningup.openai.com/en/latest/>.

networks use ReLU non-linearities for all hidden layers. In all actor-critic methods, the architecture of the shared agent network and critic network is also a MLP with 2 hidden layers with 400 and 300 units respectively, while the final output layer of the actor network is a tanh layer, to bound the actions. In all value-based methods, the agent receives its current local observation as input. In MADDPG, the centralised critic receives the global state and the joint action of all agents as input. The global state consists of the full state information returned by the original OpenAI Gym [5]. In other actor-critic methods, there is a shared critic network that approximates per-agent utilities, which receives each agent’s local observation and individual action as input.

During training and testing, we restrict each episode to have a length of 1000 time steps. Training lasts for 2 million or 4 million timesteps. To encourage exploration, we use uncorrelated, mean-zero Gaussian noise with  $\sigma = 0.1$  during training. 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. The replay buffer contains the most recent  $10^6$  transitions. We train on a batch size of 100 after every timestep. The hyper-parameter *batch\_size\_run* is set to be 1. For the soft target network updates we use  $\tau = 0.001$ . All neural networks (actor and critic) 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 independent test episodes are run with agents performing action selection greedily in a decentralised fashion.

In FACMAC, the architecture of the mixing network consists of a single hidden layer of 64 units with an ELU non-linearity. The weights of the mixing network are produced by separate hypernetworks. The hypernetworks consist of a feedforward network with a single hidden layer of 64 units with a ReLU non-linearity. The output of the hypernetwork is passed through an absolute activation function (to achieve non-negativity) and then resized into a matrix of appropriate size.

### E.3 SMAC

For baseline algorithms DOP [46], COMA [9], CentralV [9], VDAC-mix [40], QMIX [33], and QPLEX [44], we use the the same training setup as provided by their authors where the hyperparameters have been fine-tuned on the SMAC benchmark.

Most of our training hyperparameters for FACMAC and MADDPG [23] follow [33]. In both methods, the architecture of the shared actor network is a DRQN with a recurrent layer comprised of a GRU with a 64-dimensional hidden state, with a fully-connected layer before and after. The shared critic network is a MLP with 2 hidden layers of 64 units and ReLU non-linearities. Exploration is performed during training using a scheme similar to COMA [9]. Action probabilities are produced from the final layer of the actor network,  $z$ , via a bounded softmax distribution that lower-bounds the probability of any given action by  $\epsilon/|U|$ :  $P(u) = (1 - \epsilon)\text{softmax}_u + \epsilon/|U|$ , where  $|U|$  is the size of the joint action space. Throughout the training, for exploration we anneal  $\epsilon$  linearly from 0.5 to 0.05 over  $50k$  timesteps and keep it constant for the rest of the training. On map *MMM2*, we anneal  $\epsilon$  linearly from 1.0 to 0.05 over  $50k$  timesteps and keep it constant for the rest of the training for FACMAC and MADDPG (as it gives better performance). The replay buffer contains the most recent 5000 episodes. We sample batches of 32 episodes uniformly from the replay buffer and train on fully unrolled episodes. In MADDPG, we use a target network for the actor and critic, respectively. In FACMAC, we use a target network for the actor, critic, and mixing network, respectively. All target networks are periodically updated every 200 training steps. All neural networks are trained using Adam optimiser with learning rate 0.0025 for the actor network and 0.0005 for the critic network. We set  $\gamma = 0.99$  for all experiments.

The architecture of the mixing network in FACMAC follows [33]. It consists of a single hidden layer of 32 units with an ELU non-linearity. The weights of the mixing network are produced by separate hypernetworks. The hypernetworks consist of a feedforward network with a single hidden layer of 64 units with a ReLU non-linearity. The output of the hypernetwork is passed through an absolute activation function (to achieve non-negativity) and then resized into a matrix of appropriate size. In FACMAC-nonmonotonic, the architecture of the mixing network (and the hypernetworks) is similar to the one used in FACMAC, but without the constraint of monotonicity enforced by using non-negative weights.

## F Additional Results on Different Critic Factorisations

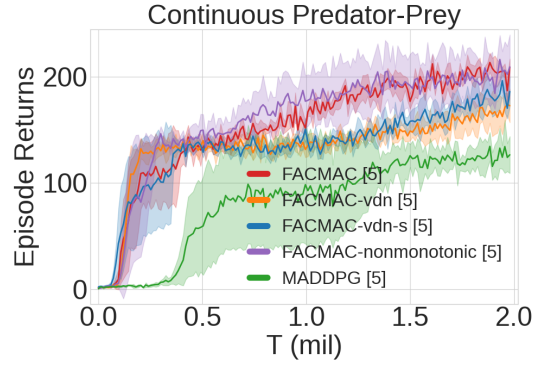


Figure 12: Mean episode return on our Continuous Predator-Prey task (with 3 agents and 1 prey).

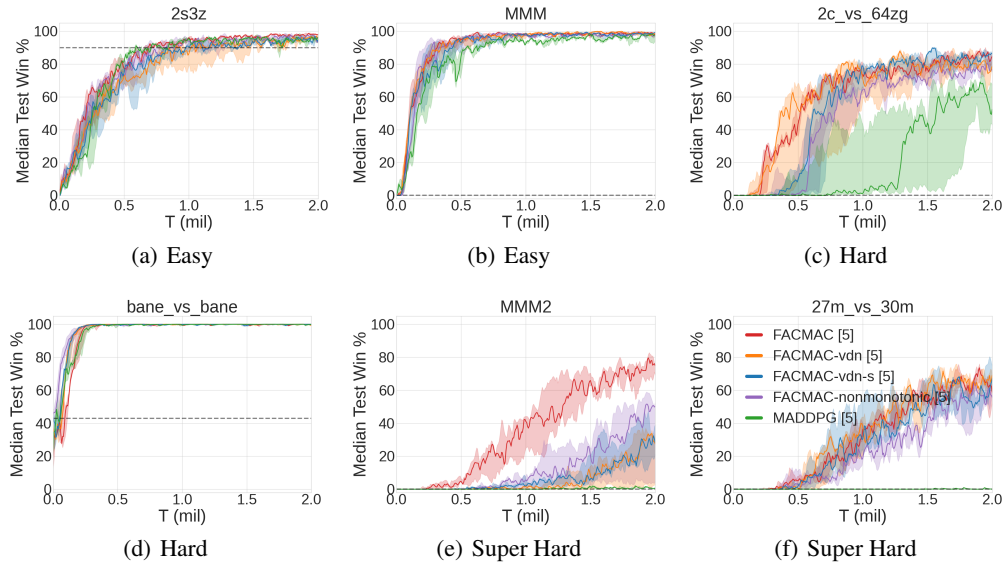


Figure 13: Median test win % on six different SMAC maps: (a) *2s3z* (easy), (b) *MMM* (easy), (c) *2c\_vs\_64zg* (hard), (d) *bane\_vs\_bane* (hard), (e) *MMM2* (super hard), and (f) *27m\_vs\_30m* (super hard), comparing FACMAC with different forms of critic factorisations.