

A Recursive Policy Gradient Algorithm for Fast Multi-Agent Reinforcement Learning

Anonymous authors
Paper under double-blind review

Abstract

Policy gradient algorithms for deep multi-agent reinforcement learning (MARL) typically employ an update that responds to the current strategies of other agents. While being straightforward, this approach does not account for the updates of other agents within the same update step, resulting in miscoordination and reduced sample efficiency. In this paper, we introduce methods that recursively refine the policy gradient by updating each agent against the updated policies of other agents within the same update step, speeding up the discovery of effective coordinated policies. We provide principled implementations of recursive reasoning in MARL by applying it to competitive multi-agent algorithms in both on and off-policy regimes. Empirically, we demonstrate superior performance and sample efficiency over existing deep MARL algorithms in StarCraft II and multi-agent MuJoCo. We theoretically prove that higher recursive reasoning in gradient-based methods with finite iterates achieves monotonic convergence to a local Nash equilibrium under certain conditions.

1 Introduction

Deep multi-agent reinforcement learning (MARL) research has made strides towards solving practical problems such as cooperative robotics (Ismail et al., 2018), transportation management (Haydari & Yilmaz, 2020) and network traffic optimization (Pi et al., 2024). While deep RL algorithms have garnered impressive results in complex single-agent control problems (Mnih et al., 2015; Tang et al., 2024), multi-agent systems present unique challenges. Roadblocks in MARL research include exploding joint state-action spaces and non-stationarity due to concurrent learning (Li et al., 2009; Barfuss & Mann, 2021). Another significant challenge caused by simultaneous learning in MARL is that each agent’s update does not account for the updates of other agents in the same update step, resulting in reduced sample efficiency (Zhang et al., 2021).

In this paper, we propose a recursive policy gradient update that allows each agent to reason about the change in behavior of other agents. Idealistic multi-agent frameworks assume mutual consistency: the assumption that each agent’s beliefs about other agents’ behavior and updates are accurate (Robertson, 1936). Due to limited computation, practical multi-agent systems update each agent as if the policy of every other agent is fixed - an assumption known as fictitious play (Foster & Young, 1998). Mutual consistency in the deep MARL setting can be approximated by using the updated policies of other agents in order to recursively refine the policy gradient. We introduce a recursive on-policy algorithm we term **ReMAPPO** (based on MAPPO (Yu et al., 2022)) which uses importance sampling for recursive updates. We also apply recursive reasoning to off-policy algorithms which utilize the deterministic policy gradient theorem (Silver et al., 2014); we term these implementations **ReFACMAC** and **ReMADDPG**, based on FACMAC (Peng et al., 2021) and MADDPG (Lowe et al., 2017) respectively. Finally, we motivate our method further by conducting a theoretical study of higher-level recursive reasoning with policy gradients and show that it results in bounded convergence to an ϵ -Nash equilibrium under certain conditions with finite iterates.

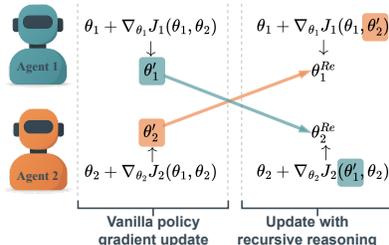


Figure 1: Illustration of recursive policy gradients for two agents.

2 Related works

Policy gradient methods in multi-agent reinforcement learning Policy gradient methods often estimate action-values of joint actions taken during training. MADDPG utilizes a multi-agent extension of the deterministic policy gradient (Silver et al., 2014) and exhibits more robust behavior than independent DDPG (Lillicrap et al., 2015). Foerster et al. (2018) propose COMA, which uses a baseline term to reduce centralized gradient noise, improving credit assignment. Similarly, Du et al. (2019) implement a framework that learns a proxy reward in order to discriminatively credit agents in multi-agent actor-critic methods. Yu et al. (2022) conduct a comprehensive study on MAPPO, a variant of Proximal Policy Optimization (PPO) (Schulman et al., 2017) which conditions its advantage estimation on global state information.

Recursive reasoning via opponent shaping While policy gradient methods can display impressive results in single-agent settings, they suffer from instability due to a lack of mutual consistency. Recursive reasoning and lookahead techniques have been used in several works to counter this. LOLA (Foerster et al., 2017) utilizes higher-order gradient terms to mutually shape the learning updates of agents in two-player reciprocity-based games. COLA (Willi et al., 2022) correct the naive-opponent assumption of LOLA by formalizing a notion of consistency and proves convergence properties of higher-order LOLA. POLA (Zhao et al., 2022) builds on LOLA further, reinterpreting it as a proximal operator by penalizing divergence over policy behavior, mitigating LOLA’s sensitivity to parameterizations. M-FOS (Khan et al., 2023) achieves model-free opponent shaping via a meta-game. While opponent shaping is a conceptually similar paradigm to our work, they differ algorithmically as they require meta-self-play formulations or higher order derivatives which assume white-box access to an opponent’s differentiable learning parameters - while ReMAPPO/ReFACMAC only require actions.

Game-theoretic frameworks of recursion The ideas of recursive reasoning and rationality hierarchies predate MARL, as various game theoretic models such as k-level reasoning, cognitive hierarchies (Camerer et al., 2004), recursive reasoning models (Gmytrasiewicz et al., 1991), fictitious play (Monderer & Shapley, 1996), and quantal hierarchies (Evans & Prokopenko, 2021) show. Notably, the cognitive hierarchies framework has been combined extensively with deep RL techniques for training self-driving vehicles to cooperate within a heterogeneous population of peer vehicles (Wang et al., 2022; Karimi et al., 2023; Dai et al., 2023). K-level reasoning has also been used for zero-shot coordination in Hanabi (Cui et al., 2021) by synchronously training multiple levels of competence with ad-hoc teamplay. GR2 (Wen et al., 2019a) models hierarchical levels of bounded rationality as a probabilistic graph model, which is shown to converge to a stationary point in 2-player games. Ma et al. (2022) uses message passing in order to recursively model opponent actions and build a centralized best-response model. The PR2 (Wen et al., 2019b) framework adopts variational Bayes methods to approximate conditional policies of opponents and proves convergence in self-play games. These methods fall out of the scope of this work as they aim to model the policies or rationalities of other agents, whereas we use opponent behavior to directly inform a recursive update.

Recursion to tackle rotational dynamics in games A notable example is Symplectic Gradient Adjustment (Balduzzi et al., 2018), which corrects rotational instability in N-player games by decomposing the game Hessian into symmetric and antisymmetric components and applying a second-order correction that dampens rotational dynamics. Liu & Pavel (2022) propose Level k gradient play, a recursive reasoning algorithm which stabilizes GAN training (Goodfellow et al., 2020).

Multi-agent performance difference lemma The Multi-Agent Performance Difference Lemma (MAPDL) is an extension of the PDL (Kakade & Langford, 2002) applied to joint action spaces. It has been used in Zhao et al. (2023) to prove local optimality in factored policy MARL updates, and A2PO (Wang et al., 2023)/HATPRO (Kuba et al., 2021) establish guarantees on sequential update convergence using the MAPDL. We utilize the MAPDL to derive the surrogate loss of ReMAPPO, but we omit these works as they do not fall within the simultaneous update scheme.

3 Preliminaries

We consider a multi-agent extension of a Markov decision process (Puterman, 2014) known as a Markov game (Littman, 1994), defined by a tuple $\mathcal{G} = \langle \mathcal{I}, \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, N, \iota \rangle$. $\mathcal{I} \equiv \{1, \dots, N\}$ is the set of agents, \mathcal{S} is the state space, $\mathcal{A} \equiv \times_{i \in \mathcal{I}} A_i$ is the joint action space of the agents. At each timestep t , agent i (the ‘self’ agent) samples an action $a_i \in A_i$ from policy $\pi_i(a_i|s)$ in state $s \in \mathcal{S}$ parameterized by $\theta_i \in \mathbb{R}^{d_i}$, where d_i is the dimensionality of the parameterization. At each timestep, the ‘non-self’ agents sample a joint action $\mathbf{a}_{-i} \in \times_{j \in \mathcal{I} \setminus i} A_j$ from joint policy $\pi_{-i}(\cdot|s)$ parameterized by joint non-self parameters θ_{-i} . The joint action of all agents \mathbf{a} from the joint policy $\pi(\cdot|s)$ determines the next state according to the joint state transition function $\mathcal{P}(s'|s, \mathbf{a})$. In the case of deterministic policies, we denote the self, non-self, and joint policies as μ_i , μ_{-i} , and μ respectively. $\mathcal{R} \equiv \{R_1, \dots, R_N\}$ are the set of agent reward functions. Each agent i has a learning rate η_i and receives a reward $r_{i,t}$ at time t according to its reward function $R_i(s, \mathbf{a}, s')$. γ is the discount factor and ι is the initial state distribution. We define the joint value function for agent i as $V_i^\pi(s) = \mathbb{E}_{\pi, \mathcal{P}} [\sum_{k=0}^{\infty} \gamma^k r_{i,k}|s]$ as the sum of discounted rewards for agent i following the joint policy π from state s . Similarly, we define the joint action-value function for agent i as $Q_i^\pi(s, \mathbf{a}) = \mathbb{E}_{\pi, \mathcal{P}} [\sum_{k=0}^{\infty} \gamma^k r_{i,k}|s, \mathbf{a}]$ as the sum of discounted rewards for agent i after taking joint action \mathbf{a} in state s and following the joint policy π thereafter. We define the advantage function for agent i as $A_i^\pi(s, \mathbf{a}) = Q_i^\pi(s, \mathbf{a}) - V_i^\pi(s)$. We denote $d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi, \iota)$ as the discounted state distribution under joint policy π and starting state distribution ι , where s_t is the state visited at time t . Each agent aims to maximize its own multi-agent objective, $J_i(\theta_i, \theta_{-i}) = \mathbb{E}_{s \sim \iota} [V_i^\pi(s)]$. Note that we condition objectives and loss functions on parameters θ to emphasize that gradient updates are made in weight space, while value functions are denoted with policies π to emphasize that actions depend on both weights and the form of the policy distribution.

4 Recursive Reasoning in multi-agent policy gradient algorithms

Multi-agent policy gradient algorithms are primarily concerned with estimating the gradient of the objective in Section 3. Typical estimations of $\nabla_{\theta_i} J_i(\theta_i, \theta_{-i})$ respond to the action distribution of other agents before they have made an update, under the assumption of fictitious play. We use θ'_i to denote agent i ’s parameters after a naive update:

$$\theta'_i \leftarrow \theta_i + \eta_i \nabla_{\theta_i} J_i(\theta_i, \theta_{-i}), \forall i \in \mathcal{I}. \quad (1)$$

Once θ'_i has been obtained for all agents, the update step can be taken once again from the *initial* parameters of each agent while considering the updated policies of the other agents. We use θ_i^{Re} to denote agent i ’s parameters after this recursive procedure:

$$\theta_i^{Re} \leftarrow \theta_i + \eta_i \nabla_{\theta_i} J_i(\theta_i, \theta'_{-i}), \forall i \in \mathcal{I}. \quad (2)$$

Successfully estimating Equation 2 in deep MARL settings is the primary aim of this work. Note that the recursive update of each agent is still only *one* gradient step away from the initial parameters; **the recursive updates are not moving further in weight space, but rather finding a gradient direction that is refined by the updated policies of the other agents.**

4.1 ReMAPPO

The performance difference lemma (PDL) (Kakade & Langford, 2002) can be applied to joint action spaces as shown in the following lemma, termed the Multi-Agent Performance Difference Lemma (MAPDL):

Lemma 4.1. *Given any joint policies π' and π , the difference in the performance of agent i under the joint policies can be expressed as:*

$$J_i(\pi') - J_i(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi'}, \mathbf{a} \sim \pi'_i, \pi'_{-i}} [A_i^\pi(s, \mathbf{a})], \quad (3)$$

The proof is a corollary of Lemma 6.1 in (Kakade & Langford, 2002) and is also provided in Appendix A.4 for completeness.

We thus derive the surrogate loss for ReMAPPO:

$$\mathcal{L}_i^{\text{ReMAPPO}} = \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_i, \pi_{-i}} [\min(r'_i(s, \mathbf{a}), \text{clip}(r'_i(s, \mathbf{a}), 1 - \epsilon, 1 + \epsilon)) \cdot A_i^\pi(s, \mathbf{a})], \quad (4)$$

where $r'_i(s, \mathbf{a}) = \frac{\pi'_i(a_i|s)}{\pi_i(a_i|s)} \cdot \frac{\pi'_{-i}(\mathbf{a}_{-i}|s)}{\pi_{-i}(\mathbf{a}_{-i}|s)}$ is the importance sampling ratio of the updated policies for actions taken in the environment and ϵ is a clipping boundary. Note that we assume the first-order approximation $d^{\pi'} \approx d^\pi$, as in Schulman et al. (2015).

This contrasts the surrogate loss of MAPPO:

$$\mathcal{L}_i^{\text{MAPPO}} = \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi_i, \pi_{-i}} [\min(r_i(s, \mathbf{a}), \text{clip}(r_i(s, \mathbf{a}), 1 - \epsilon, 1 + \epsilon)) \cdot A_i^\pi(s, \mathbf{a})], \quad (5)$$

where $r_i(s, \mathbf{a}) = \frac{\pi'_i(a_i|s)}{\pi_i(a_i|s)}$. Crucially, MAPPO *ignores* the update of non-self agents and assumes that non-self policies π_{-i} are stationary for simplicity. We emphasize that the surrogate loss in Equation 4 does not make this simplification, and is thus more appropriately informative for simultaneous-update multi-agent settings.

Intuitively, $r'_i(s, \mathbf{a}) = \frac{\pi'_i(a_i|s)}{\pi_i(a_i|s)} \cdot \frac{\pi'_{-i}(\mathbf{a}_{-i}|s)}{\pi_{-i}(\mathbf{a}_{-i}|s)}$ in Equation 4 can be seen as a weighting which places more or less weight on each sample depending on the updates of non-self agents. If the advantage is positive and other agents increase their action probabilities, the self-agent is incentivized to increase its probability even more in response. If the advantage is positive but the other agents decrease their respective action probabilities, the degree to which the self-agent will increase its action probability is reduced accordingly. The effect is similarly intuitive in the negative advantage case. We provide the pseudocode of ReMAPPO in Algorithm 1.

4.2 ReFACMAC and ReMADDPG

Off-policy gradient algorithms in MARL typically estimate the Deterministic Policy Gradient (DPG) over joint actions. The actor losses for MADDPG and FACMAC take the following respective forms:

$$\mathcal{L}_i^{\text{MADDPG}} = \mathbb{E}_{s \sim d^\beta} [Q_i^{\mu_i}(s, \boldsymbol{\mu}(s))], \quad (6)$$

$$\mathcal{L}_i^{\text{FACMAC}} = \mathbb{E}_{s \sim d^\beta} [\mathcal{F}(Q_1^{\mu_1}(s, \boldsymbol{\mu}(s)), \dots, Q_i^{\mu_i}(s, \boldsymbol{\mu}(s)), \dots, Q_N^{\mu_N}(s, \boldsymbol{\mu}(s)))], \quad (7)$$

where β is an arbitrary behavior policy, $Q_i^{\mu_i}(s, \boldsymbol{\mu}(s))$ is the action-value function for agent i over joint actions, and \mathcal{F} is a learned mixing function which factorizes a joint action-value function. By incorporating the updates of non-self agents, we reformulate these losses to define the recursive algorithms ReMADDPG and ReFACMAC:

$$\mathcal{L}_i^{\text{ReMADDPG}} = \mathbb{E}_{s \sim d^\beta} [Q_i^{\mu_i}(s, \mu_i(s), \boldsymbol{\mu}'_{-i}(s))], \quad (8)$$

$$\mathcal{L}_i^{\text{ReFACMAC}} = \mathbb{E}_{s \sim d^\beta} [\mathcal{F}(\dots, Q_i^{\mu_i}(s, \mu_i(s), \boldsymbol{\mu}'_{-i}(s)), \dots)]. \quad (9)$$

Note that while correctly estimating the DPG requires an unbiased action-value estimate, off-policy algorithms such as DDPG (Lillicrap et al., 2015) introduce bias in practice via stabilization tricks such as bootstrapping with target networks. Centralized multi-agent algorithms introduce further bias by estimating joint action value functions with off-policy target bootstrapping of non-self agents (Liu et al., 2022; Lowe et al., 2017). In other words, the centralized and factored critics used in MADDPG and FACMAC are trained on actions taken by the behavior policies of the agents, but are used to estimate the joint action-value function given

Algorithm 1 ReMAPPO

Input: Initial actor parameters θ_i , actor policies π_i , learning rates η_i , $\forall i \in \mathcal{I}$, optimizer Ω , advantage estimator $\text{GAE}()$,

for num_update_steps **do**

$\mathcal{D} \leftarrow \emptyset$

for num_rollout_steps **do**

$\mathbf{a} \sim \pi(\cdot|s)$ ▷ sample joint actions

$s' \sim \mathcal{P}(s, \mathbf{a}, s')$

$\mathcal{D} \leftarrow \mathcal{D} \cup (s, \mathbf{a})$ ▷ collect batch of samples

end for

for $i \in \mathcal{I}$ **do** ▷ initial update

$A_i^\pi(s, \mathbf{a}) \leftarrow \text{GAE}(s, \mathbf{a}), \forall (s, \mathbf{a}) \in \mathcal{D}$

$r_i(s, \mathbf{a}) \leftarrow \frac{\pi'_i(a_i|s)}{\pi_i(a_i|s)}$

$\mathcal{L}_i^{\text{MAPPO}} \leftarrow \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi} [\min(r_i(s, \mathbf{a}), \text{clip}(r_i(s, \mathbf{a}), 1 - \epsilon, 1 + \epsilon)) \cdot A_i^\pi(s, \mathbf{a})]$

$\theta'_i \leftarrow \theta_i + \eta_i \Omega(\mathcal{L}_i^{\text{MAPPO}})$

$\mathbf{a} \sim \pi'(\cdot|s) \forall s \in \mathcal{S}$ ▷ sample updated actions

end for

for $i \in \mathcal{I}$ **do** ▷ recursive update

$r'_i(s, \mathbf{a}) = \frac{\pi'_i(a_i|s)}{\pi_i(a_i|s)} \cdot \frac{\pi'_{-i}(\mathbf{a}_{-i}|s)}{\pi_{-i}(\mathbf{a}_{-i}|s)}$

$\mathcal{L}_i^{\text{ReMAPPO}} \leftarrow \mathbb{E}_{s \sim d^\pi, \mathbf{a} \sim \pi} [\min(r'_i(s, \mathbf{a}), \text{clip}(r'_i(s, \mathbf{a}), 1 - \epsilon, 1 + \epsilon)) \cdot A_i^\pi(s, \mathbf{a})]$

$\theta''_i \leftarrow \theta_i + \eta_i \Omega(\mathcal{L}_i^{\text{ReMAPPO}})$

$\theta_i \leftarrow \theta''_i$ ▷ assign parameters for next update step

end for

end for

arbitrary actions from the non-self agents during the actor update. By using recursive non-self actions for each agent’s joint action-value function, we also estimate the value of joint actions that are arbitrarily distinct from those present during data collection (with the added benefit of mutual consistency with the updates of other agents). Thus, we maintain that ReMADDPG and ReFACMAC introduce no additional bias over their non-recursive counterparts.

5 Experiments

In this section, we demonstrate the effectiveness of recursive reasoning algorithms across three challenging benchmarks: StarCraft II in JaxMARL (SMAX) (Rutherford et al., 2023), the StarCraft II Multi-Agent Challenge (SMAC) (Samvelyan et al., 2019), and Multi-Agent Multi-Joint dynamics with Contact (MAMuJoCo) (Peng et al., 2021). Note that SMAX and SMAC have different dynamics and are considered separate benchmarks. We benchmark 11 SMAX maps (8 SMACv1-based and 3 SMACv2-based), 8 SMAC maps, and 4 MAMuJoCo maps. These environments present a diverse set of cooperative tasks on which we can compare the sample efficiency and performance of our methods against competitive MARL algorithms.

5.1 SMAX

On SMAX, we compare ReMAPPO against Independent Proximal Policy Optimization (IPPO) (Schulman et al., 2017) and Multi-Agent Proximal Policy Optimization (MAPPO) as on-policy gradient-based methods. We also benchmark Independent Q-Learning (Watkins & Dayan, 1992) (IQL), Value Decomposition Networks (VDN) (Sunehag et al., 2017), and QMIX (Rashid et al., 2020b). Finally, we benchmark POLA (Zhao et al., 2022) as it is the most well-known MARL algorithm which incorporates higher-order opponent updates (in particular, we implement Outer POLA with advantage estimation). We do not benchmark FACMAC and MADDPG in SMAX due to a lack of existing well-tuned implementations in JAX-based environments. A lack of satisfactory results across hyperparameter sweeps (both nominal and those of similar algorithms tuned on SMAX) leads us to believe a comparison in CPU-based environments is more fair.

These baseline are run with the settings seen in Rutherford et al. (2023): each is trained for $1e7$ total training steps against the ‘HeuristicEnemySMAX’ AI, updated every 128 steps, and uses a γ of 0.99. **Off-policy algorithms** (QMIX, VDN, IQL) are trained with 16 parallel environments with a buffer size of 5000 and a batch size of 32. Each uses Adam optimizers with a learning rate of $5e - 5$, and performs ϵ -greedy exploration during training time with and ϵ that decays from 1 to 0.05 over the first 10% of total steps (learning is also paused until the ϵ decay is concluded). Neural networks use a hidden size of 512 and relu activations, hard target updates every 10 updates, 8 update epochs, and a reward scale of 10 (the reward scale of the original SMAC environments). The maximum gradient norm is constrained to be 10. In QMIX, the mixer embedding dimension is 64, the mixer hypernet hidden dimension is 256, and the initial scale of the kernel weights of the mixer weights is set to 0.001. Baselines are evaluated every 5% of total steps for 128 steps across 512 environments. **On-policy algorithms** (ReMAPPO, MAPPO, IPPO, POLA) are trained with 64 parallel environments. Each uses Adam optimizers with a learning rate of $4e - 3$ which is annealed to 0 over the entire course of training. Neural networks use a hidden size of 128 and relu activations, 2 minibatch updates, and 2 update epochs. The maximum gradient norm is constrained to be 0.5. The value of λ for the GAE is set to 0.95, the value of ϵ for surrogate clipping is 0.2, the value loss coefficient is 0.5, and no entropy bonus is provided.

5.2 SMAC and MAMuJoCo

In SMAC and MAMuJoCo, we compare ReFACMAC against FACMAC, QMIX/COMIX, MAPPO, MADDPG, and POLA. We choose FACMAC as the algorithm with which to demonstrate the benefits of recursion as FACMAC achieves SOTA performance in these environments. We choose not to benchmark ReMAPPO in these environments, as the low parallelizability and low-data regime of CPU-based environments is not as conducive to on-policy algorithms (as seen by MAPPO’s performance in these benchmarks). We therefore consider it more fair to benchmark ReMAPPO in SMAX as a fair comparison of sample efficiency. Results for ReMADDPG are present in Section 6.3. For each algorithm, we evaluate the performance by pausing training after 10,000 steps and running a fixed number of independent test episodes (10 for MAMuJoCo and 32 for SMAC). During these test episodes, each agent acts greedily in a decentralized fashion. The mean performance of the agents is reported in MAMuJoCo (the performance for each agent is identical since they share a common objective) and the mean success rate is reported for SMAC. Note that our results sometimes appear different to previous works that benchmark SMAC and MAMuJoCo (Rashid et al., 2020b;a) since these works report median win rates. We choose to report mean win rates as this highlights the impact of seeds which fail completely on particularly difficult maps such as Corridor, which are ignored when the median is reported. Since our results are primarily obtained by applying recursive reasoning to FACMAC, we mostly kept the algorithmic implementation standards used in FACMAC for reproducibility. We use parameter sharing for all actor and critic networks to speed up learning. All actor, critic, mixer, and Q-networks have target networks. We set $\gamma = 0.99$ for all experiments. Further benchmark-specific training details are as follows:

MAMuJoCo All MAMuJoCo environments and agents are configured according to the default configurations used in Peng et al. (2021) where they were introduced. Each agent observes its own joint positions, control its own joints, and receives a common team reward. The exact configurations and rewards can be seen at <https://robotics.farama.org/envs/MaMuJoCo/>. The architecture of all deep Q-networks is an MLP with 2 hidden layers with 400 and 300 units respectively. In all actor-critic methods, the architecture of the shared actor and critic networks is an MLP with 2 hidden layers with 400 and 300 hidden units respectively. All hidden layers for all networks use ReLU activations. All critic networks provide raw outputs while actor networks have a tanh activation at the output. Actor networks and DQNs receive the local observations of that agent as an input, appended with a one-hot vector due to the parameter sharing. All centralized critics and mixing networks are conditioned on the global state provided by the environment.

Each episode has a maximum length of 1000 steps. The total training time for each algorithm is set to 2 million steps. To improve initial exploration, each agent takes 10,000 random steps at the beginning of each run. During training we apply uncorrelated, mean-zero noise with a standard deviation of 0.1 to further encourage exploration. Each agent has a replay buffer with a maximum size of 1 million and trains with a

batch size of 100 after every new sample. Target networks are updated using Polyak averaging with $\tau = 0.001$. All neural networks are trained using the Adam optimizer (Kingma, 2014) with a learning rate of 0.001.

SMAC All experiments using SMAC used the default team configurations, rewards, and observations as the SMAC benchmark (Samvelyan et al., 2019). The state space includes the last actions of each agent (an inbuilt feature in StarCraft II) as this was found to stabilize learning for all algorithms. The architecture of all shared deep Q-networks is a DRQN with a recurrent layer comprised of a GRU with a 64-dimensional hidden state, with fully connected layers on either side. In all actor-critic methods, the architecture of all shared actors is a recurrent MLP comprised of a GRU with a 64-dimensional hidden state, with fully connected layers on either side. We train the GRU networks on batches of 32 fully unrolled episodes (with 0-padding to account for temporal mismatch between episodes). The architecture of all shared critic networks is an MLP with 2 hidden layers with 64 units. All networks use ReLU activations for the hidden layers. All actor critic methods select discrete actions using the Gumbel-Softmax estimator (De Boer et al., 2005) in order to turn continuous softmaxed logits into discrete one-hot actions while retaining the ability to backpropagate through the network.

Actor networks and DRQNs receive the local observations of that agent as an input, appended with the last action taken by the agent, as well as a one-hot vector due to the parameter sharing. All agents use ϵ -greedy action selection and we anneal ϵ from 0.5 to 0.05 over 50k training steps. The replay buffer contains the most recent 5000 episodes. All target networks are updated hard every 200 training steps. All networks are trained using Adam with a learning rate of 0.0025 for the actor network and 0.0005 for the critic network (except for QMIX which uses the learning rates specified in Samvelyan et al. (2019) as they have already been tuned for SMAC).

5.3 Experimental Results

ReMAPPO outperforms baselines in SMAX Figure 2 compares the test win rate of ReMAPPO against baselines on SMAX. Notably, ReMAPPO performs equal or better than other baselines on 9 out of 11 maps, and is the only algorithm to achieve a 100% win rate in 3s5z and 27m_vs_30m. ReMAPPO also maintains a consistent advantage over MAPPO, demonstrating the viability of recursive reasoning. The effects of recursive reasoning are seen in higher overall performance as well as faster convergence, as ReMAPPO often reaches its maximum performance with fewer samples than other methods. Despite utilizing opponent shaping, POLA fails to match the performance of ReMAPPO in every map except 3s5z_vs_3s6z. This suggests that POLA’s formulation for 2-player reciprocity-based games likely does not generalize well to more complex environments with many agents.

ReFACMAC outperforms baselines in SMAC and MAMuJoCo Figure 3 compares the test win rate of ReFACMAC and related baselines in SMAC with a focus on Hard and Super Hard maps. ReFACMAC achieves higher or equal final success rates compared to baselines and typically converges to its maximum win rate with fewer samples. ReFACMAC particularly stands out in Corridor and 3s5z_vs_3s6z, two notoriously difficult maps in which it is the only algorithm to surpass a 50% success rate. Note that MAPPO performs much worse on SMAC maps due to SMAC being CPU-based and less parallelizable, resulting in far less sample availability (it was only possible to obtain 2e6 samples for each map rather than 1e7 in SMAX).

Figure 4 compares the performance of ReFACMAC against baselines on four selected MAMuJoCo environments. In all four environments, ReFACMAC achieves superior performance by the end of training and tends to reach its peak performance earlier. In Ant 2x4 (Figure 4c), learning a solution is difficult due to the asymmetric positioning of the agents which control opposing diagonal halves of a 4-legged agent. ReFACMAC and FACMAC both achieve the same maximum performance during training, but only ReFACMAC maintains an advantage despite the instability of the problem. ReFACMAC also achieves the *only* policy in Walker 2x3 (Figure 4d) which solves the task, obtaining SOTA performance on this very difficult control benchmark.

We note the wall-clock times of our main experiments in Tables 1 and 2. While ReMAPPO and ReFACMAC take longer to run than MAPPO and FACMAC respectively, we note that the relative increase is not extreme due to implementation optimizations such as saved optimizer states, JIT compilation, environment parallelization, and GPU tensor processing. Please refer to the supplementary material to see how the code

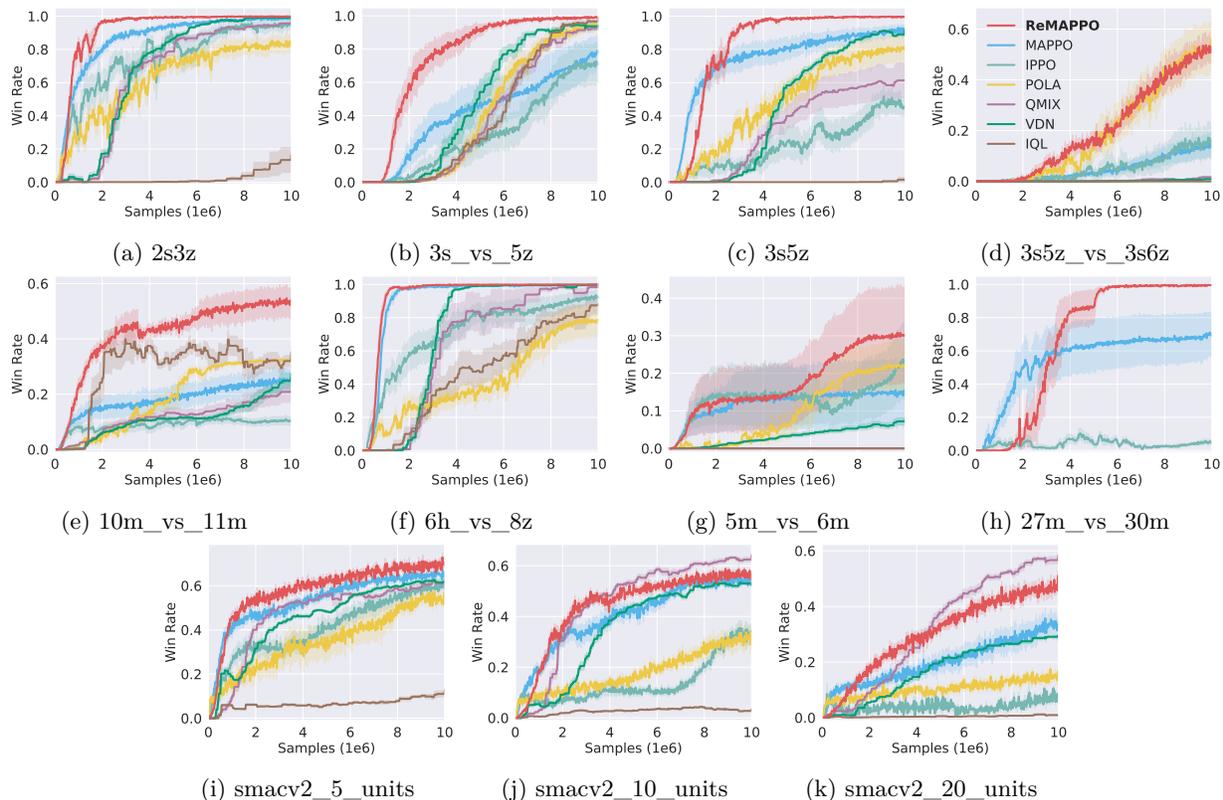


Figure 2: Mean win rate and standard error of ReMAPPO and baselines on SMAX maps across 10 seeds. Note that 27m_vs_30m is very large and could only be benchmarked with PPO-based methods due to computational constraints.

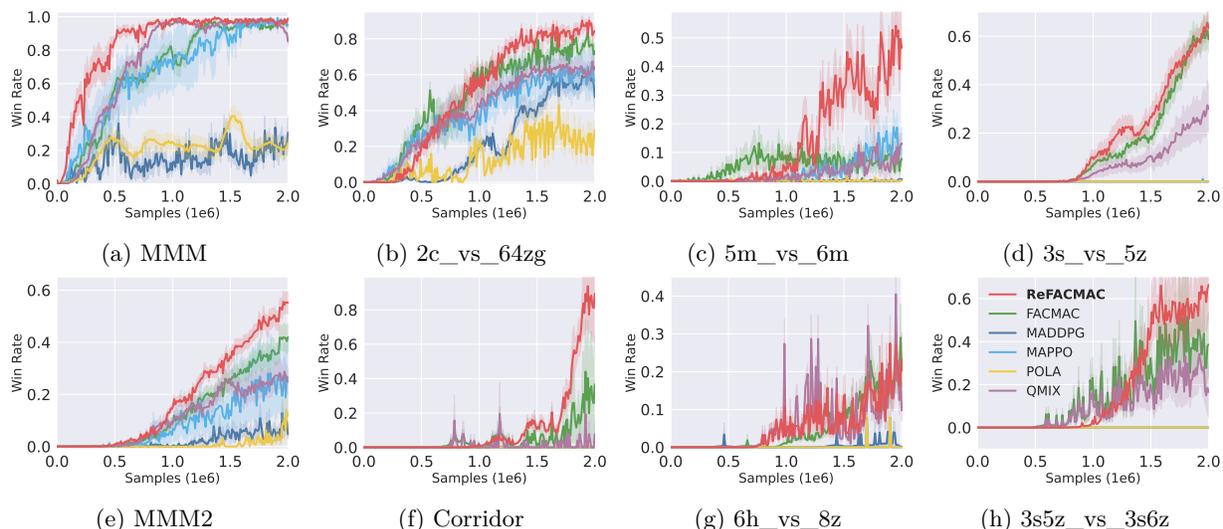


Figure 3: Mean win rate and standard error of ReFACMAC and baselines on SMAC maps across 10 seeds.

for this project implements these optimizations. Additionally, note the extreme disparity between on-policy (MAPPO, IPPO, ReMAPPO, POLA) methods and off-policy (IQL, QMIX, VDN) methods in SMAX - in the worst case, QMIX takes over 8 times as long as MAPPO to finish running a single experiment on average. This is because of the bottleneck of high VRAM usage and frequent updates of replay buffer-based reinforcement

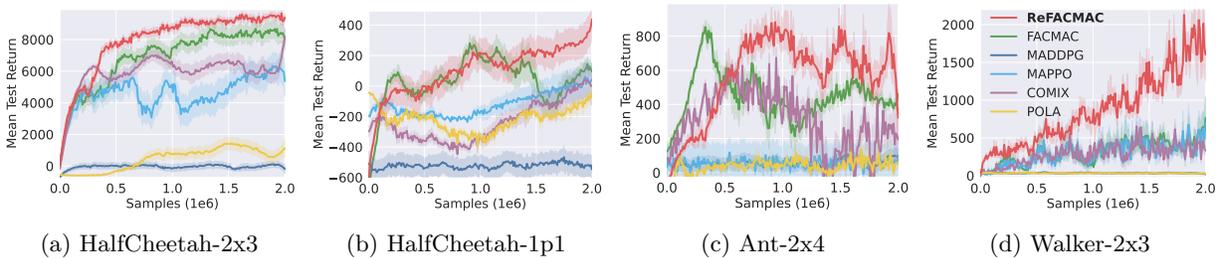


Figure 4: Mean performance and standard error of ReFACMAC and baselines on four Multi-Agent MuJoCo environments across 10 seeds.

Algorithm	SMAX
MAPPO	+0.0 ± 12%
IPPO	-24 ± 15%
ReMAPPO	+11 ± 12%
POLA	+36 ± 24%
IQL	+273 ± 10%
QMIX	+876 ± 9.3%
VDN	+243 ± 4.4%

Table 1: Mean and standard deviation of experimental wall-clock times relative to MAPPO (SMAX)

Algorithm	SMAC	MaMuJoCo
FACMAC	+0.0 ± 12%	+0.0 ± 5.6%
ReFACMAC	+13 ± 11%	+15 ± 13%
MADDPG	-13 ± 12%	-9 ± 9.9%
MAPPO	-24 ± 4.0%	-21 ± 3.9%
POLA	+21 ± 7.6%	+19 ± 6.3%
QMIX/COMIX	+12 ± 8.0%	+16 ± 8.0%

Table 2: Mean and standard deviation of experimental wall-clock times relative to FACMAC (SMAC and MaMuJoCo)

learning algorithms, which mitigate the benefits of JIT compilation and environment parallelization relative to on-policy methods.

6 Higher Recursive Reasoning

Thus far, we have considered practical implementations of recursive reasoning with policy gradient algorithms. While the results for ReMAPPO and ReFACMAC generally exhibit competitive sample-efficiency with just one level of recursive reasoning, our framework allows for further levels of recursion. In the framework of higher recursive reasoning, we consider the application of Equation 2 an arbitrary number of times within one update step; that is, we repeatedly refine the policy gradient using the updates of the non-self agents at the previous recursion level. Note that at each recursion level, no new environment data is being collected - the policy gradient is being refined using the action distributions alone.

We represent the parameters of agent i after k repeated recursions by $\theta_i^{(k)}$. We define the pre-update parameters to be recursion level 0. Thus, the vanilla update parameters denoted θ_i in Equation 1 would be represented by $\theta_i^{(1)}$, and the recursive reasoning parameters θ_i^{Re} in Equation 2 would be represented by $\theta_i^{(2)}$. Figure 5 further illustrates our higher recursion framework. Note also that recursive gradient steps are taken from the starting point of the update, hence this is not equivalent to simply taking more learning steps.

6.1 Theoretical study

In this section, we present a theoretical study which shows that repeated recursive reasoning with policy gradients converges monotonically to a local Nash equilibrium under certain conditions with finite iterates. **Note that Assumption 6.1 typically only holds under strong conditions (Pirota et al., 2015).** Firstly, we demonstrate how an unbiased, infinite application of recursive reasoning leads to perfect anticipation of other agents’ future strategies.

Assumption 6.1. *The gradient $\nabla_{\theta_i} J_i(\theta_i, \theta_{-i})$ is L_i -Lipschitz with respect to θ_{-i} , i.e.*

$$\|\nabla_{\theta_i} J_i(\theta_i, \theta_{-i,1}) - \nabla_{\theta_i} J_i(\theta_i, \theta_{-i,2})\| \leq L_i \|\theta_{-i,1} - \theta_{-i,2}\|, \forall \theta_i, \theta_{-i}, \quad (10)$$

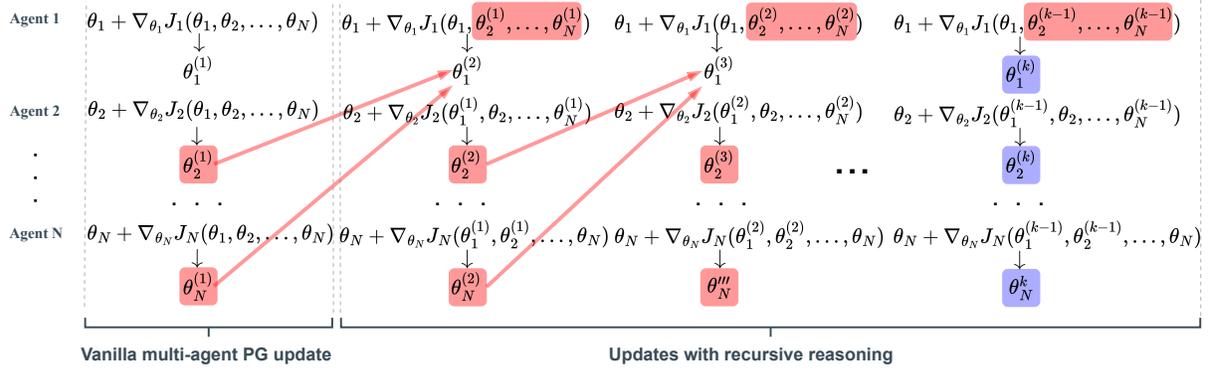


Figure 5: Higher recursive reasoning policy gradient updates for N agents. The non-self parameters used for recursive updates of Agent 1 are red and the final parameters of all agents after k steps of recursion are blue.

where $\theta_{-i,1}$ and $\theta_{-i,2}$ are two arbitrary points in the joint parameter space of the non-self agents. We define the maximum objective function Lipschitzness $L := \max_i \{L_i\}$ and the maximum agent learning rate $\eta := \max_i \{\eta_i\}$. We define the maximum objective function gradient across all agent parameters and objective functions $\nabla_{max} := \max_{i, \theta_i, \theta_{-i}} \|\nabla_{\theta_i} J_i(\theta_i, \theta_{-i})\|$.

Theorem 6.2. *Suppose Assumption 6.1 holds. Then, the update step at the k 'th level of reasoning is bounded:*

$$\|\theta^{(k)} - \theta^{(k-1)}\| \leq \eta(\eta L)^{k-1} N(N-1)^{k-1} \nabla_{max}. \quad (11)$$

Assume the maximum learning rate η satisfies $\eta < \frac{1}{L(N-1)}$. Then, the sequence $\{\theta^{(k)}\}_{k=0}^{\infty}$ is a convergent sequence. Since θ exists in a complete subspace of $\mathbb{R}^{\sum_i d_i}$, the convergent sequence $\{\theta^{(k)}\}_{k=0}^{\infty}$ is Cauchy, i.e.,

$$\exists C \in \mathbb{N} : \forall \epsilon > 0, (a > b > C \implies \|\theta^{(a)} - \theta^{(b)}\| < \epsilon). \quad (12)$$

Since every Cauchy sequence has a limit, we denote the limit of $\{\theta^{(k)}\}_{k=0}^{\infty}$ as $\lim_{k \rightarrow \infty} \theta^{(k)} = \theta^{(\infty)}$.

According to Theorem 6.2, applying the recursive update with $k=\infty$ defines the following implicit algorithm:

$$\theta_i^{(\infty)} \leftarrow \theta_i + \eta_i \nabla_{\theta_i} J_i(\theta_i, \theta_{-i}^{(\infty)}), \forall i \in \mathcal{I}, \quad (13)$$

which we denote the Generalized Semi-Proximal Point Method (GSPPM). The implication of the GSPPM is that the update of each agent responds exactly to the updated strategies of the other agents, maintaining mutual consistency.

Following from Theorem 6.2, we show that the convergence of GSPPM iterates in a non-convex non-concave strategy space can be analyzed via the game Jacobian around a local stationary point:

Theorem 6.3. *Let θ^* be a stationary point in an N -player general sum game. This stationary point is a local Nash equilibrium, i.e. a point at which no agent's objective function has a non-zero gradient under a unilateral change in policy. Let η be a block matrix of the agent learning rates η_i . Let the components of the Hessian of each objective function at θ^* be denoted*

$$\begin{pmatrix} \nabla_{\theta_i \theta_i}^2 J_i(\theta_i^* \theta_{-i}^*) & \nabla_{\theta_i \theta_{-i}}^2 J_i(\theta_i^* \theta_{-i}^*) \\ \nabla_{\theta_{-i} \theta_i}^2 J_i(\theta_i^* \theta_{-i}^*) & \nabla_{\theta_{-i} \theta_{-i}}^2 J_i(\theta_i^* \theta_{-i}^*) \end{pmatrix} = \begin{pmatrix} A_i & B_i \\ B_i^\top & C_i \end{pmatrix}.$$

Furthermore, let \mathbf{A} be the diagonal block matrix of all \mathbf{A}_i matrices and \mathbf{B} be the diagonal block matrix of all \mathbf{B}_i matrices: $\mathbf{A} = \text{diag}(\mathbf{A}_1, \dots, \mathbf{A}_n)$, $\mathbf{B} = \text{diag}(\mathbf{B}_1, \dots, \mathbf{B}_n)$. Let \mathbf{D} be a complement-selection matrix for each set of agent parameters θ_i such that $\mathbf{D}\theta = [\theta_{-1}, \dots, \theta_{-n}]^\top$

Suppose $\eta < \frac{1}{L(N-1)}$ such that the GSPPM iterates $\{\theta_t^{(k)}\}_{k=0}^\infty$ form a Cauchy sequence. Let $\lambda_{\max}(\mathbf{Z})$ and $\lambda_{\min}(\mathbf{Z})$ refer to the maximum and minimum eigenvalues of a matrix \mathbf{Z} respectively. Then, there exists a neighborhood $\mathcal{U} \in \mathbb{R}^{\sum_i d_i}$ around θ^* such that if GSPPM starts in \mathcal{U} , the iterates $\{\theta_t^{(k)}\}_{k=0}^\infty$ satisfy:

$$\|\theta^{(\infty)} - \theta^*\| \leq \frac{\lambda_{\max}(I + \eta\mathbf{A})^2}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})^2} \|\theta - \theta^*\|. \quad (14)$$

Moreover, for any η satisfying $\frac{\lambda_{\max}(I + \eta\mathbf{A})^2}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})^2} < 1$, the iterates converge asymptotically to θ^* . Hence, GSPPM iterates reach an ϵ -Nash equilibrium.

Theorem 6.4. Suppose the conditions of Theorem 6.3 apply. Then, the finite iterates $\{\theta_t^{(k)}\}$ satisfy:

$$\|\theta^{(k)} - \theta^*\|^2 \leq \left(\frac{\lambda_{\max}(I + \eta\mathbf{A})^2}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})^2} + 2 \frac{\lambda_{\max}(I + \eta\mathbf{A})}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})} \frac{\lambda_{\max}(\eta\mathbf{B}\mathbf{D})}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})} \right) \|\theta - \theta^*\|^2 \quad (15)$$

$$+ 2 \frac{\lambda_{\max}(I + \eta\mathbf{A})}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})} \frac{\lambda_{\max}(\eta\mathbf{B}\mathbf{D})}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})} (\|\theta - \theta^*\| \nabla_{\max}) \quad (16)$$

$$+ \frac{\lambda_{\max}(\eta\mathbf{B}\mathbf{D})^2}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})^2} \|\theta^{(k-1)} - \theta^*\|^2. \quad (17)$$

Moreover, for any η satisfying $\frac{\lambda_{\max}(\eta\mathbf{B}\mathbf{D})^2}{\lambda_{\min}(I - \eta\mathbf{B}\mathbf{D})^2} < 1$, the finite iterates converge asymptotically to θ^* .

6.2 Illustrative example

We further examine higher recursive reasoning using a didactic example of a simple cooperative game with two point agents taking continuous actions in a 2D space. The agents have one parameter each and produce a one-dimensional action (the angle of their next move). The highest reward is achieved when the chosen direction of each agent points towards the *future* location of the other agent. Assuming the agents move kinematically, the optimal solution is for the agents to move towards each other in a straight line (see Appendix B for details). However, the naive policy update for this game under fictitious play is for each agent to choose its next action to intercept with the *previous* action of the other agent. The top two figures of Figure 7a illustrate this problem: each agent's new action (red arrows) points towards the destination of the other agent under the other agent's old policy (yellow arrows). Hence, the naive update leads to a lack of mutual consistency.

Figure 7a bottom left shows the policy update after 2 levels of recursion: now the agents update to intercept the other agent *after* the other agent's naive policy update, resulting in better coordination. As the number of recursions increases, the policies converge on an ϵ -bound of the optimal solution. Figure 7b shows the progression of the agent parameters with gradient ascent and momentum; we use the true gradient and objective function $J(\theta_1, \theta_2)$. The benefits of recursive reasoning are evinced by the fact that increasing k -levels (darker points) exhibit monotonic convergence to the optimal parameters θ_1^*, θ_2^* at each update step, as supported by Theorem 6.4. Figure 6 shows the distance from the Nash equilibrium for one update step near the stationary point in Figure 7b.

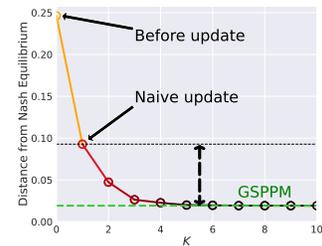
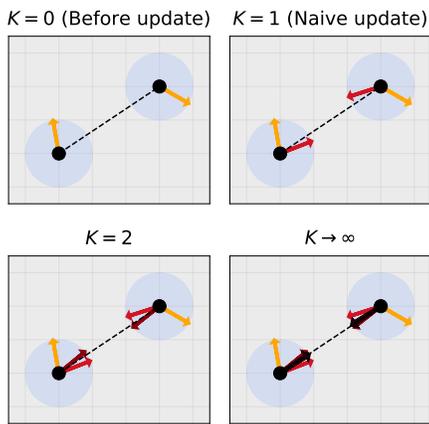
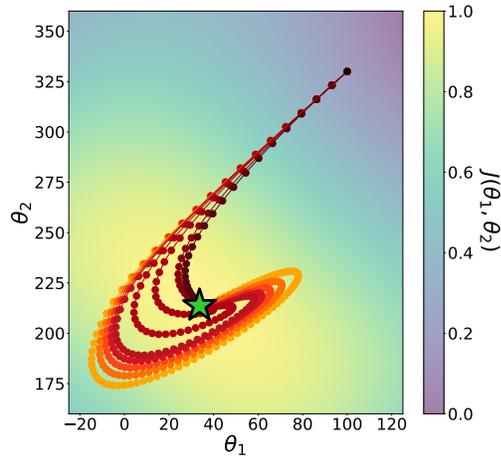


Figure 6: Convergence to GSPPM.



(a) Agents taking continuous actions with increasing recursive reasoning levels (darkening arrows) converging on the optimal actions (dashed line).



(b) Gradient ascent to the optimal parameters (green star) with recursion (darker colors for higher recursive reasoning levels).

Figure 7: An illustrative continuous cooperative game with two point agents using recursive reasoning.

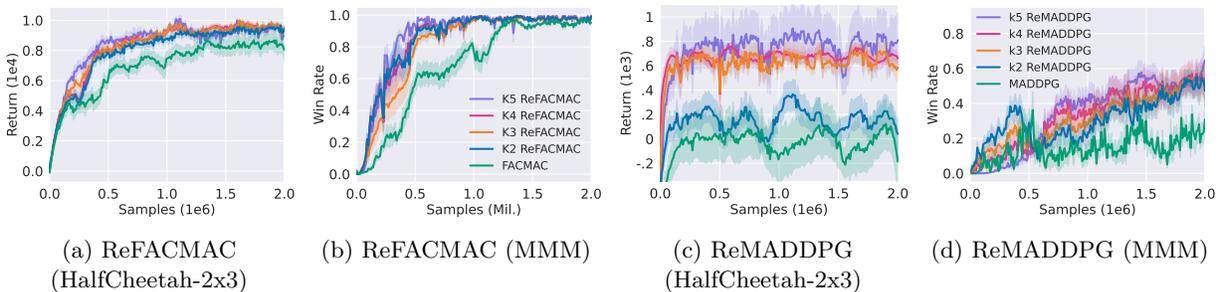


Figure 8: Ablations for higher recursive reasoning with ReFACMAC and ReMADDPG on SMAC (MMM) and MAMuJoCo (HalfCheetah-2x3).

6.3 Higher recursive reasoning in deep multi-agent reinforcement learning

Consolidating with the theoretical study in Section 6.1, we compare the sample efficiency of higher recursive reasoning in ReFACMAC and ReMADDPG up to $k=5$ in HalfCheetah-2x3 (MAMuJoCo) and MMM (SMAC) in Figure 8. While higher recursion exhibits stable performance increases, most of the benefits appear to materialize by $k=2$ in both the ablation and didactic example in Section 6.2. Similar lookahead methods in optimizers (Suh & Ma, 2025) and GANs (Liu & Pavel, 2022) exhibit polynomial convergence rates, achieving the majority of benefits within the first few iterates. We leave an analysis of the convergence rates of higher recursive reasoning with finite data to future investigations. Interestingly, Figure 8c demonstrates a situation where recursions above $k=2$ find a better policy mode.

7 Conclusion

We presented a framework for recursive reasoning for multi-agent policy gradient algorithms. We introduce our recursive policy gradient formulation and realize it in both on and off-policy regimes, achieving SOTA performance against competitive baselines in challenging MARL benchmarks. We prove theoretical convergence properties of finite and infinite recursive policy gradient iterates with respect to local equilibria. We leave it to future work to understand the convergence rates of recursive iterates in the finite data setting and to address the assumption of access to non-self agent policy distributions.

References

- David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel. The mechanics of n-player differentiable games. In *International Conference on Machine Learning*, pp. 354–363. PMLR, 2018.
- W. Barfuss and R. Mann. Modeling the effects of environmental and perceptual uncertainty using deterministic reinforcement learning dynamics with partial observability. *Physical review. E*, 105 3-1:034409, 2021. doi: 10.1103/PhysRevE.105.034409.
- Colin F Camerer, Teck-Hua Ho, and Juin-Kuan Chong. A cognitive hierarchy model of games. *The Quarterly Journal of Economics*, 119(3):861–898, 2004.
- Brandon Cui, Hengyuan Hu, Luis Pineda, and Jakob Foerster. K-level reasoning for zero-shot coordination in hanabi. *Advances in Neural Information Processing Systems*, 34:8215–8228, 2021.
- Siyu Dai, Sangjae Bae, and David Isele. Game theoretic decision making by actively learning human intentions applied on autonomous driving. *arXiv preprint arXiv:2301.09178*, 2023.
- Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. A tutorial on the cross-entropy method. *Annals of operations research*, 134:19–67, 2005.
- Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. Liir: Learning individual intrinsic reward in multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Benjamin Patrick Evans and Mikhail Prokopenko. Bounded rationality for relaxing best response and mutual consistency: The quantal hierarchy model of decision-making. *arXiv preprint arXiv:2106.15844*, 2021.
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. *arXiv preprint arXiv:1709.04326*, 2017.
- Johannes Forkel and Jakob Foerster. Entropy is all you need for inter-seed cross-play in hanabi. *arXiv preprint arXiv:2511.22581*, 2025.
- Dean P Foster and H Peyton Young. On the nonconvergence of fictitious play in coordination games. *Games and Economic Behavior*, 25(1):79–96, 1998.
- Piotr J Gmytrasiewicz, Edmund H Durfee, and David K Wehe. A decision-theoretic approach to coordinating multi-agent interactions. In *IJCAI*, volume 91, pp. 63–68, 1991.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11): 139–144, 2020.
- Ammar Haydari and Yasin Yilmaz. Deep reinforcement learning for intelligent transportation systems: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(1):11–32, 2020.
- Zool Hilmi Ismail, Nohaidda Sariff, and E Gorrostieta Hurtado. A survey and analysis of cooperative multi-agent robot systems: challenges and directions. *Applications of Mobile Robots*, 5:8–14, 2018.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the nineteenth international conference on machine learning*, pp. 267–274, 2002.
- Shahab Karimi, Arash Karimi, and Ardalan Vahidi. Level- k reasoning, deep reinforcement learning, and monte carlo decision process for fast and safe automated lane change and speed management. *IEEE Transactions on Intelligent Vehicles*, 8(6):3556–3571, 2023.

- Akbar Khan, Timon Willi, Newton Kwan, Andrea Tacchetti, Chris Lu, Edward Grefenstette, Tim Rocktäschel, and Jakob Foerster. Scaling opponent shaping to high dimensional games. *arXiv preprint arXiv:2312.12568*, 2023.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint arXiv:2109.11251*, 2021.
- Hui Li, X. Liao, and L. Carin. Multi-task reinforcement learning in partially observable stochastic environments. *J. Mach. Learn. Res.*, 10:1131–1186, 2009. doi: 10.5555/1577069.1577109.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pp. 157–163. Elsevier, 1994.
- Bing Liu, Yuxuan Xie, Lei Feng, and Ping Fu. Correcting biased value estimation in mixing value-based multi-agent reinforcement learning by multiple choice learning. *Engineering Applications of Artificial Intelligence*, 116:105329, 2022.
- Zichu Liu and Lacra Pavel. Recursive reasoning in minimax games: A level k gradient play method. *Advances in Neural Information Processing Systems*, 35:16903–16917, 2022.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- Xiaobai Ma, David Isele, Jayesh K Gupta, Kikuo Fujimura, and Mykel J Kochenderfer. Recursive reasoning graph for multi-agent reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7664–7671, 2022.
- 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.
- Dov Monderer and Lloyd S Shapley. Fictitious play property for games with identical interests. *Journal of economic theory*, 68(1):258–265, 1996.
- Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Böhmer, and Shimon Whiteson. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems*, 34:12208–12221, 2021.
- Yue Pi, Wang Zhang, Yong Zhang, Hairong Huang, Baoquan Rao, Yulong Ding, and Shuanghua Yang. Applications of multi-agent deep reinforcement learning communication in network management: A survey. *arXiv preprint arXiv:2407.17030*, 2024.
- Matteo Pirodda, Marcello Restelli, and Luca Bascetta. Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2):255–283, 2015.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:10199–10210, 2020a.

- 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. *Journal of Machine Learning Research*, 21(178):1–51, 2020b.
- D Robertson. General theory of employment, interest and money. *QJ Econ*, 51:791–795, 1936.
- Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Gardar Ingvarsson, Timon Willi, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, et al. Jaxmarl: Multi-agent rl environments in jax. *arXiv preprint arXiv:2311.10090*, 2023.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *International conference on machine learning*, pp. 387–395. Pmlr, 2014.
- Jaewook J Suh and Shiqian Ma. An adaptive and parameter-free nesterov’s accelerated gradient method for convex optimization. *arXiv preprint arXiv:2505.11670*, 2025.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- Chen Tang, Ben Abbatematteo, Jiaheng Hu, Rohan Chandra, Roberto Martín-Martín, and Peter Stone. Deep reinforcement learning for robotics: A survey of real-world successes. *Annual Review of Control, Robotics, and Autonomous Systems*, 8, 2024.
- Xihuai Wang, Zheng Tian, Ziyu Wan, Ying Wen, Jun Wang, and Weinan Zhang. Order matters: Agent-by-agent policy optimization. *arXiv preprint arXiv:2302.06205*, 2023.
- Xinpeng Wang, Songan Zhang, and Huei Peng. Comprehensive safety evaluation of highly automated vehicles at the roundabout scenario. *IEEE Transactions on Intelligent Transportation Systems*, 23(11):20873–20888, 2022.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3):279–292, 1992.
- Ying Wen, Yaodong Yang, Rui Luo, and Jun Wang. Modelling bounded rationality in multi-agent interactions by generalized recursive reasoning. *arXiv preprint arXiv:1901.09216*, 2019a.
- Ying Wen, Yaodong Yang, Rui Luo, Jun Wang, and Wei Pan. Probabilistic recursive reasoning for multi-agent reinforcement learning. *arXiv preprint arXiv:1901.09207*, 2019b.
- Timon Willi, Alistair Hp Letcher, Johannes Treutlein, and Jakob Foerster. Cola: consistent learning with opponent-learning awareness. In *International Conference on Machine Learning*, pp. 23804–23831. PMLR, 2022.
- Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information processing systems*, 35:24611–24624, 2022.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pp. 321–384, 2021.

Stephen Zhao, Chris Lu, Roger B Grosse, and Jakob Foerster. Proximal learning with opponent-learning awareness. *Advances in Neural Information Processing Systems*, 35:26324–26336, 2022.

Yulai Zhao, Zhuoran Yang, Zhaoran Wang, and Jason D Lee. Local optimization achieves global optimality in multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 42200–42226. PMLR, 2023.

A Proofs

A.1 Proof of Theorem 6.2

Recall Assumption 6.1. We analyze the pattern in successive updates of $\boldsymbol{\theta}$ as k increases.

Consider level $k=1$:

$$\boldsymbol{\theta}_i^{(1)} = \boldsymbol{\theta}_i + \eta_i \nabla_{\boldsymbol{\theta}_i} J(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}) \quad \forall i \in \mathcal{I}. \quad (18)$$

The jump between $\boldsymbol{\theta}_i^{(1)}$ and $\boldsymbol{\theta}_i$ is

$$\|\boldsymbol{\theta}_i^{(1)} - \boldsymbol{\theta}_i\| = \eta_i \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \quad \forall i \in \mathcal{I}. \quad (19)$$

Thus, for all agents

$$\begin{aligned} \|\boldsymbol{\theta}^{(1)} - \boldsymbol{\theta}\| &\leq \sum_{i=1}^N \|\boldsymbol{\theta}_i^{(1)} - \boldsymbol{\theta}_i\| = \sum_{i=1}^N \eta_i \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \\ &\leq \eta \sum_{i=1}^N \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \\ &\leq \eta N \max_i \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \\ &\leq \eta N \nabla_{\max}. \end{aligned} \quad (20)$$

Now consider level $k=2$:

$$\boldsymbol{\theta}_i^{(2)} = \boldsymbol{\theta}_i + \eta_i \nabla_{\boldsymbol{\theta}_i} J(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(1)}) \quad \forall i \in \mathcal{I}. \quad (21)$$

The jump between $\boldsymbol{\theta}_i^{(2)}$ and $\boldsymbol{\theta}_i^{(1)}$ is

$$\begin{aligned} \|\boldsymbol{\theta}_i^{(2)} - \boldsymbol{\theta}_i^{(1)}\| &= \|\boldsymbol{\theta}_i + \eta_i \nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(1)}) - \boldsymbol{\theta}_i - \eta_i \nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \\ &= \eta_i \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(1)}) - \nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i})\| \\ &\leq \eta_i L_i \|\boldsymbol{\theta}_{-i}^{(1)} - \boldsymbol{\theta}_{-i}\| \quad \forall i \in \mathcal{I}. \end{aligned} \quad (22)$$

Thus, for all agents

$$\begin{aligned}
\|\boldsymbol{\theta}^{(2)} - \boldsymbol{\theta}^{(1)}\| &\leq \sum_{i=1}^N \|\boldsymbol{\theta}_i^{(2)} - \boldsymbol{\theta}_i^{(1)}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \|\boldsymbol{\theta}_{-i}^{(1)} - \boldsymbol{\theta}_{-i}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \|\boldsymbol{\theta}_j^{(1)} - \boldsymbol{\theta}_j\| \\
&= \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \eta_j \|\nabla_{\boldsymbol{\theta}_j} J_j(\boldsymbol{\theta}_j, \boldsymbol{\theta}_{-j})\| \\
&\leq \sum_{i=1}^N \eta_i L_i (N-1) \eta \nabla_{max} \\
&\leq \eta^2 L N (N-1) \nabla_{max}.
\end{aligned} \tag{23}$$

Now consider level $k=3$:

$$\boldsymbol{\theta}_i^{(3)} = \boldsymbol{\theta}_i + \eta_i \nabla_{\boldsymbol{\theta}_i} J(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(2)}) \quad \forall i \in \mathcal{I}. \tag{24}$$

The jump between $\boldsymbol{\theta}_i^{(3)}$ and $\boldsymbol{\theta}^{(2)}$ is

$$\begin{aligned}
\|\boldsymbol{\theta}_i^{(3)} - \boldsymbol{\theta}_i^{(2)}\| &= \eta_i \|\nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(2)}) - \nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}^{(1)})\| \\
&\leq \eta_i L_i \|\boldsymbol{\theta}_{-i}^{(2)} - \boldsymbol{\theta}_{-i}^{(1)}\| \quad \forall i \in \mathcal{I}.
\end{aligned} \tag{25}$$

Thus, for all agents

$$\begin{aligned}
\|\boldsymbol{\theta}^{(3)} - \boldsymbol{\theta}^{(2)}\| &\leq \sum_{i=1}^N \|\boldsymbol{\theta}_i^{(3)} - \boldsymbol{\theta}_i^{(2)}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \|\boldsymbol{\theta}_i^{(2)} - \boldsymbol{\theta}_i^{(1)}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \|\boldsymbol{\theta}_j^{(2)} - \boldsymbol{\theta}_j^{(1)}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \eta_j L_j \|\boldsymbol{\theta}_{-j}^{(1)} - \boldsymbol{\theta}_{-j}\| \\
&\leq \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \eta_j L_j \sum_{l \neq j} \|\boldsymbol{\theta}_m^{(1)} - \boldsymbol{\theta}_m\| \\
&= \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \eta_j L_j \sum_{l \neq j} \eta_m \|\nabla_{\boldsymbol{\theta}_m} J_m(\boldsymbol{\theta}_m, \boldsymbol{\theta}_{-m})\| \\
&= \sum_{i=1}^N \eta_i L_i \sum_{j \neq i} \eta_j L_j \sum_{l \neq j} \eta_m \|\nabla_{\boldsymbol{\theta}_m} J_m(\boldsymbol{\theta}_m, \boldsymbol{\theta}_{-m})\| \\
&\leq \eta^3 L^2 N (N-1)^2 \nabla_{max}.
\end{aligned} \tag{26}$$

We see by induction that any consecutive states during the recursive procedure are bounded by

$$\|\boldsymbol{\theta}^k - \boldsymbol{\theta}^{k-1}\| \leq \eta(\eta L)^{k-1} N(N-1)^{k-1} \nabla_{max}. \quad (27)$$

Let $\eta < \frac{1}{L(n-1)}$ such that the difference between two recursive steps is a contraction. Consider the difference $\|\boldsymbol{\theta}^{(a)} - \boldsymbol{\theta}^{(b)}\|$, where $a > b > 0$:

$$\begin{aligned} \|\boldsymbol{\theta}^{(a)} - \boldsymbol{\theta}^{(b)}\| &= \left\| \sum_{j=b+1}^a \left(\boldsymbol{\theta}^{(j)} - \boldsymbol{\theta}^{(j-1)} \right) \right\| \\ &\leq \sum_{j=b+1}^a \|\boldsymbol{\theta}^{(j)} - \boldsymbol{\theta}^{(j-1)}\| \\ &\leq \sum_{j=b+1}^a \eta(\eta L)^{j-1} N(N-1)^{j-1} \nabla_{max} \\ &\leq N \nabla_{max} \left(\sum_{j=b+1}^a (N-1)^{j-1} L^{j-1} \right) \eta \left(\sum_{j=b+1}^a \eta^{j-1} \right) \\ &\leq N \nabla_{max} \left(\sum_{j=b+1}^a (N-1)^{j-1} L^{j-1} \right) \eta \left(\sum_{j=b+1}^a \eta^{b-1} \right) \approx \mathcal{O}(\eta^b). \end{aligned} \quad (28)$$

Thus, for any $\epsilon > 0$, we can solve for b such that $\eta(\eta L)^{k-1} N(N-1)^{k-1} \nabla_{max} < \epsilon$, or

$$\exists C \in \mathbb{N} : \forall \epsilon > 0, (a > b > C \implies \|\boldsymbol{\theta}^{(a)} - \boldsymbol{\theta}^{(b)}\| < \epsilon). \quad (29)$$

Hence $\{\boldsymbol{\theta}^{(k)}\}_{k=0}^{\infty}$ is a Cauchy sequence. Since $\boldsymbol{\theta}$ lies in a complete subspace of $\mathbb{R}^{\sum_i d_i}$, the Cauchy sequence has a limit: $\lim_{k \rightarrow \infty} \boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}^{\infty}$. \square

A.2 Proof of Theorem 6.3

We abuse notation by denoting agent i 's parameters at update step t with $\boldsymbol{\theta}_{t,i}$ and the parameters one update step later at $t+1$ with $\boldsymbol{\theta}_{t+1,i}$, where an arbitrary number of recursive reasoning steps were taken in between. Let us define $\hat{\boldsymbol{\theta}}_{t,i} = \boldsymbol{\theta}_{t,i} - \boldsymbol{\theta}_i^*$ and $\hat{\boldsymbol{\theta}}_t = [\hat{\boldsymbol{\theta}}_{t,1}, \dots, \hat{\boldsymbol{\theta}}_{t,N}]^T$ for all $i \in \mathcal{I}$. It follows by linearizing the system about the stationary point $\boldsymbol{\theta}^*$,

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{t+1,i} &= \boldsymbol{\theta}_{t,i} + \eta_i \nabla_{\boldsymbol{\theta}_i} J_i(\boldsymbol{\theta}_{t,i}, \boldsymbol{\theta}_{t,-i}) - \boldsymbol{\theta}^* \\ &\approx \left(I + \eta_i \nabla_{\boldsymbol{\theta}_i}^2 J_i(\boldsymbol{\theta}_i^*, \boldsymbol{\theta}_{-i}^*), \eta_i \nabla_{\boldsymbol{\theta}_i, \boldsymbol{\theta}_{-i}}^2 J_i(\boldsymbol{\theta}_i^*, \boldsymbol{\theta}_{-i}^*) \right) \begin{pmatrix} \hat{\boldsymbol{\theta}}_{t,i} \\ \hat{\boldsymbol{\theta}}_{t+1,-i} \end{pmatrix} \quad \text{First order Taylor expansion} \\ &= \hat{\boldsymbol{\theta}}_{t,i} + \eta_i \mathbf{A}_i \hat{\boldsymbol{\theta}}_{t,i} + \eta_i \mathbf{B}_i \hat{\boldsymbol{\theta}}_{t+1,-i} \\ \therefore \hat{\boldsymbol{\theta}}_{t+1} &= \hat{\boldsymbol{\theta}}_t + \boldsymbol{\eta} \mathbf{A} \hat{\boldsymbol{\theta}}_t + \boldsymbol{\eta} \mathbf{B} \mathbf{D} \hat{\boldsymbol{\theta}}_{t+1}. \end{aligned} \quad (30)$$

By analyzing the distance r of the GSPPM iterates from the stationary point,

$$\begin{aligned}
r_{t+1}^2 &= \|\hat{\boldsymbol{\theta}}_{t+1}\|^2 \\
&= \hat{\boldsymbol{\theta}}_t^T (I + \boldsymbol{\eta}\mathbf{A})^T (I - \boldsymbol{\eta}\mathbf{B}\mathbf{D})^{-T} (I - \boldsymbol{\eta}\mathbf{B}\mathbf{D})^{-1} (I + \boldsymbol{\eta}\mathbf{A}) \hat{\boldsymbol{\theta}}_t \\
&\leq \frac{\lambda \max (I + \boldsymbol{\eta}\mathbf{A})^2}{\lambda \min (I - \boldsymbol{\eta}\mathbf{B}\mathbf{D})^2} r_t^2.
\end{aligned} \tag{31}$$

Thus, for any $\{\eta_i\}$ satisfying $\frac{\lambda \max (I + \boldsymbol{\eta}\mathbf{A})^2}{\lambda \min (I - \boldsymbol{\eta}\mathbf{B}\mathbf{D})^2} < 1$, GSPPM iterates converge asymptotically to the local Nash equilibrium. \square

A.3 Proof of Theorem 6.4

Let us define $\hat{\boldsymbol{\theta}}_{t,i}^{(k)} = \boldsymbol{\theta}_{t,i}^{(k)} - \boldsymbol{\theta}_{t,i}^*$ and $\hat{\boldsymbol{\theta}}^{(k)} = [\hat{\boldsymbol{\theta}}_1^{(k)}, \dots, \hat{\boldsymbol{\theta}}_n^{(k)}]^T$ for all $i \in \mathcal{I}$. It follows by linearizing the system about the stationary point $\boldsymbol{\theta}^*$,

$$\begin{aligned}
\hat{\boldsymbol{\theta}}_{t+1,i}^{(k)} &= \boldsymbol{\theta}_{t,i} + \eta_i \nabla_{\boldsymbol{\theta}_{t,i}} J_i(\boldsymbol{\theta}_{t,i}, \boldsymbol{\theta}_{t,-i}^{(k-1)}) - \boldsymbol{\theta}_i^* \\
&\approx (I + \eta_i \nabla_{\boldsymbol{\theta}_{t,i}}^2 J_i(\boldsymbol{\theta}_{t,i}^*, \boldsymbol{\theta}_{t,-i}^*), \eta_i \nabla_{\boldsymbol{\theta}_{t,-i}}^2 J_i(\boldsymbol{\theta}_{t,i}^*, \boldsymbol{\theta}_{t,-i}^*)) \begin{pmatrix} \hat{\boldsymbol{\theta}}_{t,i} \\ \hat{\boldsymbol{\theta}}_{t,-i}^{(k-1)} \end{pmatrix} \quad \text{First order Taylor expansion} \\
&= \hat{\boldsymbol{\theta}}_{t,i} + \eta_i \mathbf{A}_i \hat{\boldsymbol{\theta}}_{t,i} + \eta_i \mathbf{B}_i \hat{\boldsymbol{\theta}}_{t,-i}^{(k-1)} \\
\therefore \hat{\boldsymbol{\theta}}_t^{(k)} &= \hat{\boldsymbol{\theta}}_t + \boldsymbol{\eta}\mathbf{A}\hat{\boldsymbol{\theta}}_t + \boldsymbol{\eta}\mathbf{B}\mathbf{D}\hat{\boldsymbol{\theta}}_t^{(k-1)}.
\end{aligned} \tag{32}$$

By analyzing the distance $r^{(k)}$ of the iterates from the stationary point,

$$\begin{aligned}
\left(r_t^{(k)}\right)^2 &= \|\hat{\boldsymbol{\theta}}_t^{(k)}\|^2 \\
&= \hat{\boldsymbol{\theta}}_t^T (I + \boldsymbol{\eta}\mathbf{A})^T (I + \boldsymbol{\eta}\mathbf{A}) \hat{\boldsymbol{\theta}}_t + \hat{\boldsymbol{\theta}}_t^T (I + \boldsymbol{\eta}\mathbf{A})^T \boldsymbol{\eta}\mathbf{B}\mathbf{D} \hat{\boldsymbol{\theta}}_t^{(k-1)} \\
&\quad + (\hat{\boldsymbol{\theta}}_t^{(k-1)})^T \mathbf{D}^T \mathbf{B}^T \boldsymbol{\eta}^T (I + \boldsymbol{\eta}\mathbf{A}) \hat{\boldsymbol{\theta}}_t + (\hat{\boldsymbol{\theta}}_t^{(k-1)})^T \mathbf{D}^T \mathbf{B}^T \boldsymbol{\eta}^T \boldsymbol{\eta}\mathbf{B}\mathbf{D} \hat{\boldsymbol{\theta}}_t^{(k-1)} \\
&\leq \left(\lambda \max (I + \boldsymbol{\eta}\mathbf{A})^2 + 2 \lambda \max (I + \boldsymbol{\eta}\mathbf{A}) \lambda \max (\boldsymbol{\eta}\mathbf{B}\mathbf{D}) \right) \left(r_t^{(0)}\right)^2 + \\
&\quad 2 \lambda \max (I + \boldsymbol{\eta}\mathbf{A}) \lambda \max (\boldsymbol{\eta}\mathbf{B}\mathbf{D}) \left(r_t^{(0)} \nabla_{\max}\right) + \lambda \max (\boldsymbol{\eta}\mathbf{B}\mathbf{D})^2 \left(r_t^{(k-1)}\right)^2.
\end{aligned} \tag{33}$$

defining the bound of the finite-k iterates to the stationary point $\hat{\boldsymbol{\theta}}^*$.

Hence, for any $\{\eta_i\}$ satisfying $\lambda \max (\boldsymbol{\eta}\mathbf{B}\mathbf{D})^2 < 1$, the iterates converge asymptotically to the local Nash Equilibrium. \square

A.4 Proof of the Multi-Agent Performance Difference Lemma

Here we prove the Multi-Agent Performance Difference Lemma (PDL) used in the surrogate loss of ReMAPPO. Recall the definitions of the state s , joint action \mathbf{a} , transition dynamics \mathcal{P} , state occupancy distribution d^π under joint policy π , and reward function for agent i , $\mathcal{R}_i(s, \mathbf{a}, s')$ in Section 3.

Recall the relation between the value and action-value function for agent i under joint policy π ,

$$Q_i^\pi(s, \mathbf{a}) = \mathbb{E}_{\mathcal{P}} [R_i(s, \mathbf{a}, s') + \gamma V_i^\pi(s') | s, \mathbf{a}]. \tag{34}$$

Thus, the advantage function can be written as such:

$$A_i^\pi(s, \mathbf{a}) = \mathbb{E}_{\mathcal{P}} [R_i(s, \mathbf{a}, s') + \gamma V_i^\pi(s') - V_i^\pi(s) | s, \mathbf{a}]. \quad (35)$$

Taking the expectation of the advantage under the updated joint policy π' and transition dynamics \mathcal{P} ,

$$\mathbb{E}_{\pi', \mathcal{P}} [A_i^\pi(s, \mathbf{a})] = \mathbb{E}_{\pi', \mathcal{P}} [R_i(s, \mathbf{a}, s') + \gamma V_i^\pi(s') - V_i^\pi(s) | s, \mathbf{a}]. \quad (36)$$

Multiplying by γ^t and summing over all timesteps $t \geq 0$:

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [A_i^\pi(s, \mathbf{a})] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [R_i(s, \mathbf{a}, s')] + \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s')] - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s)]. \quad (37)$$

The last two sums are reduced by telescoping:

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s')] - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s)] = -\mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s_0)] + \lim_{t \rightarrow \infty} \gamma^{t+1} \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s)], \quad (38)$$

where we have slightly abused notation to denote s_0 as the initial state. Since $\gamma < 1$ and $R_i(s, \mathbf{a}, s') < \infty$, the limit term tends to 0. Hence,

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [A_i^\pi(s, \mathbf{a})] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [R_i(s, \mathbf{a}, s')] - \mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s_0)] \quad (39)$$

Recall that $\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [R_i(s, \mathbf{a}, s')] = J_i(\pi')$ and $\mathbb{E}_{\pi', \mathcal{P}} [V_i^\pi(s_0)] = J_i(\pi)$. Hereafter we have,

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [A_i^\pi(s, \mathbf{a})] = J_i(\pi') - J_i(\pi) \quad (40)$$

Recall the conversion from the discounted time-series expectation to state-action occupancy expectation,

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi', \mathcal{P}} [A_i^\pi(s, \mathbf{a})] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi'}, \mathbf{a} \sim \pi'} [A_i^\pi(s, \mathbf{a})], \quad (41)$$

providing the key result:

$$J_i(\pi') - J_i(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi'}, \mathbf{a} \sim \pi'} [A_i^\pi(s, \mathbf{a})], \quad (42)$$

Importance sampling changes the expectation:

$$J_i(\pi') - J_i(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi'}, \mathbf{a} \sim \pi} \left[\frac{\pi'_i}{\pi_i} A_i^\pi(s, \mathbf{a}) \right], \quad (43)$$

which can be represented by self and non-self actions relative to agent i ,

$$J_i(\pi') - J_i(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi'}, \mathbf{a} \sim \pi_i, \pi'_{-i}} \left[\frac{\pi'_i}{\pi_i} \frac{\pi'_{-i}}{\pi_{-i}} A_i^\pi(s, \mathbf{a}) \right], \quad (44)$$

Note that the surrogate loss of MAPPO treats the non-self action distribution as constant for simplicity, while ReMAPPO uses the full form of the Multi-Agent PDL.

B Details of the illustrative example

We report the full details of the toy problem introduced in Sec. 6.2, which we here refer to as the *Meet-up* problem.

Environment properties. The problem is designed as a simple 2-player continuous cooperative game in a 2D space. The state of the game $s = (s_1, s_2) \in \mathbb{R}^4$ encodes the location of the two players, with $s_i \in \mathbb{R}^2$. For the sake of simplicity, agents can only move by a fixed distance step of 1 around their current position, towards a chosen direction. The initial state of the two agents is deterministic and fixed to $\iota = (\iota_1 = (0, 0), \iota_2 = (3, 2))$. We assume undiscounted returns ($\gamma = 1$) and terminate an episode when the agents effectively meet each other as a result of their actions.

Policy parameterization. Although one-dimensional continuous actions are trivially tractable for one-step games, sequential decision making problems demand finding policies that respond optimally for any possible configuration s of the game. Here, we reduce the complexity of the problem by conveniently parameterizing each agent as single-parameter policies. In particular, we define an agent action as a 1-DoF unit vector $a_i \in \mathbb{R}^2$, and parameterize the deterministic policy of agent i with $\theta_i \in \mathbb{R}$, as

$$\pi_i(s) = \begin{cases} (\cos \theta_i, \sin \theta_i)^\top & \text{if, } s = \iota \\ \pi_i^*(s) & \text{if, } s \neq \iota \end{cases} \quad (45)$$

where, $\pi_i^*(s) = \frac{s_{-i} - s_i}{\|s_{-i} - s_i\|}$ is the optimal policy that goes straight towards the other agent. In other words, we assume that both agents will act optimally after taking the first action, and we only parametrize the agents decisions at the starting state. This design choice allows to easily study the joint policy space directly, as well as computing the closed-form solution of the return $J(\cdot)$ (see below).

Solving the Meet-up problem. We design the reward function of the Meet-up problem to reward each agent for getting closer to the other agent after the effect of both actions. We achieve this by computing the cosine similarity between the agent’s action a_i and the actual direction that would have led closest to the other agent:

$$R_i(s, a, s') = a_i \cdot \pi_i^*(s_i, s'_{-i}) - 1 = a_i \cdot \frac{s'_{-i} - s_i}{\|s'_{-i} - s_i\|} - 1. \quad (46)$$

Here, we denote the joint action as $a = (a_1, a_2)$, and the next state as $s' = (s_1 + a_1, s_2 + a_2)$. Note that a -1 offset is added so that both the reward signal and the return $J_i(\theta_i, \theta_{-i})$ of each agent is always ≤ 0 . In turn, this makes the computation of the optimal value function $V^*(s)$ of this game trivial: the strategy of moving towards each other in a straight line leads to returns of 0 from any state s ; since this is the maximum return, this joint policy must also be optimal, and $V^*(s) = 0 \forall s$ is the unique optimal value function. We now derive the analytical form of $J_i(\theta_i, \theta_{-i})$, as needed to compute the recursive gradient updates. Given that both agents are assumed to act optimally in any state besides the starting state, we can conveniently write the return as

$$\begin{aligned} J_i(\theta_i, \theta_{-i}) &= R_i(\iota, a, s') + V_i(s') \\ &= R_i(\iota, a, s') + V^*(s') = \\ &= R_i(\iota, a, s') \end{aligned} \quad (47)$$

where s' is the resulting state after the players’ first actions $a_1 = (\cos \theta_1, \sin \theta_1)$ and $a_2 = (\cos \theta_2, \sin \theta_2)$. Following this, we may therefore compute the gradient of the return for any pair of agent policies θ_1, θ_2 in closed form:

$$\begin{aligned}
\nabla_{\theta_i} J_i(\theta_i, \theta_{-i}) &= \nabla_{\theta_i} R_i(l, a, s') \\
&= \nabla_{\theta_i} \left(a_i \cdot \pi_i^*(l_i, s'_{-i}) \right) \\
&= \nabla_{\theta_i} \left(a_i \right) \cdot \pi_i^*(l_i, s'_{-i}) \\
&= \nabla_{\theta_i} \begin{pmatrix} \cos \theta_i \\ \sin \theta_i \end{pmatrix} \cdot \pi_i^*(l_i, s'_{-i}) \\
&= \begin{pmatrix} -\sin \theta_i \\ \cos \theta_i \end{pmatrix} \cdot \pi_i^*(l_i, s'_{-i})
\end{aligned} \tag{48}$$

In conclusion, Eq. 48 allows us to compute the recursive reasoning steps with the true analytical gradient of the return.

C Additional Results

C.1 The Stability of the ReMAPPO Importance Sampling Ratio

As seen in Equation 4, the surrogate loss for ReMAPPO contains an importance sampling (IS) ratio which could theoretically grow exponentially large in environments with many agents, resulting in instability. We empirically investigate this by comparing the win rate, IS ratio, and clipping fraction (the fraction of minibatch updates at each update step which reach the clipping boundary ϵ) of the 3 PPO variants in the main paper. Additionally, we test a variant of ReMAPPO which geometrically normalizes the non-self portion of the IS ratio by the number of non-self agents. The aim is to investigate whether this proportionally reduces the clipping fraction and how that effects performance. Table 3 summarizes each variant of PPO tested in this ablation and their respective IS ratio.

Figure 9 compares the win rates, IS ratios, and clipping statistics of the 4 largest SMAX maps tested in the main paper: 27m_vs_30m, smacv2_20_units, smacv2_10_units, and 10m_vs_11m. As expected, the IS ratio (and therefore the clipping fraction) of ReMAPPO far exceeds that of IPPO and MAPPO due to the compounding IS ratio of each non-self agent (notice, however, that it is still not close to the clipping boundary of 1 ± 0.2 on average). Geometric normalization brings the IS ratio and clipping fraction of ReMAPPO_GEOMEAN closer to that of MAPPO. Crucially, we observe that normalization has a detrimental effect on performance: the mean win rate of ReMAPPO_GEOMEAN falls somewhere between the win rates of ReMAPPO and MAPPO in 3 of the maps tested (and falls short of MAPPO in 1 map). This suggests that the large IS ratio multiplier of ReMAPPO is the cause of increased performance in SMAX in recursive updates, and remains relatively stable even in environments with up to 27 agents. Future work will consider whether these effects are still observed in environments with many more agents (e.g. 100-1000), as well as the sensitivity of the IS ratio to the individual learning rates.

Algorithm	Importance Sampling Ratio
IPPO	$\frac{\pi'_i(a_i s)}{\pi_i(a_i s)}$
MAPPO	$\frac{\pi'_i(a_i s)}{\pi_i(a_i s)}$
ReMAPPO	$\frac{\pi'_i(a_i s)}{\pi_i(a_i s)} \cdot \frac{\pi'_{-i}(\mathbf{a}_{-i} s)}{\pi_{-i}(\mathbf{a}_{-i} s)}$
ReMAPPO_GEOMEAN	$\frac{\pi'_i(a_i s)}{\pi_i(a_i s)} \cdot \sqrt[N-1]{\frac{\pi'_{-i}(\mathbf{a}_{-i} s)}{\pi_{-i}(\mathbf{a}_{-i} s)}}$

Table 3: Importance Sampling Ratio of each PPO Variant

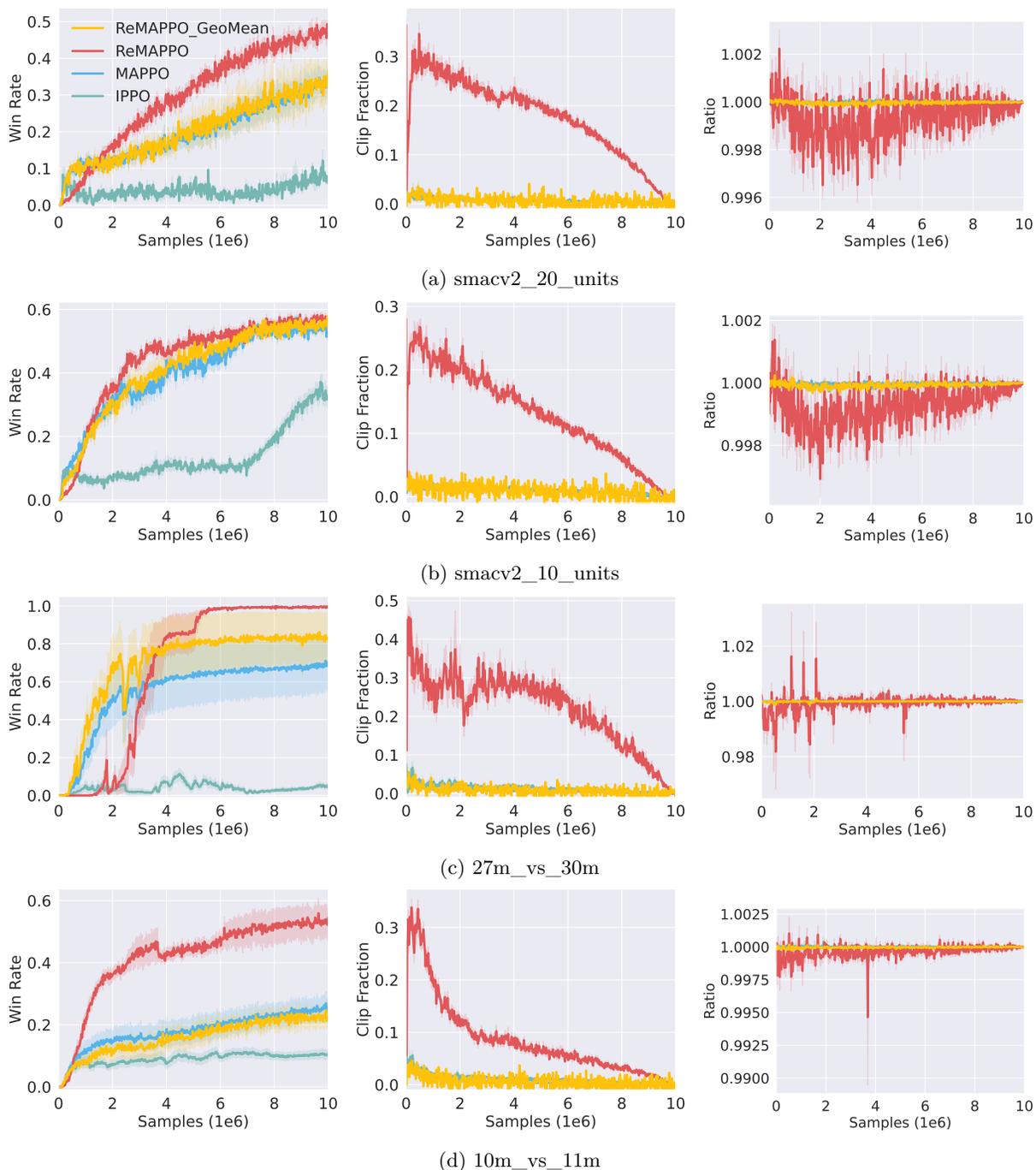


Figure 9: Mean win rate, clipping fraction, and importance sampling ratio with standard errors of IPPO, MAPPO, ReMAPPO, and ReMAPPO with geometric ratio normalization on the 4 largest SMAX maps.

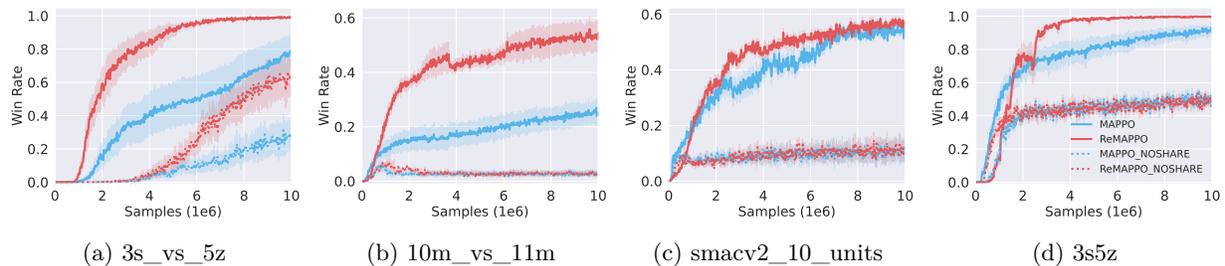


Figure 10: Mean win rate and standard error of MAPPO and ReMAPPO, with and without parameter sharing

C.2 Parameter Sharing

Parameter sharing between agents is a common technique used in MARL research in order to speed up convergence in cooperative settings (Rutherford et al., 2023; Peng et al., 2021; Forkel & Foerster, 2025). In this section, we ablate the design choice of parameter sharing in 4 SMAX maps from the main paper.

Figure 10 clearly demonstrates the importance of parameter sharing between agents in SMAX; when agents use separate parameters, the final win rates and sample efficiency are drastically reduced. This is because parameter sharing allows agents to utilize shared useful representations, effectively allowing them to learn indirectly from the experiences of others. Notably, in 2 out of 4 maps tested, ReMAPPO retains some advantage in final performance and/or sample efficiency (3s_vs_5z and 3s5z), whereas in the other 2 maps, the lack of shared parameters leads to near-total performance collapse.