

VARIATIONAL INEQUALITY METHODS FOR MULTI-AGENT REINFORCEMENT LEARNING: PERFORMANCE AND STABILITY GAINS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-agent reinforcement learning (MARL) poses distinct challenges as agents learn strategies through experiences. Gradient-based methods often fail to converge in MARL, and performances are highly sensitive to initial random seeds, contributing to what has been termed the *MARL reproducibility crisis*. Concurrently, significant advances have been made in solving Variational Inequalities (VIs)—which include equilibrium-finding problems—particularly in addressing the non-converging rotational dynamics that impede convergence of traditional gradient-based optimization methods. This paper explores the potential of leveraging VI-based techniques to improve MARL training. Specifically, we study the integration of VI methods—namely, Nested-Lookahead VI (nLA-VI) and Extragradient (EG)—into the *multi-agent deep deterministic policy gradient* (MADDPG) algorithm. We present a VI reformulation of the actor-critic algorithm for both single- and multi-agent settings. We introduce three algorithms that use nLA-VI, EG, and a combination of both, named *LA-MADDPG*, *EG-MADDPG*, and *LA-EG-MADDPG*, respectively. Our empirical results show that these VI-based approaches yield significant performance improvements in benchmark environments, such as the zero-sum games: *rock-paper-scissors and matching pennies*, where equilibrium strategies can be quantitatively assessed, and the *Multi-Agent Particle Environment: Predator-prey* benchmark, where VI-based methods also yield balanced participation of agents from the same team, further highlighting the substantial impact of advanced optimization techniques on MARL performance.

1 INTRODUCTION

Multi-agent reinforcement learning (MARL) is a powerful machine learning approach for solving complex, multi-player problems in diverse domains. It has been applied to tasks such as coordinating multi-robot and multi-drone systems for instance for search and warehouse automation, optimizing traffic flow and vehicle platooning in autonomous driving, managing energy distribution in smart grids, simulating financial markets and automated trading, improving patient management and drug discovery in healthcare, enhancing network performance in telecommunications, training intelligent agents in games (e.g., Omidshafiei et al., 2017; Vinyals et al., 2017; Spica et al., 2018; Zhou et al., 2021; Bertsekas, 2021, among others). In MARL, a system of N agents seeks to jointly optimize a shared objective, with each agent operating based on its own policy, derived from its observations of the environment. Depending on their reward objectives, agents may exhibit cooperative, competitive, or mixed behaviors. These interactions introduce complex learning dynamics, making MARL significantly more challenging and distinct from single and actor-only reinforcement learning (RL).

Despite the applicability of MARL to a wide range of problems, their deployment and research development face significant challenges. Key issues include: (i) The iterative training process in data-driven MARL is notoriously difficult, often failing to start to converge. This lack of convergence remains a fundamental obstacle. (ii) Performance is highly sensitive to small changes in hyperparameters or initial random seed variations, leading to unpredictable outcomes. This hinders reproducibility. These challenges are also evident in single-agent reinforcement learning with actor-critic structure (Konda & Tsitsiklis, 1999)—an instance of a two-player game. Such methods require meticulous hyperparameter tuning, and their outcomes can vary significantly based on the

random seeds used for sampling and model initialization (Wang et al., 2022; Eimer et al., 2023). This variability undermines reproducibility, complicating both research progress and real-world deployment (Henderson et al., 2019; Lynnerup et al., 2019). In MARL, these challenges are even more pronounced, contributing to what is often referred to as the *reproducibility crisis* (Bettini et al., 2024). Small changes to hyperparameter values can drastically alter results, as demonstrated by Gorsane et al. (2022). Their study reveals significant performance variability across different seeds in popular MARL benchmarks such as the *StarCraft* multi-agent challenge (Samvelyan et al., 2019). Additionally, gradient-based optimization methods in MARL face unique challenges, such as difficulties in exploring the joint policy space of multiple agents (Li et al., 2023; Christianos et al., 2021), often leading to suboptimal solutions. Some MARL structures also exhibit inherent cycling effects (Zheng et al., 2021), further exacerbating the problem of convergence. In this work, we focus primarily on addressing challenge (i) by improving training stability. In doing so, we also aim to mitigate the high variability caused by random seed and hyperparameter sensitivity, thereby partially addressing (ii).

A concurrent line of works focuses on the *Variational Inequality* (VI) problem, a general class of problems that encompasses both equilibria- and optima-finding problems; see Section 3 for definition. VIs generalize standard minimization problems. In this case, the operator F is a gradient field $F \equiv \nabla f$. However, by allowing F to be a more general vector field, VIs also model problems such as finding equilibria in zero-sum and general-sum games (Cottle & Dantzig, 1968; Rockafellar, 1970). It has been observed that standard optimization methods that perform well in minimization tasks of the form $\min_z f(z)$ —where $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a real-valued loss function—often fail to solve some simple instances of VIs. This failure is due to the *rotational* component inherent in the gradient dynamics of these settings, which causes the latest iterate to cycle around the solution, leading to non-convergence (Mescheder et al., 2018; Balduzzi et al., 2018). More precisely, the Jacobian of the associated vector field (see def. in Section 3) can be decomposed into a symmetric and antisymmetric component (Balduzzi et al., 2018), where each behaves as a *potential* (Monderer & Shapley, 1996) and a *Hamiltonian* (purely rotational) game, resp. For instance, the gradient descent method for the simple $\min_{z_1 \in \mathbb{R}^{d_1}} \max_{z_2 \in \mathbb{R}^{d_2}} z_1 \cdot z_2$ game, which simultaneously updates z_1, z_2 rotates around the solution for infinitesimally small learning rates, and diverges away from it for practical choices of its value. As a result, all variation methods based on gradient descent, such as *Adam* (Kingma & Ba, 2015) have no hope of converging for some broad problem classes of VIs. This problematic behavior is particularly pronounced when the separate sets of parameters are neural networks, as in generative adversarial networks (GANs, Goodfellow et al., 2014). As a result, when GANs were first introduced, substantial computational resources were required to fine-tune hyperparameters (Radford et al., 2016). In addition, even for highly tuned hyperparameters, training with minimization methods often (eventually) diverges away (Chavdarova et al., 2021), in sharp contrast to standard minimization. The original GAN formulation (Goodfellow et al., 2014) and practical GAN implementations that are not necessarily zero-sum, are all instances of VIs. The above training difficulties inspired numerous recent research efforts to develop numerical methods to approximately solve variational inequalities (VIs) and to study how VI optimization differs from minimization. Various algorithms have been proposed and studied; reviewed in Sections 2 and 3 and Appendix A.1.1.

This paper hypothesizes that training instabilities in MARL and actor-critic RL primarily arise from the rotational dynamics inherent in competitive learning objectives. This issue is compounded by the prevalent reliance in practice on gradient descent-based methods, such as Adam (Kingma & Ba, 2015), which are known to induce similar issues in the Variational Inequality (VI) literature and have been observed to cause analogous challenges in other VI applications (Goodfellow et al., 2014; Gidel et al., 2019a; Chavdarova et al., 2021). Since VI optimization methods are specifically designed to address rotational vector fields, this paper raises the following question:

Can MARL algorithms benefit from the application of Variational Inequality optimization methods?

To address this question, we focus on the *multi-agent deep deterministic policy gradient* (MADDPG) method (Lowe et al., 2017) and integrate it with the (combination of) *nested-Lookahead-VI* (nLA-VI) (Chavdarova et al., 2021) and *Extragradient* (EG) (Korpelevich, 1976) methods for solving variational inequalities (VIs). Our main contributions are as follows:

- Primarily, we present a VI perspective for multi-agent reinforcement learning (MARL) problems.

- We propose the *LA-MADDPG*, *EG-MADDPG*, and *LA-EG-MADDPG* algorithms, which extend MADDPG by combining it with nLA-VI, and with EG and a mix of both (respectively) in the actor-critic parameter optimization for all agents.
- We evaluate the proposed methods against standard optimization approaches in several two-player games and benchmarks from the Multi-Agent Particle Environment (MPE, [Lowe et al., 2017](#)).
- We also discuss additional insights into the use of rewards as a performance metric in MARL.

2 RELATED WORKS

Our work draws mainly from two lines of work that we review next.

Multi-Agent Reinforcement Learning (MARL). Various MARL algorithms have been developed ([Lowe et al., 2017](#); [Iqbal & Sha, 2018](#); [Ackermann et al., 2019](#); [Yu et al., 2021](#)), with some extending existing single-agent reinforcement learning (RL) methods ([Rashid et al., 2018](#); [Son et al., 2019](#); [Yu et al., 2022](#); [Kuba et al., 2022](#)). [Lowe et al. \(2017\)](#) extend the actor-critic algorithm to the MARL setting using the *centralized training decentralized execution* framework. In the proposed algorithm, named *multi-agent deep deterministic policy gradient (MADDPG)*, each agent in the game consists of two components: an *actor* and a *critic*. The actor is a policy network that has access only to the local observations of the corresponding agent and is trained to output appropriate actions. The critic is a value network that receives additional information about the policies of other agents and learns to output the Q-value; see Section 3. After a phase of experience collection, a batch is sampled from a replay buffer and used for training the agents. To our knowledge, all deep MARL implementations rely on either stochastic gradient descent or *Adam* optimizer ([Kingma & Ba, 2015](#)) to train all networks. Game theory and MARL share many foundational concepts, and several studies explore the relationships between the two fields ([Yang & Wang, 2021](#); [Fan, 2024](#)), with some using game-theoretic approaches to model MARL problems ([Zheng et al., 2021](#)). This work proposes incorporating game-theoretic techniques into the optimization process of existing MARL methods to determine if these techniques can enhance MARL optimization.

Variational Inequalities (VIs). VIs were first formulated to understand the equilibrium of a dynamical system ([Stampacchia, 1964](#)). Since then, they have been studied extensively in mathematics including operational research and network games (see [Facchinei & Pang, 2003](#), and references therein). More recently, after the shown training difficulties of GANs ([Goodfellow et al., 2014](#))—which are an instance of VIs—an extensive line of works in machine learning studies the convergence of iterative gradient-based methods to solve VIs numerically. Since the last and average iterates can be far apart when solving VIs (see e.g., [Chavdarova et al., 2019](#)), these works primarily aimed at obtaining last-iterate convergence for special cases of VIs that are important in applications, including bilinear or strongly monotone games (e.g., [Tseng, 1995](#); [Malitsky, 2015](#); [Facchinei & Pang, 2003](#); [Daskalakis et al., 2018](#); [Liang & Stokes, 2019](#); [Gidel et al., 2019b](#); [Azizian et al., 2020](#); [Thekumparampil et al., 2022](#)), VIs with cocoercive operators ([Diakonikolas, 2020](#)), or monotone operators ([Chavdarova et al., 2023](#); [Gorbunov et al., 2022](#)). Several works (*i*) exploit continuous-time analyses ([Ryu et al., 2019](#); [Bot et al., 2020](#); [Rosca et al., 2021](#); [Chavdarova et al., 2023](#); [Bot et al., 2022](#)), (*ii*) establish lower bounds for some VI classes (e.g., [Golowich et al., 2020b;a](#)), and (*iii*) study the constrained setting ([Daskalakis & Panageas, 2019](#); [Cai et al., 2022](#); [Yang et al., 2023](#); [Chavdarova et al., 2024](#)), among other. Due to the computational complexities involved in training neural networks, iterative methods that rely solely on first-order derivative computation are the most commonly used approaches for solving variational inequalities (VIs). However, standard gradient descent and its momentum-based variants often fail to converge even on simple instances of VIs. As a result, several alternative methods have been developed to address this issue. Some of the most popular first-order methods for solving VIs include the *extragradient* method ([Korpelevich, 1976](#)), *optimistic gradient* method ([Popov, 1980](#)), *Halpern* method ([Diakonikolas, 2020](#)), and (nested) *Lookahead-VI* method ([Chavdarova et al., 2021](#)); these are discussed in detail in Section 3 and Appendix A.1.1. In this work, we primarily focus on the nested Lookahead-VI method, which has achieved state-of-the-art results on the CIFAR-10 ([Krizhevsky, 2009](#)) benchmark for generative adversarial networks ([Goodfellow et al., 2014](#)).

3 PRELIMINARIES

Notation. Bold small letters denote vectors, and curly capital letters denote sets. Let \mathcal{Z} be a convex and compact set in the Euclidean space, with inner product $\langle \cdot, \cdot \rangle$.

Setting: multi-agent deep deterministic policy gradient. *Markov Games* (MGs) extend Markov Decision Processes to the multi-agent setting. In a Markov Game, N agents interact within an environment characterized by a set of states \mathcal{S} . Agents receive observations $\mathbf{o}_i, i = 1, \dots, N$ of the current environment state $\mathbf{s} \in \mathcal{S}$. Based on their policies π_i , each agent i chooses an action $a_i \in \mathcal{A}_i$ from predefined finite action sets $\mathcal{A}_i, i = 1, \dots, N$. These actions, collectively represented as \mathbf{a} , are then applied to the environment, which transitions to a new state $\hat{\mathbf{s}} \in \mathcal{S}$ according to a transition function $\mathcal{T}: \mathcal{S} \rightarrow \mathcal{S}$. Each agent receives a reward $r_i, i = 1 \dots N$, and a new observation $\hat{\mathbf{o}}_i$. In the MARL setting herein, each agent has its own Q-value that is, how much reward it expects to get from a state when joint action \mathbf{a} is performed.

Multi-agent deep deterministic policy gradient (MADDPG, Lowe et al., 2017), extends Deep deterministic policy gradient (DDPG, Lillicrap et al., 2015) to multi-agent setting using the framework of *centralized training decentralized execution*. Each agent i has (i) a *critic network*— Q_i —which acts as a centralized action-value function, (ii) a target critic network Q'_i that is less frequently updated with the most recent Q_i parameters for learning stability, (iii) an *actor network* μ_i which represents the policy to be updated, and (iv) a target actor network μ'_i from which it selects its actions and is periodically updated with the learned policy μ_i . Both the Critic and Actor networks are modeled using feedforward networks, parameterized by \mathbf{w} and $\boldsymbol{\theta}$ respectively.

The VI framework. Broadly speaking, VIs formalize equilibrium-seeking problems. The goal is to find an equilibrium \mathbf{z}^* from the domain of continuous strategies \mathcal{Z} , such that:

$$\langle \mathbf{z} - \mathbf{z}^*, F(\mathbf{z}^*) \rangle \geq 0, \quad \forall \mathbf{z} \in \mathcal{Z}, \quad (\text{VI})$$

where the so-called *operator* $F: \mathcal{Z} \rightarrow \mathbb{R}^n$ is continuous, and \mathcal{Z} is a subset of the Euclidean d -dimensional space \mathbb{R}^d . Thus, VIs are defined by the tuple F, \mathcal{Z} , denoted herein as $\text{VI}(F, \mathcal{Z})$. This problem is equivalent to standard minimization, when $F \equiv \nabla f$, where f is a real-valued function $f: \mathbb{R}^d \rightarrow \mathbb{R}$. We refer the reader to (Facchinei & Pang, 2003) for an introduction and examples. To illustrate the relevance of VIs to multi-agent problems, consider the following example. Suppose we have N agents, each with a strategy $\mathbf{z}_i \in \mathbb{R}^{d_i}$, and let us denote the joint strategy with $\mathbf{z} \equiv [\mathbf{z}_1^\top, \dots, \mathbf{z}_N^\top]^\top \in \mathbb{R}^d, d = \sum_{i=1}^N d_i$. Each agent aims to optimize its objective $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$. Then, finding an equilibrium in this game is equivalent to solving a VI where F corresponds to: $F(\mathbf{z}) \equiv [\nabla_{\mathbf{z}_1} f_1(\mathbf{z}), \dots, \nabla_{\mathbf{z}_N} f_N(\mathbf{z})]^\top$.

Methods for solving VIs. The *gradient descent* method naturally extends for the VI problem as follows:

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \eta F(\mathbf{z}_t), \quad (\text{GD})$$

where t denotes the iteration count, and $\eta \in (0, 1)$ the step size or learning rate. The *nested-Lookahead-VI* algorithm for VI problems (Chavdarova et al., 2021), originally proposed for minimization by Zhang et al. (2019), is a general wrapper of a “base” optimizer where, at every step t : (i) a copy of the current iterate $\tilde{\mathbf{z}}_t$ is made: $\tilde{\mathbf{z}}_t \leftarrow \mathbf{z}_t$, (ii) $\tilde{\mathbf{z}}_t$ is updated $k \geq 1$ times, yielding $\tilde{\omega}_{t+k}$, and finally (iii) the actual update \mathbf{z}_{t+1} is obtained as a *point that lies on a line between the current \mathbf{z}_t iterate and the predicted one $\tilde{\mathbf{z}}_{t+k}$* :

$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t + \alpha(\tilde{\mathbf{z}}_{t+k} - \mathbf{z}_t), \quad \alpha \in [0, 1]. \quad ((\text{nested})\text{LA-VI})$$

Notice that we can apply this idea recursively, and when the base optimizer is *(nested)LA-VI* (at some level), then we have *nested LA-VI*, as proposed in Algorithm 3 in (Chavdarova et al., 2021).

Extragradient (Korpelevich, 1976) uses a “prediction” step to obtain an extrapolated point $\mathbf{z}_{t+\frac{1}{2}}$ using GD: $\mathbf{z}_{t+\frac{1}{2}} = \mathbf{z}_t - \eta F(\mathbf{z}_t)$, and the gradients at the *extrapolated* point are then applied to the *current* iterate \mathbf{z}_t as follows:

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \eta F\left(\mathbf{z}_t - \eta F(\mathbf{z}_t)\right), \quad (\text{EG})$$

where $\eta > 0$ is a learning rate (step size). Unlike gradient descent, EG converges on some simple game instances, such as in games linear in both players (Korpelevich, 1976).

Algorithm 1 Procedure Pseudocode for nLA-VI, called from Algorithm 2.

```

216 1: procedure NESTEDLOOKAHEAD:
217
218 2: Input: Number of agents  $N$ , current episode  $e$ , current actor weights and snapshots
219    $\{(\theta_i, \theta_{i,s}, \theta_{i,ss})\}_{i=1}^N$ , current critic weights and snapshots  $\{(\mathbf{w}_i, \mathbf{w}_{i,s}, \mathbf{w}_{i,ss})\}_{i=1}^N$ , lookahead
220   hyperparameters  $k_s, k_{ss}$  (where  $k_{ss}$  can be  $\emptyset$ ) and  $\alpha_\theta, \alpha_w$ 
221 3: Result: Updated actor and critic weights and snapshots for all agents
222 4: if  $e\%k_s == 0$  then
223 5:   for all agent  $i \in 1, \dots, N$  do
224 6:      $\mathbf{w}_i \leftarrow \mathbf{w}_{i,s} + \alpha_w(\mathbf{w}_i - \mathbf{w}_{i,s})$  Apply lookahead (1st level)
225 7:      $\theta_i \leftarrow \theta_{i,s} + \alpha_\theta(\theta_i - \theta_{i,s})$ 
226 8:      $(\theta_{i,s}, \mathbf{w}_{i,s}) \leftarrow (\theta_i, \mathbf{w}_i)$  Update snapshots (1st level)
227 9:   end for
228 10: end if
229 11: if  $k_{ss}$  is not  $\emptyset$  and  $e\%k_{ss} == 0$  then
230 12:   for all agent  $i \in 1, \dots, N$  do
231 13:      $\mathbf{w}_i \leftarrow \mathbf{w}_{i,ss} + \alpha_w(\mathbf{w}_i - \mathbf{w}_{i,ss})$  Apply lookahead (2nd level)
232 14:      $\theta_i \leftarrow \theta_{i,ss} + \alpha_\theta(\theta_i - \theta_{i,ss})$ 
233 15:      $(\theta_{i,s}, \theta_{i,ss}, \mathbf{w}_{i,s}, \mathbf{w}_{i,ss}) \leftarrow (\theta_i, \theta_i, \mathbf{w}_i, \mathbf{w}_i)$  Update snapshots (1st & 2nd level)
234 16:   end for
235 17: end if
236 18: end procedure

```

4 A VI PERSPECTIVE & OPTIMIZATION METHODS FOR MARL

Herein, we describe our proposed approach, which utilizes VI methods in combination with a MARL algorithm. Specifically, we delve into MADDPG, give a VI perspective of it, and describe its combination with extragradient (Korpelevich, 1976) and nested Lookahead (Chavdarova et al., 2021).

4.1 A VI PERSPECTIVE OF MADDPG

Recall that for each $i = 1, \dots, N$ agent, we have:

1. Q -Network, $\mathbf{Q}_i^\mu(\mathbf{x}, a_1, \dots, a_N; \mathbf{w}_i)$: central critic network for agent i ;
2. Policy network, $\mu_i(o_i; \theta_i)$: policy network for agent i ;
3. Target Q -network, $\mathbf{Q}_i^{\mu'}(\mathbf{x}, a_1, \dots, a_N; \mathbf{w}'_i)$;
4. Target policy network, $\mu'_i(o_i; \theta'_i)$.

These networks (maps) are parametrized by $\mathbf{w}_i, \theta_i, \mathbf{w}'_i, \theta'_i$, respectively; with $\mathbf{w}_i, \mathbf{w}'_i \in \mathbb{R}^{d_i^Q}$ and $\theta_i, \theta'_i \in \mathbb{R}^{d_i^\mu}$. The latter two— \mathbf{w}'_i, θ'_i for agent i —are running averages computed as:

$$\begin{aligned} \theta'_i &\leftarrow \tau \theta_i + (1 - \tau) \theta'_i \\ \mathbf{w}'_i &\leftarrow \tau \mathbf{w}_i + (1 - \tau) \mathbf{w}'_i \end{aligned} \quad (\text{Target-Nets})$$

Given a batch of experiences $(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \hat{\mathbf{x}}^j)$ —sampled from a replay buffer (\mathcal{D})—the goal is to find an equilibrium by solving the VI problem with the operator F defined as:

$$F_{\text{MADDPG}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \theta_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \frac{1}{S} \sum_j \left(r_i^j + \gamma \mathbf{Q}_i^{\mu'}(\hat{\mathbf{x}}^j, a'_1, \dots, a'_N; \mathbf{w}'_i) \Big|_{a'_k = \mu'_k(o_k^j)} - \mathbf{Q}_i^\mu(\mathbf{x}^j, \mathbf{a}^j; \mathbf{w}_i) \right)^2 \\ \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(o_i^j; \theta_i) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_N^j; \mathbf{w}_i) \Big|_{a_i = \mu_i(o_i^j)} \\ \vdots \end{bmatrix}, \quad (F_{\text{MADDPG}})$$

and $\mathcal{Z} \equiv \mathbb{R}^d$, where $d = \sum_{i=1}^N (d_i^Q + d_i^\mu)$. Even if $N = 1$, there is still a game between the actor and critic—the update of \mathbf{w}_i depends on θ_i and vice versa.

4.2 PROPOSED METHODS

To solve the VI problem with the operator as defined in (F_{MADDPG}), we propose the *LA-MADDPG*, and *EG-MADDPG* methods, described in detail in this section.

Algorithm 2 Pseudocode for LA-MADDPG: MADDPG with (Nested)-Lookahead-VI.

```

1: Input: Environment  $\mathcal{E}$ , number of agents  $N$ , number of episodes  $T$ , action spaces  $\{\mathcal{A}_i\}_{i=1}^N$ ,
2: random steps  $T_{\text{rand}}$ , learning interval  $T_{\text{learn}}$ , actor networks  $\{\mu_i\}_{i=1}^N$  with weights  $\theta \equiv \{\theta_i\}_{i=1}^N$ ,
3: critic networks  $\{Q_i\}_{i=1}^N$  with weights  $w \equiv \{w_i\}_{i=1}^N$ , target actor networks  $\{\mu'_i\}_{i=1}^N$  with weights
4:  $\theta' \equiv \{\theta'_i\}_{i=1}^N$ , target critic networks  $\{Q'_i\}_{i=1}^N$  with weights  $w' \equiv \{w'_i\}_{i=1}^N$ , learning rates  $\eta_\theta, \eta_w$ ,
5: optimizer  $\mathcal{B}$ , discount factor  $\gamma$ , lookahead parameters  $k_s, k_{ss}, \alpha_\theta, \alpha_w$ , soft update parameter  $\tau$ .
6: Initialize:
7:   Replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
8:   Weights snapshots  $(\theta_s, \theta_{ss}, w_s, w_{ss}) \leftarrow (\theta, \theta, w, w)$ 
9:   for all episode  $e = 1$  to  $T$  do
10:    Sample initial state  $\mathbf{x}$  from  $\mathcal{E}$ 
11:     $step \leftarrow 1$ 
12:    repeat
13:      if  $step \leq T_{\text{rand}}$  then
14:        Randomly select actions for each agent  $i$ 
15:      else
16:        Select actions using policy for each agent  $i$ 
17:      end if
18:      Execute actions  $\mathbf{a}$ , observe rewards  $\mathbf{r}$  and new state  $\hat{\mathbf{x}}$ 
19:      Store  $(\mathbf{x}, \mathbf{a}, \mathbf{r}, \hat{\mathbf{x}})$  in replay buffer  $\mathcal{D}$ 
20:       $\mathbf{x} \leftarrow \hat{\mathbf{x}}$ 
21:      if  $step \% T_{\text{learn}} == 0$  then
22:        Sample a batch  $B$  from  $\mathcal{D}$ 
23:        Use  $B$  and update to solve VI( $F_{\text{MADDPG}}, \mathbb{R}^d$ ) using  $\mathcal{B}$ 
24:        Update target networks:
25:           $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$ 
26:           $w' \leftarrow \tau w + (1 - \tau)w'$ 
27:        end if
28:       $step \leftarrow step + 1$ 
29:    until environment terminates
30:    NESTEDLOOKAHEAD( $N, e, \theta, w, k_s, k_{ss}, \alpha_\theta, \alpha_w$ )
31:  end for
32: Output:  $\theta, w$ 

```

LA-MADDPG. Algorithm 2 describes the *LA-MADDPG* method. Critically, the (nested)LA-VI method is used in the joint strategy space of all players. In this way, the averaging steps address the rotational component of the associated vector field defined by F_{MADDPG} resulting from the adversarial nature of the agents’ objectives. In particular, it is necessary not to use an agent whose parameters have already been averaged at that iteration.

The LA-MADDPG algorithm saves snapshots of the actor and critic networks for all agents, periodically averaging them with the current networks during training. While the MADDPG algorithm (Algorithm 4) runs normally using a base optimizer (e.g., Adam), at every interval k , a lookahead averaging step is performed between the current networks (denoted θ, w), and their saved snapshots θ_s, w_s , as detailed in Algorithm 1. This method updates both the current networks and snapshots with the α -averaged values. Multiple nested lookahead levels can be applied, where each additional level updates its snapshot after a longer interval; see Algorithm 1. We denote lookahead update intervals (episodes) with k subscripted by s and a larger number of s in subscript implies outer lookahead level, e.g., k_s, k_{ss}, k_{sss} for three levels. All agents undergo lookahead updates at the same step, applying this to both the actor and critic parameters simultaneously. An extended version of the algorithm with more detailed notations can be found in appendix in algorithm 5.

(**LA-EG-MADDPG**). For **EG-MADDPG**, **EG** is used for both the actor and critic networks and for all agents; see Algorithm 6 for details. Algorithm 2 can also be used with **EG** as the base optimizer—an option abstracted by \mathcal{B} in Algorithm 2—resulting in **LA-EG-MADDPG**.

On the convergence. Under standard assumptions on the agents’ reward functions, such as convexity, the above VI becomes monotone (see Appendix A.1.1 for definition). In this case, the above methods have convergence guarantees, that is EG-MADDPG (Korpelevich, 1976; Gorbunov et al., 2022), LA-MADDPG (Pethick et al., 2023), and LA-EG-MADDPG (Chavdarova et al., 2021, Thm. 3) have convergence guarantees. Contrary to these, the standard gradient descent method does not converge for this problem (Korpelevich, 1976).

5 EXPERIMENTS

5.1 SETUP

We build upon the open-source *PyTorch* implementation of MADDPG (Lowe et al., 2017)¹. We use the same hyperparameter settings as specified in the original paper; detailed in Appendix A.2. For our experiments, we use two zero-sum games: the *Rock-Paper-Scissors* (RPS) game and *Matching pennies*. We then apply the methods to two of the *Multi-agent Particle Environments* (MPE) (Lowe et al., 2017). We used versions of the games from the *PettingZoo* (Terry et al., 2021) library. We used five different random seeds for training for all games and trained for 50000 episodes per seed for Matching pennies and 60000 for the rest.

2-player game: rock–paper–scissors. Rock–paper–scissors is a widely studied game in multi-agent settings because, in addition to its analytically computable Nash equilibrium that allows for a precise performance measure, it demonstrates interesting cyclical behavior (Zhou, 2015; Wang et al., 2014). The game, with $M = 3$ actions, has a mixed Nash equilibrium where each action is played with equal probability. At equilibrium, each agent’s action distribution is $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. This equilibrium allows us to assess the alignment of learned policies with the optimal strategy. In our experiments version, N players compete in an M -action game over t steps. At each step, players receive an observation of their opponent’s last action. Once all players have selected their actions for the current step, rewards are assigned to each player: as of -1 for a losing, 0 for tie and $+1$ for winning the game. we used $N = 2$ players, $M = 3$ actions, and a time horizon of $t = 25$ steps.

2-player game: matching pennies. The game has $M = 2$ actions, $N = 2$ players: even and odd that compete over t steps. At each step, the players must choose between two actions: Heads or Tails. Even player wins with a reward of $+1$ if the players chose the same action and loses with a -1 otherwise, and vice versa. We used $t = 25$ steps. Similar to Rock–paper–scissors, this game also has mixed Nash equilibrium where each action is played with equal probability. At equilibrium, each agent’s action distribution is $(\frac{1}{2}, \frac{1}{2})$.

We measured and plotted the squared norm of the learned policy probabilities relative to the equilibrium for both *rock–paper–scissors* and *matching pennies*.

MPE: Predator-prey— from the *Multi-Agent Particle Environments* (MPE) benchmark (Lowe et al., 2017). It consists of N *good* agents, L landmarks, and M *adversary* agents. The good agents are faster and receive negative rewards if caught by adversaries, while the slower adversary agents are rewarded for catching a good agent. All agents can observe the positions of other agents, and adversaries also observe the velocities of the good agents. Additionally, good agents are penalized for going out of bounds. This environment combines elements of both competition and collaboration. While all adversaries are rewarded when one of them catches a good agent, their slower speed typically requires them to collaborate, especially since there are usually more adversaries than good agents. For our experiments, we set $N = 1$, $M = 2$, and $L = 2$.

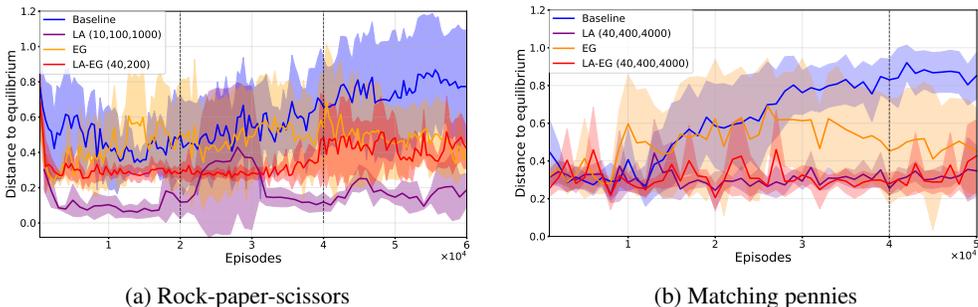
MPE: Physical deception, (Lowe et al., 2017). The game has N good agents, one adversary agent, and N landmarks, with one designated as the target. The adversary does not observe the target and must infer which of the N landmarks is the target one, aiming to get as close as possible and receiving rewards based on its distance from the target. The good agents can observe the target and aim to

¹Available at <https://github.com/Git-123-Hub/maddpg-pettingzoo-pytorch/tree/master>.

378 deceive the adversary while also staying as close as possible to the target. All good agents share the
 379 same reward, based on a combination of their minimum distance to the target and the adversary’s
 380 distance. This game has no “competitive component” for the adversary: its reward depends solely on
 381 its own policy. In our experiments, we set $N = 2$.

382 **Methods.** We evaluate our proposed methods by comparing them to the baseline, which is the
 383 original MADDPG algorithm using Adam (Kingma & Ba, 2015) as the optimizer for all networks.
 384 Throughout the section, we will refer to the LA-MADDPG, EG-MADDPG, and LA-EG-MADDPG
 385 methods as LA, EG, and LA-EG, respectively. When referring to nLA-based methods, we will
 386 indicate the k values for each lookahead level in brackets. For example, LA (10, 1000) represents a
 387 two-level lookahead with $k_s = 10$ and $k_{ss} = 1000$. We also use Adam in combination with the VI
 388 methods for consistency with the baseline.

389 Details on the remaining hyperparameters can be found in Appendix A.2.

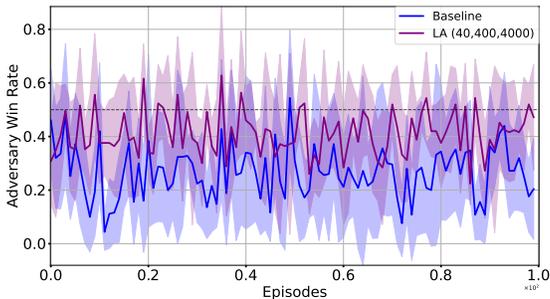


392
393
394
395
396
397
398
399
400
401
402
403 **Figure 1: Comparison on the rock–paper–scissors game and matching pennies game between the *GD-MADDPG*, *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG* methods, denoted as *Baseline*, *LA*, *EG*, *LA-EG*, resp.** x -axis: training episodes. y -axis: total distance of agents’ policies to the equilibrium policy, averaged over 5 seeds. The dotted line depicts the start of the “shifting” (in first-in-first-out order) of the experiences in the buffer.

409 5.2 RESULTS

411 **2-player games: rock–paper–scissors and matching pennies.** Figures 1a and 1b depict the average
 412 distance of the agents’ learned policies from the equilibrium policy. The baseline method eventually
 413 diverges. In contrast, LA-MADDPG consistently converges to a near optimal policy, outperforming
 414 the baseline. While EG-MADDPG behaves similarly to the baseline, combining it with Lookahead
 415 stabilizes the performances. Additionally, Adam exhibits high variance across different seeds, while
 416 Lookahead significantly reduces variance, providing more stable and reliable results—an important
 417 factor in MARL experiments.

418 For LA, we used 0.5 for the α hyperparameter, and after experimenting with several values for k , we
 419 observed that smaller k -values for the innermost LA-averaging works better. Refer to Appendix A.2.1



423
424
425
426
427
428
429
430
431 **Figure 2: Comparison on the MPE:Predator-prey game between the *GD-MADDPG* and *LA-MADDPG*, optimization methods, denoted as *Baseline*, *LA*, resp.** x -axis: evaluation episodes. y -axis: average win rate of adversary agents, averaged over 5 runs with different seeds. The dotted line depicts the desired win rate (0.5) if both agents learn good policies.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Method	Adversary Win Rate
Baseline	$0.45 \pm .16$
LA-MADDPG	$0.53 \pm .11$
EG-MADDPG	$0.56 \pm .27$
LA-EG-MADDPG	$0.51 \pm .14$

Table 1: Means and standard deviations (over 5 seeds) of **adversary win rate on last training episode for MPE: Physical deception**, on 100 test environments. The *win rate* is the fraction of times the adversary was closer to the target. *Closer to 0.5 is better*.

for further discussion. Despite relatively little hyperparameter tuning, the results indicate consistent improvement.

MPE: Predator-prey. Figure 2 depicts the win rate of the adversary against the good agents. While typical training monitors average rewards to indicate convergence, we observed that after training, one adversary learns to chase the good agent while the other’s policy diverges, causing it to move away or wander aimlessly. This suggests a convergence issue in the joint policy space, where one agent’s strategy is affected by the other’s. Our results in Figure 2 demonstrate that using Algorithm 2 improves this behavior, with both adversaries learning to chase the good agent, reflected in a higher win rate. Full method comparisons are provided in Appendix A.3.3.

MPE: Physical deception. Table 1 lists the mean and standard deviation of the adversary’s win rate, indicating how often it managed to be closer to the target. Agents reach equilibrium when both teams win with equal probability across multiple instances. Thus, we used 100 test environments per method per seed. Given the game’s cooperative nature, the baseline performs relatively well, with EG-MADDPG showing similar performance. Both LA-MADDPG and LA-EG-MADDPG outperform their respective base optimizers—baseline and EG-MADDPG.

Summary. Overall, our results indicate the following. (i) With hyperparameter tuning, the proposed VI-based methods achieve significant performance improvements; see Figure 1a. (ii) With informed guesses for hyperparameters (details in Appendix A.2.1), our VI-based methods consistently outperform the baseline methods. (iii) Overall, the proposed VI methods do not yield worse performance than their respective baseline methods, as they effectively address the rotational dynamics.

Comparison among our VI methods & contrasting with GAN conclusions. The widely used Extragradient (EG) method for solving VIs—known for its convergence for monotone VIs—overall performs close to the baseline. EG only introduces a minor local adjustment compared to GD. As such, the results align with expectations: while EG occasionally outperforms GD (the baseline), its performance is often similar. In contrast, nested LA applies a significantly stronger contraction, with the degree of contraction increasing as the number of nested levels increases. This leads to substantial performance gains, particularly in terms of stability, as it prevents the last iterate from diverging. However, if the number of nested levels is too high, the steps can become overly conservative or slow. Based on our experiments, three levels of nested LA yielded the best results (see Fig. 1-a). The results also confirm that the MARL vector field in these games is highly rotational. For scenarios with highly competitive reward structures among agents, we recommend using VI methods with higher contractiveness, such as employing multiple levels of nested LA.

These observations are consistent with results from GANs settings (Chavdarova et al., 2021), while EG offers slight improvements over the baseline, more contractive methods consistently achieve better results.

On the rewards as a metric in MARL. While saturating rewards are commonly used as a performance metric in MARL, our experiments suggest otherwise, consistent with observations made in some previous works such as (Bowling, 2004). In multi-agent games like Rock-paper-scissors, rewards may converge to a target value even with suboptimal policies, leading to misleading evaluations. For instance, in Figure 3 (top row), agents repeatedly choose similar actions, resulting in ties that yield the correct reward but fail to reach equilibrium—leaving them vulnerable to exploitation by a more skilled opponent. Conversely, LA-MADDPG (bottom row) did not fully converge to the maximum reward, but agents learned near-optimal policies by randomizing over their actions, which is the desired equilibrium. This underscores the need for stronger evaluation metrics in multi-agent reinforcement learning, particularly when the true equilibrium remains unknown. Refer to Appendix A.5 for additional discussion.

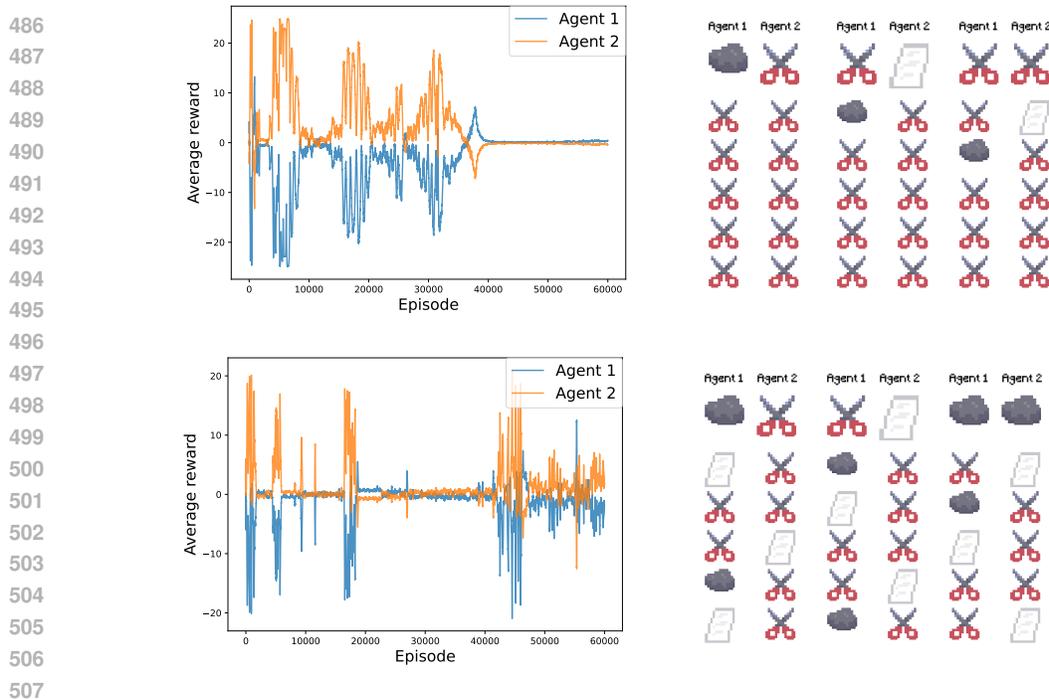


Figure 3: **Saturating rewards (left) versus actions of the learned policies at the end (right) in the rock-paper-scissors game. Top row: *GD-MADDPG*; bottom row: *LA-MADDPG*.** In the left column, blue and orange show the running average of rewards through a window of 100 episodes. In the right column, we depict actions from the respective learned policies evaluation after training is completed, where each row represents what actions players have chosen in one step of the episode. Saturating rewards do not imply good performance, as evidenced by the top row; refer to Section 5.2 for discussion.

6 CONCLUSION

This paper addresses the fundamental optimization challenges in multi-agent reinforcement learning (MARL). By framing MARL as an instance of a Variational Inequality (VI) problem, we highlight its inherent optimization difficulties, which resemble those encountered in solving VIs. These challenges manifest in practice as notoriously difficult training, significant performance variability across random seeds, and other issues that hinder MARL reproducibility, development, and deployment.

To address these challenges, we leverage VI optimization techniques to enhance the convergence and stability of MARL methods. We introduced the *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG* algorithms that combine the multi-agent deep deterministic policy gradient (MADDPG) method with nested Lookahead-VI (Chavdarova et al., 2021), Extragradient (Korpelevich, 1976), and a combination of both, respectively. Our experiments on the *rock-paper-scissors*, *matching pennies* and two *MPE* environments (Lowe et al., 2017) consistently demonstrated the effectiveness of the VI variants of MADDPG in improving performance and stabilizing training compared to the standard baseline method. These findings point toward promising opportunities for further development of VI-based methods in MARL, particularly in leveraging the structure of the MARL optimization landscape.

REFERENCES

- Johann J.H. Ackermann, Volker Gabler, Takayuki Osa, and Masashi Sugiyama. Reducing overestimation bias in multi-agent domains using double centralized critics. *ArXiv:1910.01465*, 2019.
- Wäiss Azizian, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. A tight and unified analysis of gradient-based methods for a whole spectrum of differentiable games. In *AISTATS*, pp. 2863–2873, 2020.

- 540 David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel.
541 The mechanics of n-player differentiable games. In *ICML*, 2018.
- 542
- 543 D. Bertsekas. *Rollout, Policy Iteration, and Distributed Reinforcement Learning*. Athena scientific
544 optimization and computation series. Athena Scientific, 2021. ISBN 9781886529076.
- 545 Matteo Bettini, Amanda Prorok, and Vincent Moens. Benchmark: Benchmarking multi-agent
546 reinforcement learning. *arXiv:2312.01472*, 2024.
- 547
- 548 Radu Ioan Bot, Ernő Robert Csetnek, and Phan Tu Vuong. The forward-backward-forward method
549 from continuous and discrete perspective for pseudo-monotone variational inequalities in Hilbert
550 spaces. *arXiv:1808.08084*, 2020.
- 551 Radu Ioan Bot, Ernő Robert Csetnek, and Dang-Khoa Nguyen. Fast OGDA in continuous and
552 discrete time. *arXiv:2203.10947*, 2022.
- 553
- 554 Michael Bowling. Convergence and no-regret in multiagent learning. In *NIPS*, volume 17. MIT Press,
555 2004.
- 556 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Tight last-iterate convergence of the extragrad-
557 ient method for constrained monotone variational inequalities. *arXiv:2204.09228*, 2022.
- 558
- 559 Tatjana Chavdarova, Gauthier Gidel, François Fleuret, and Simon Lacoste-Julien. Reducing noise in
560 GAN training with variance reduced extragradient. In *NeurIPS*, 2019.
- 561 Tatjana Chavdarova, Matteo Pagliardini, Sebastian U Stich, François Fleuret, and Martin Jaggi.
562 Taming GANs with Lookahead-Minmax. In *ICLR*, 2021.
- 563
- 564 Tatjana Chavdarova, Michael I. Jordan, and Manolis Zampetakis. Last-iterate convergence of saddle
565 point optimizers via high-resolution differential equations. In *Minimax Theory and its Applications*,
566 2023.
- 567 Tatjana Chavdarova, Tong Yang, Matteo Pagliardini, and Michael I. Jordan. A primal-dual approach
568 for solving variational inequalities with general-form constraints. In *ICLR*, 2024.
- 569
- 570 Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Shared experience actor-critic for
571 multi-agent reinforcement learning. *arXiv:2006.07169*, 2021.
- 572
- 573 Richard W. Cottle and George B. Dantzig. Complementary pivot theory of mathematical programming.
574 *Linear Algebra and its Applications*, 1(1):103–125, 1968. ISSN 0024-3795.
- 575
- 576 Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and
577 constrained min-max optimization. In *ITCS*, 2019.
- 578
- 579 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with
580 optimism. In *ICLR*, 2018.
- 581
- 582 Jelena Diakonikolas. Halpern iteration for near-optimal and parameter-free monotone inclusion and
583 strong solutions to variational inequalities. In *COLT*, volume 125, 2020.
- 584
- 585 Theresa Eimer, Marius Lindauer, and Roberta Raileanu. Hyperparameters in reinforcement learning
586 and how to tune them. *arXiv:2306.01324*, 2023.
- 587
- 588 Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional Variational Inequalities and Complementarity Problems*. Springer, 2003.
- 589
- 590 Wentao Fan. A comprehensive analysis of game theory on multi-agent reinforcement. *Highlights in
591 Science, Engineering and Technology*, 85:77–88, 03 2024. doi: 10.54097/gv6fpz53.
- 592
- 593 Gauthier Gidel, Hugo Berard, Pascal Vincent, and Simon Lacoste-Julien. A variational inequality
perspective on generative adversarial nets. In *ICLR*, 2019a.
- Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel Huang,
Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics.
In *AISTATS*, 2019b.

- 594 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates
595 for no-regret learning in multi-player games. In *NeurIPS*, 2020a.
- 596
597 Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is
598 slower than averaged iterate in smooth convex-concave saddle point problems. In *COLT*, pp.
599 1758–1784, 2020b.
- 600 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
601 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- 602
603 Eduard Gorbunov, Nicolas Loizou, and Gauthier Gidel. Extragradient method: $\mathcal{O}(1/K)$ last-iterate
604 convergence for monotone variational inequalities and connections with cocoercivity. In *AISTATS*,
605 2022.
- 606 Rihab Gorsane, Omayma Mahjoub, Ruan de Kock, Roland Dubb, Siddarth Singh, and Arnu Pretorius.
607 Towards a standardised performance evaluation protocol for cooperative marl. *arXiv:2209.10485*,
608 2022.
- 609 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
610 Deep reinforcement learning that matters. *arXiv:1709.06560*, 2019.
- 611
612 Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *ICML*,
613 2018.
- 614
615 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- 616 Vijay Konda and John Tsitsiklis. Actor-critic algorithms. In S. Solla, T. Leen, and K. Müller (eds.),
617 *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999.
- 618
619 Galina Michailovna Korpelevich. The extragradient method for finding saddle points and other
620 problems. *Matecon*, 1976.
- 621
622 Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Master’s thesis, 2009.
- 623 Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong
624 Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv:2109.11251*,
625 2022.
- 626
627 Marc Lanctot, John Schultz, Neil Burch, Max Olan Smith, Daniel Hennes, Thomas Anthony, and
628 Julien Perolat. Population-based evaluation in repeated rock-paper-scissors as a benchmark for
629 multiagent reinforcement learning. *arXiv:2303.03196*, 2023.
- 630 Pengyi Li, Jianye Hao, Hongyao Tang, Yan Zheng, and Xian Fu. Race: improve multi-agent
631 reinforcement learning with representation asymmetry and collaborative evolution. In *Proceedings*
632 *of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- 633
634 Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence
635 of generative adversarial networks. *AISTATS*, 2019.
- 636
637 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Manfred Otto Heess, Tom Erez,
638 Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement
639 learning. *CoRR*, abs/1509.02971, 2015.
- 640
641 Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-
642 critic for mixed cooperative-competitive environments. *NIPS*, 2017.
- 643
644 Nicolai A. Lynnerup, Laura Nolling, Rasmus Hasle, and John Hallam. A survey on reproducibility by
645 evaluating deep reinforcement learning algorithms on real-world robots. *arXiv:1909.03772*, 2019.
- 646
647 Yu. Malitsky. Projected reflected gradient methods for monotone variational inequalities. *SIAM*
Journal on Optimization, 25:502–520, 2015.
- Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which Training Methods for GANs do
actually Converge? In *ICML*, 2018.

- 648 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-
649 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,
650 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
651 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.
652 *Nature*, 518:529–533, 2015.
- 653 Dov Monderer and Lloyd S Shapley. Potential games. *Games and economic behavior*, 1996.
- 654 Shayegan Omidshafiei, Jason Pazis, Chris Amato, Jonathan P. How, and John Vian. Deep de-
655 centralized multi-task multi-agent reinforcement learning under partial observability. In *ICML*,
656 2017.
- 658 Thomas Pethick, Wanyun Xie, and Volkan Cevher. Stable nonconvex-nonconcave training via linear
659 interpolation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 660 Leonid Denisovich Popov. A modification of the arrow–hurwicz method for search of saddle points.
661 *Mathematical Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- 662 Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep
663 convolutional generative adversarial networks. In *ICLR*, 2016.
- 665 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder de Witt, Gregory Farquhar, Jakob Foerster,
666 and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent
667 reinforcement learning. *arXiv:1803.11485*, 2018.
- 668 Ralph Tyrrell Rockafellar. Monotone operators associated with saddle-functions and minimax
669 problems. *Nonlinear functional analysis*, 18(part 1):397–407, 1970.
- 670 Mihaela Rosca, Yan Wu, Benoit Dherin, and David G. T. Barrett. Discretization drift in two-player
671 games. In *ICML*, 2021.
- 673 Ernest K. Ryu, Kun Yuan, and Wotao Yin. Ode analysis of stochastic gradient methods with optimism
674 and anchoring for minimax problems. *arXiv:1905.10899*, 2019.
- 675 Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli,
676 Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The
677 starcraft multi-agent challenge. *arXiv:1902.04043*, 2019.
- 679 Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran:
680 Learning to factorize with transformation for cooperative multi-agent reinforcement learning.
681 *arXiv:1905.05408*, 2019.
- 682 Riccardo Spica, Davide Falanga, Eric Cristofalo, Eduardo Montijano, Davide Scaramuzza, and
683 Mac Schwager. A real-time game theoretic planner for autonomous two-player drone racing.
684 *arXiv:1801.02302*, 2018.
- 685 Guido Stampacchia. Formes bilineaires coercitives sur les ensembles convexes. *Académie des
686 Sciences de Paris*, 258:4413–4416, 1964.
- 688 J Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S
689 Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym
690 for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:
691 15032–15043, 2021.
- 692 Kiran Koshy Thekumparampil, Niao He, and Sewoong Oh. Lifted primal-dual method for bilinearly
693 coupled smooth minimax optimization. In *AISTATS*, 2022.
- 694 Paul Tseng. On linear convergence of iterative methods for the variational inequality problem.
695 *Journal of Computational and Applied Mathematics*, 60:237–252, 1995. ISSN 0377-0427.
- 697 Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle
698 Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, John Quan, Stephen
699 Gaffney, Stig Petersen, Karen Simonyan, Tom Schaul, Hado van Hasselt, David Silver, Tim-
700 othy Lillicrap, Kevin Calderone, Paul Keet, Anthony Brunasso, David Lawrence, Anders Ek-
701 ermo, Jacob Repp, and Rodney Tsing. Starcraft ii: A new challenge for reinforcement learning.
arXiv:1708.04782, 2017.

702 Han Wang, Archit Sakhadeo, Adam White, James Bell, Vincent Liu, Xutong Zhao, Puer Liu, Tadashi
703 Kozuno, Alona Fyshe, and Martha White. No more pesky hyperparameters: Offline hyperparameter
704 tuning for rl. *arXiv:2205.08716*, 2022.
705
706 Zhijian Wang, Bin Xu, and Hai-Jun Zhou. Social cycling and conditional responses in the rock-paper-
707 scissors game. *Scientific Reports*, 4(1), 2014. ISSN 2045-2322. doi: 10.1038/srep05830.
708
709 Tong Yang, Michael I. Jordan, and Tatjana Chavdarova. Solving constrained variational inequalities
710 via an interior point method. In *ICLR*, 2023.
711
712 Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game
713 theoretical perspective. *arXiv:2011.00583*, 2021.
714
715 Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre M. Bayen, and Yi Wu. The surprising
716 effectiveness of mappo in cooperative, multi-agent games. *ArXiv:2103.01955*, 2021.
717
718 Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The
719 surprising effectiveness of ppo in cooperative, multi-agent games. *arXiv:2103.01955*, 2022.
720
721 Michael Zhang, James Lucas, Jimmy Ba, and Geoffrey E Hinton. Lookahead optimizer: k steps
722 forward, 1 step back. In *NeurIPS*, 2019.
723
724 Liyuan Zheng, Tanner Fiez, Zane Alumbaugh, Benjamin Chasnov, and Lillian J. Ratliff. Stackelberg
725 actor-critic: Game-theoretic reinforcement learning algorithms. *arXiv:2109.12286*, 2021.
726
727 Hai-Jun Zhou. The rock–paper–scissors game. *Contemporary Physics*, 57(2):151–163, March 2015.
728 ISSN 1366-5812. doi: 10.1080/00107514.2015.1026556.
729
730 M. Zhou, Y. Guan, M. Hayajneh, K. Niu, and C. Abdallah. Game theory and machine learning in
731 uavs-assisted wireless communication networks: A survey. *arXiv:2108.03495*, 2021.
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX

A.1 ADDITIONAL BACKGROUND

A.1.1 VI CLASSES AND ADDITIONAL METHODS

The following VI class is often referred to as the generalized class for VIs to that of convexity in minimization.

Definition 1 (monotonicity) An operator $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is monotone if $\langle z - z', F(z) - F(z') \rangle \geq 0$, $\forall z, z' \in \mathbb{R}^d$. F is μ -strongly monotone if: $\langle z - z', F(z) - F(z') \rangle \geq \mu \|z - z'\|^2$ for all $z, z' \in \mathbb{R}^d$.

In addition to those presented in the main part, we describe the following popular VI method.

Optimistic Gradient Descent (OGD). The update rule of Optimistic Gradient Descent OGD ((OGD) Popov, 1980) is:

$$z_{t+1} = z_t - 2\eta F(z_t) + \eta F(z_{t-1}), \quad (\text{OGD})$$

where $\eta \in (0, 1)$ is the learning rate.

A.1.2 PSEUDOCODE FOR NESTED LOOKAHEAD FOR A TWO-PLAYER GAME

For completeness, in Algorithm 3 we give the details of the nested Lookahead-Minmax algorithm proposed in (Algorithm 6, Chavdarova et al., 2021) with two-levels.

Algorithm 3 Pseudocode of Two-Level Nested Lookahead–Minmax. (Chavdarova et al., 2021)

```

1: Input: Stopping time  $T$ , learning rates  $\eta_\theta, \eta_\varphi$ , initial weights  $\theta, \varphi$ , lookahead hyperparameters
    $k_s, k_{ss}$  and  $\alpha$ , losses  $\mathcal{L}^\theta, \mathcal{L}^\varphi$ , update ratio  $r$ , real-data distribution  $p_d$ , noise-data distribution  $p_z$ .
2:  $(\theta_s, \theta_{ss}, \varphi_s, \varphi_{ss}) \leftarrow (\theta, \theta, \varphi, \varphi)$  (store copies for slow and super-slow)
3: for  $t \in 1, \dots, T$  do
4:   for  $i \in 1, \dots, r$  do
5:      $x \sim p_d, z \sim p_z$ 
6:      $\varphi \leftarrow \varphi - \eta_\varphi \nabla_\varphi \mathcal{L}^\varphi(\theta, \varphi, x, z)$  (update  $\varphi$   $r$  times)
7:   end for
8:    $z \sim p_z$ 
9:    $\theta \leftarrow \theta - \eta_\theta \nabla_\theta \mathcal{L}^\theta(\theta, \varphi, z)$  (update  $\theta$  once)
10:  if  $t \% k_s == 0$  then
11:     $\varphi \leftarrow \varphi_s + \alpha_\varphi(\varphi - \varphi_s)$  (backtracking on interpolated line  $\varphi_s, \varphi$ )
12:     $\theta \leftarrow \theta_s + \alpha_\theta(\theta - \theta_s)$  (backtracking on interpolated line  $\theta_s, \theta$ )
13:     $(\theta_s, \varphi_s) \leftarrow (\theta, \varphi)$  (update slow checkpoints)
14:  end if
15:  if  $t \% k_{ss} == 0$  then
16:     $\varphi \leftarrow \varphi_{ss} + \alpha_\varphi(\varphi - \varphi_{ss})$  (backtracking on interpolated line  $\varphi_{ss}, \varphi$ )
17:     $\theta \leftarrow \theta_{ss} + \alpha_\theta(\theta - \theta_{ss})$  (backtracking on interpolated line  $\theta_{ss}, \theta$ )
18:     $(\theta_{ss}, \varphi_{ss}) \leftarrow (\theta, \varphi)$  (update super-slow checkpoints)
19:     $(\theta_s, \varphi_s) \leftarrow (\theta, \varphi)$  (update slow checkpoints)
20:  end if
21: end for
22: Output:  $\theta_{ss}, \varphi_{ss}$ 

```

A.1.3 DETAILS ON THE MADDPG ALGORITHM

The MADDPG algorithm is outlined in Algorithm 4. An empty replay buffer \mathcal{D} is initialized to store experiences (line 3). In each episode, the environment is reset and experiences in the form of (state, action, reward, next state) are saved to \mathcal{D} . After a predetermined number of random iterations, learning begins by sampling batches from \mathcal{D} .

The critic of agent i receives the sampled joint actions \mathbf{a} of all agents and the state information of agent i to output the predicted Q_i -value of agent i . Deep Q-learning (Mnih et al., 2015) is then used to update the critic network; lines 21-22. Then, the agents’ policy network is optimized using policy gradient; refer to 24. Finally, following each learning iteration, the target networks are updated towards current actor and critic networks using a fraction τ .

All networks are optimized using the Adam optimizer (Kingma & Ba, 2015). Once training is complete, each agent’s actor operates independently during execution. This approach is applicable across cooperative, competitive, and mixed environments.

Algorithm 4 Pseudocode for MADDPG (Lowe et al., 2017).

```

1: Input: Environment  $\mathcal{E}$ , number of agents  $N$ , number of episodes  $T$ , action spaces  $\{\mathcal{A}_i\}_{i=1}^N$ ,
2: number of random steps  $T_{\text{rand}}$  before learning, learning interval  $T_{\text{learn}}$ , actor networks  $\{\mu_i\}_{i=1}^N$ ,
3: with initial weights  $\theta \equiv \{\theta_i\}_{i=1}^N$ , critic networks  $\{Q_i\}_{i=1}^N$  with initial weights  $w \equiv \{w_i\}_{i=1}^N$ ,
4: learning rates  $\eta_\theta, \eta_w$ , optimizer  $\mathcal{B}$  (e.g., Adam), discount factor  $\gamma$ , soft update parameter  $\tau$ .
5: Initialize:
6: 3: Replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
7: 4: for all episode  $e \in 1, \dots, T$  do
8: 5:  $\mathbf{x} \leftarrow \text{Sample}(\mathcal{E})$  (sample from environment  $\mathcal{E}$ )
9: 6:  $step \leftarrow 1$ 
10: 7: repeat
11: 8: if  $e \leq T_{\text{rand}}$  then
12: 9: for each agent  $i, a_i \sim \mathcal{A}_i$  (sample actions randomly)
13: 10: else
14: 11: for each agent  $i$ , select action  $a_i = \mu_i(o_i) + \mathcal{N}_t$  using current policy and exploration
15: 12: end if
16: 13: (apply actions and record results)
17: 14: Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$ , observe rewards  $\mathbf{r}$  and new state  $\hat{\mathbf{x}}$ 
18: 15: replay buffer  $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \hat{\mathbf{x}})$ 
19: 16:  $\mathbf{x} \leftarrow \hat{\mathbf{x}}$ 
20: 17: (apply learning step if applicable)
21: 18: if  $step \% T_{\text{learn}} = 0$  then
22: 19: for all agent  $i \in 1, \dots, N$  do
23: 20: sample batch  $\{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \hat{\mathbf{x}}^j)\}_{j=1}^B$  of size  $B$  from  $\mathcal{D}$ 
24: 21:  $y^j \leftarrow r_i^j + \gamma \mathbf{Q}^{\mu'}(\hat{\mathbf{x}}^j, a_1^j, \dots, a_N^j)$ , where  $\mathbf{a}'_k = \{\mu'_k(o_k^j)\}$ 
25: 22: Update critic by minimizing the loss (using optimizer  $\mathcal{B}$ ):
26: 23: 
$$\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j (y^j - \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_N^j))^2$$

27: 24: Update actor policy using policy gradient formula and optimizer  $\mathcal{B}$ 
28: 25:  $\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(o_i^j) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_N^j)$ , where  $a_i = \mu_i(o_i^j)$ 
29: 26: end for
30: 27: for all agent  $i \in 1, \dots, N$  do
31: 28:  $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$  (update target networks)
32: 29:  $w'_i \leftarrow \tau w_i + (1 - \tau) w'_i$ 
33: 30: end for
34: 31: end if
35: 32:  $step \leftarrow step + 1$ 
36: 33: until environment terminates
37: 34: end for
38: 34: Output:  $\theta, w$ 

```

A.1.4 EXTENDED VERSION OF LA-MADDPG PSEUDOCODE

We include an extended version for the LA-MADDPG algorithm without VI notations in algorithm 5.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Algorithm 5 Pseudocode for LA-MADDPG: MADDPG with (Nested) Lookahead.

```

1: Input: Environment  $\mathcal{E}$ , number of agents  $N$ , number of episodes  $T$ , action spaces  $\{\mathcal{A}_i\}_{i=1}^N$ , number
of random steps  $T_{\text{rand}}$  before learning, learning interval  $T_{\text{learn}}$ , actor networks  $\{\mu_i\}_{i=1}^N$ , with
initial weights  $\theta \equiv \{\theta_i\}_{i=1}^N$ , critic networks  $\{Q_i\}_{i=1}^N$  with initial weights  $w \equiv \{w_i\}_{i=1}^N$ , learning
rates  $\eta_\theta, \eta_w$ , base optimizer  $\mathcal{B}$  (e.g., Adam), discount factor  $\gamma$ , lookahead hyperparameters
 $k_s, k_{ss}$  (where  $k_{ss}$  can be  $\emptyset$ ) and  $\alpha_\theta, \alpha_w$ , soft update parameter  $\tau$ .
2: Initialize:
3:   Replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
4:   for all agent  $i \in 1, \dots, N$  do
5:      $(\theta_{i,s}, \theta_{i,ss}, w_{i,s}, w_{i,ss}) \leftarrow (\theta_i, \theta_i, \theta_i, w_i, w_i, w_i)$ 
6:     (store snapshots for nLA)
7:   end for
8: for all episode  $e \in 1, \dots, T$  do
9:    $\mathbf{x} \leftarrow \text{Sample}(\mathcal{E})$  (sample from environment  $\mathcal{E}$ )
10:   $\text{step} \leftarrow 1$ 
11:  repeat
12:    if  $e \leq T_{\text{rand}}$  then
13:      for each agent  $i, a_i \sim \mathcal{A}_i$  (sample actions randomly)
14:    else
15:      for each agent  $i$ , select action  $a_i$  using current policy and exploration
16:    end if
17:    (apply actions and record results)
18:    Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$ , observe rewards  $\mathbf{r}$  and new state  $\hat{\mathbf{x}}$ 
19:    replay buffer  $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \hat{\mathbf{x}})$ 
20:     $\mathbf{x} \leftarrow \hat{\mathbf{x}}$ 
21:    (apply learning step if applicable)
22:    if  $\text{step} \% T_{\text{learn}} = 0$  then
23:      for all agents  $i \in 1, \dots, N$  do
24:        sample batch  $\{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \hat{\mathbf{x}}^j)\}_{j=1}^B$  of size  $B$  from  $\mathcal{D}$ 
25:         $y^j \leftarrow r_i^j + \gamma \mathbf{Q}^{\mu'}(\hat{\mathbf{x}}^j, a'_1, \dots, a'_N)$ , where  $\mathbf{a}'_k = \{\mu'_k(\sigma_k^j)\}$ 
26:        Update critic by minimizing the loss  $\mathcal{L}(w_i) = \frac{1}{S} \sum_j (y^j - \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_N^j))^2$ 
        using  $\mathcal{B}$ 
27:        Update actor policy using policy gradient formula and  $\mathcal{B}$ 
28:         $\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(\sigma_i^j) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_N^j)$ , where  $a_i = \mu_i(\sigma_i^j)$ 
29:      end for
30:      for all agents  $i \in 1, \dots, N$  do
31:         $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$  (update target networks)
32:         $w'_i \leftarrow \tau w_i + (1 - \tau) w'_i$ 
33:      end for
34:    end if
35:     $\text{step} \leftarrow \text{step} + 1$ 
36:  until environment terminates
37:  NESTEDLOOKAHEAD( $N, e, \Theta, \mathbf{W}, k_s, k_{ss}, \alpha_\theta, \alpha_w$ )
38:  where:
39:     $\Theta = \{(\theta_i, \theta_{i,s}, \theta_{i,ss})\}_{i=1}^N$  (all actor weights and snapshots)
40:     $\mathbf{W} = \{(w_i, w_{i,s}, w_{i,ss})\}_{i=1}^N$  (all critic weights and snapshots)
41: end for
42: Output:  $\theta, w$ 

```

A.1.5 PSEUDOCODE FOR EXTRAGRADIENT

In Algorithm 6 outlines the *Extragradient* optimizer (Korpelevich, 1976), which we employ in EG-MADDPG. This method uses a gradient-based optimizer to compute the extrapolation iterate, then applies the gradient at the extrapolated point to perform an actual update step. The extragradient optimizer is used to update all agents’ actor and critic networks. In our experiments, we use Adam for both the extrapolation and update steps, maintaining the same learning intervals and parameters as in the baseline algorithm.

Algorithm 6 Extragradient optimizer; Can be used as \mathcal{B} in algorithm 2.

```

1: Input: learning rate  $\eta_\varphi$ , initial weights  $\varphi$ , loss  $\mathcal{L}^\varphi$ , extrapolation steps  $t$ 
2:  $\varphi^{copy} \leftarrow \varphi$  (Save current parameters)
3: for  $i \in 1, \dots, t$  do
4:    $\varphi = \varphi - \eta_\varphi \nabla_\varphi \mathcal{L}^\varphi(\varphi)$  (Compute the extrapolated  $\varphi$ )
5: end for
6:  $\varphi = \varphi^{copy} - \eta_\varphi \nabla_\varphi \mathcal{L}^\varphi(\varphi)$  (update  $\varphi$ )
7: Output:  $\varphi$ 

```

A.2 DETAILS ON THE IMPLEMENTATION

As mentioned earlier, we followed the configurations and hyperparameters from the original MADDPG paper for our implementation. For completeness, these are listed in Table 2. We ran $T = 60000$ for all environments except Matching Pennies where we ran for 50000 training episodes, with a maximum of 25 environment steps (s) per episode.

In all Rock-Paper-Scissors and Matching pennies experiments, we used a 2-layer MLP with 64 units per layer, while for MPE: Predator-prey, we used a 2-layer MLP with 128 units per layer. ReLU activation was applied between layers for both the policy and value networks of all agents.

A.2.1 HYPERPARAMETER SELECTION FOR NESTED-LOOKAHEAD

In this section, we discuss and share guidelines for hyperparameter selection based on our experiments.

Summary.

- We observed two- or three-level of nested-Lookahead outperform single-level Lookahead.
- Each level has different k , denoted here with k_s, k_{ss}, k_{sss} as in the main part. These should be selected as multiple of the selected k for the level before, that is, $k_{ss} \equiv c_{ss}k_s$, and $k_{sss} \equiv c_{sss}k_{ss}$, where c_{ss}, c_{sss} are positive integers.
- We observed that for the innermost lookahead, small values for k_s , such as smaller than 50, perform better than using large values. For the outer k_{ss}, k_{sss} large values work well, such as in the range between 5 – 10 for the c_{ss}, c_{sss} .
- We typically used $\alpha = 0.5$, and we observed lower values, such as $\alpha = 0.3$, give better performances than $\alpha > 0.5$.

Discussion.

- To give an intuition regarding the above-listed conclusions, small values for k_s help because the MARL setting is very noisy and the vector field is rotational. If large values are used for k_s , then the algorithm will diverge away. It is known that the combination of noise and rotational vector field can cause methods to diverge away (Chavdarova et al., 2019).
- Relative to the analogous conclusions for GANs (Chavdarova et al., 2021), the differences is that:
 - The better-performing values for k_s are of a similar range as for Lookahead with GD for GANs; however they are smaller than those used for Lookahead with EG for GANs.

Table 2: Hyperparameters used for LA-MADDPG experiments.

Name	Description
Adam lr	0.01
Adam β_1	0.9
Adam β_2	0.999
Batch-size	1024
Update ratio τ	0.01
Discount factor γ	0.95
Replay Buffer	10^6
learning step T_{learn}	100
T_{rand}	1024
Lookahead α	0.5

A.3 ADDITIONAL RESULTS

A.3.1 ROCK-PAPER-SCISSORS: BUFFER STRUCTURE

For the Rock-Paper-Scissors (RPS) game, using a buffer size of 1M wasn’t sufficient to store all experiences from the 60K training episodes. We observed a change in algorithm behavior around 40K episodes. To explore the impact of buffer configurations, we experimented with different sizes and structures, as experience storage plays a critical role in multi-agent reinforcement learning.

Full buffer. The buffer is configured to store all experiences from the beginning to the end of training without any loss.

Buffer clearing. In this setup, a smaller buffer is used, and once full, the buffer is cleared completely, and new experiences are stored from the start.

Buffer shifting. Similar to the small buffer setup, but once full, old experiences are replaced by new ones in a first-in-first-out (FIFO) manner.

Results. Figure 4 depicts the results when using different buffer options for the RPS game.

A.3.2 ROCK-PAPER-SCISSORS: SCHEDULED LEARNING RATE

We experimented with gradually decreasing the learning rate (LR) during training to see if it would aid convergence to the optimal policy in RPS. While this approach reduced noise in the results, it also led to increased variance across all methods except for LA-MADDPG.

Figure 5 depicts the average distance to the equilibrium policy over 5 different seeds for each methods, using periodically decreased step sizes.

A.3.3 MPE: PREDATOR-PREY FULL RESULTS

While in the main part in Figure 2 we showed only two methods for clarity, Figure 6 depicts all methods.

We also evaluated the trained models of all methods on an instance of the environment that runs for 50 steps to compare learned policies. We present snapshots from it in Figure 7. Here, you can clearly anticipate the difference between the policies from baseline and our optimization methods. As in the baseline, only one agent will chase at the beginning of episode. Moreover, for the baseline (topmost row), the agents move further away from the landmarks and the good agent, which is suboptimal. This can be noticed from the decreasing agents’ size in the figures. While in ours, both adversary agents engage in chasing the good agent until the end.

A.3.4 MPE: PREDATOR-PREY AND PHYSICAL DECEPTION TRAINING FIGURES

In figures 8a and 8b we include the rewards achieved during the training of GD-MADDPG and LA-MADDPG resp. for MPE: Predator-prey. The figures show individual rewards for the agent

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

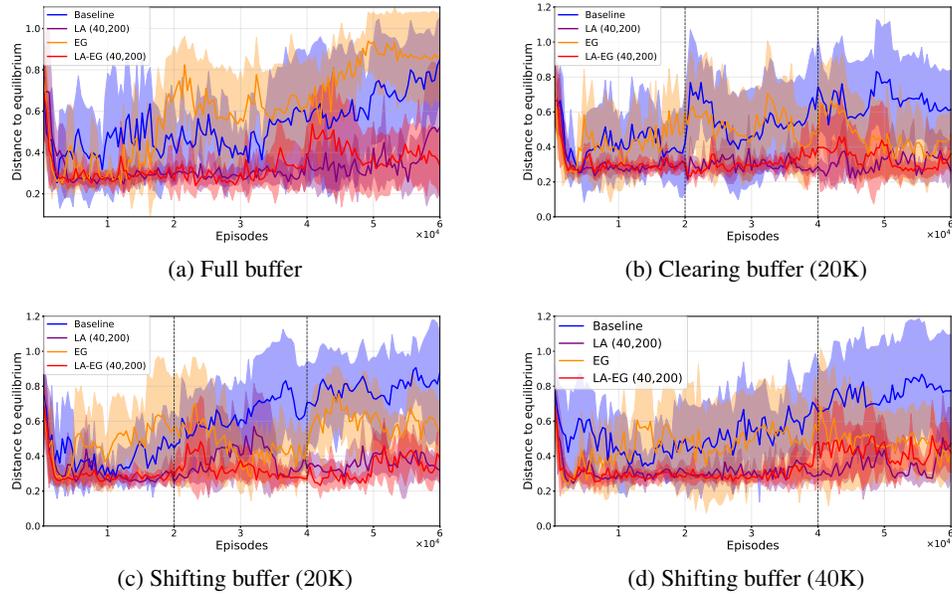


Figure 4: Comparison of different buffer configurations (see Appendix A.3.1) and methods on Rock-paper-scissors game. x -axis: training episodes. y -axis: 5-seed average norm between the two players’ policies and equilibrium policy $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$. The dotted line indicates the point at which the buffer begins to change, either through shifting or clearing.

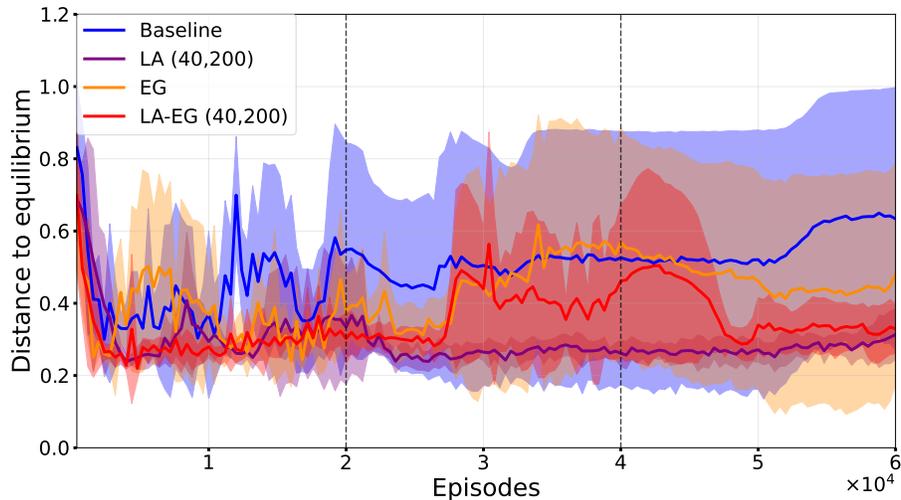
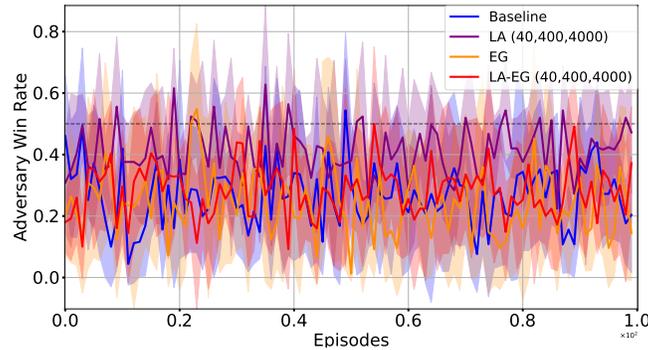


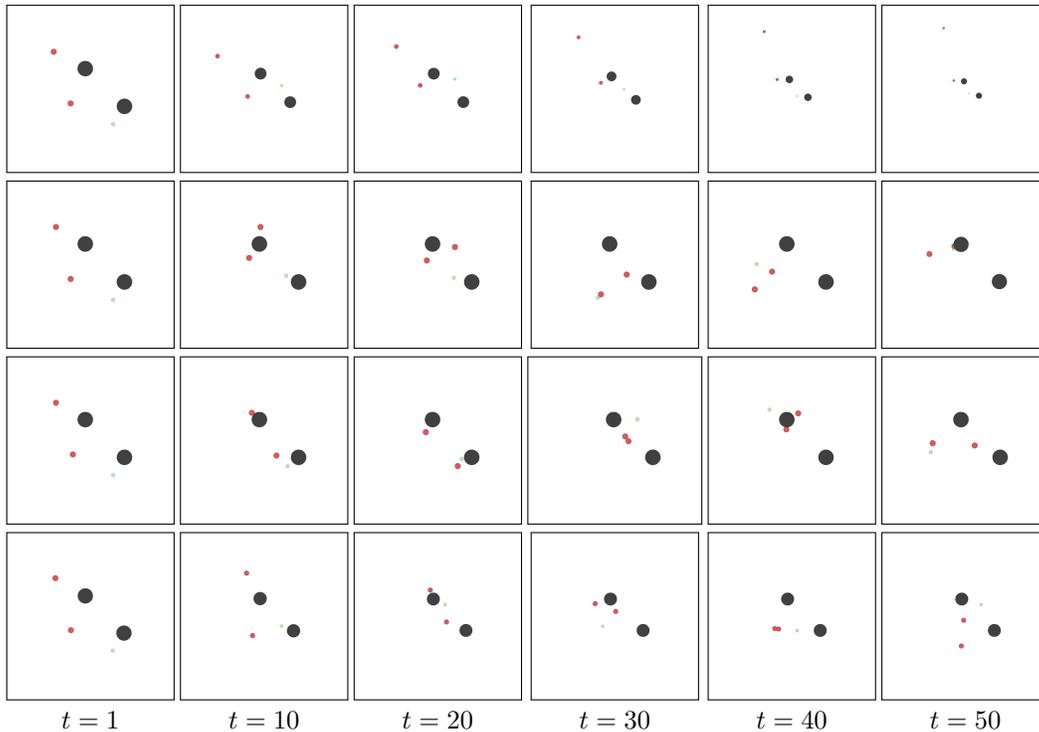
Figure 5: Compares MADDPG with different LA-MADDPG configurations to the baseline MADDPG with (Adam) in rock-paper-scissors. x -axis: training episodes. y -axis: 5-seed average norm between the two players’ policies and equilibrium policy $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$. The dotted lines depict the times when the learning rate was decreased by a factor of 10.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092



1093 **Figure 6: Comparison on the MPE–Predator-prey game between the *GD-MADDPG*, *LA-MADDPG*, *EG-***
1094 ***MADDPG* and *LA-EG-MADDPG* optimization methods, denoted as *Baseline*, *LA*, *EG*, *LA-EG*, resp. x -axis:**
1095 **evaluation episodes. y -axis: mean adversaries win rate, averaged over 5 runs with different seeds.**

1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120



1121 **Figure 7: Agents' trajectories of fully trained models with all considered optimization methods on the**
1122 **same environment seed of MPE: Predator-prey.** Snapshots show the progress of agents as time progresses in a
1123 50 steps long environment. Each row contains snapshots of one method, from top to bottom: *GD-MADDPG*,
1124 *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG*. Big dark circles represent landmarks, small red circles are
1125 adversary agents and green one is the good agent.

1126
1127
1128
1129
1130
1131

(prey) and one adversary (predator). Blue and green show the individual rewards received at each episode while the orange and red lines are the respective running averages with window size of 100 of those rewards.

1132 Figures 9a and 9b demonstrate same results but for MPE: Physical deception. In this game, We have
1133 two good agents, 'Agent 0 and 1' but since they are both receive same rewards, we only show agent 0.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

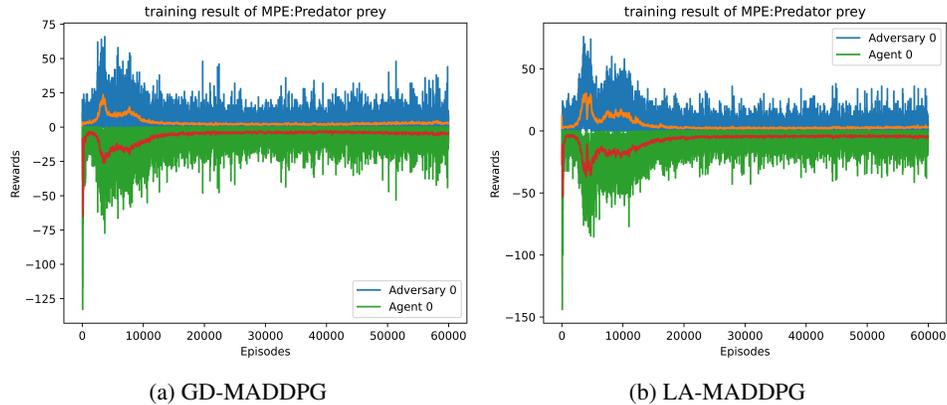


Figure 8: The figure shows the learning curves during training of GD-MADDPG and LA-MADDPG for MPE: Predator-Prey. x -axis: training episodes. y -axis: agents rewards and their moving average with a window size of 100, calculated over 5-seeds over 5 seeds.

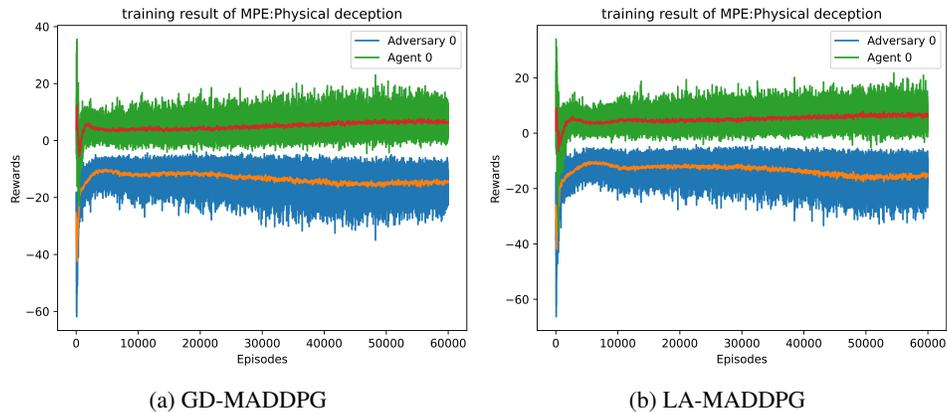


Figure 9: The figure shows the learning curves during training of GD-MADDPG and LA-MADDPG for MPE: Physical deception. x -axis: training episodes. y -axis: agents rewards and their moving average with a window size of 100, calculated over 5-seeds over 5 seeds.

A.4 MATD3 EXPERIMENTS

We compared Multi-agent TD3 (MATD3), (Ackermann et al., 2019) with MADDPG on the MPE benchmark, *Physical deception*. From Figure 10, we observe that the performances are similar: both algorithms fluctuate between high and low rewards. Hence, optimization methods dealing with rotational dynamics would benefit both.

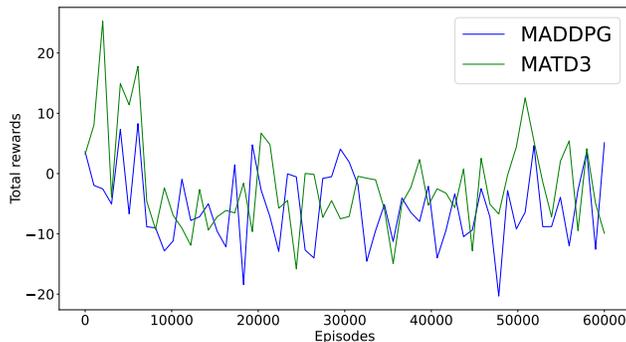


Figure 10: Comparison on the MPE-Physical deception game between the MADDPG, and MATD3 algorithms. x -axis: training episodes. y -axis: total rewards.

A.5 ON THE REWARDS AS CONVERGENCE METRIC

Based on our experiments and findings from the multi-agent literature (Bowling, 2004), we observe that average rewards offer a weaker measure of convergence compared to policy convergence in multi-agent games. This implies that rewards can reach a target value even when the underlying policy is suboptimal. For example, in the Rock-paper-scissors game, the Nash equilibrium policy leads to nearly equal wins for both players, resulting in a total reward of zero. However, this same reward can also be achieved if one player always wins while the other consistently loses, or if both players repeatedly select the same action, leading to a tie. As such, relying solely on rewards during training can be misleading.

Figure 3 (top row) depicts a case with the baseline where, despite rewards converging during training, the agents ultimately learned to play the same action repeatedly, resulting in ties. Although this matched the expected reward, it falls far short of equilibrium and leaves the agents vulnerable to exploitation by more skilled opponents. In contrast, the same figure shows results from LA-MADDPG under the same experimental conditions. Notably, while the rewards did not fully converge, the agents learned a near-optimal policy during evaluation, alternating between all three actions as expected. These results also align with the findings shown in Figure 1a.

We explored the use of gradient norms as a potential metric in these scenarios but found them to be of limited utility, as they provided no clear indication of convergence for either method. We include those results in Figure 11, where we compare the gradient norms of Adam and LA across the networks of different players.

This work highlights the need for more robust evaluation metrics in multi-agent reinforcement learning, a point also emphasized in (Lanctot et al., 2023), as reward-based metrics alone may be inadequate, particularly in situations where the true equilibrium is unknown.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

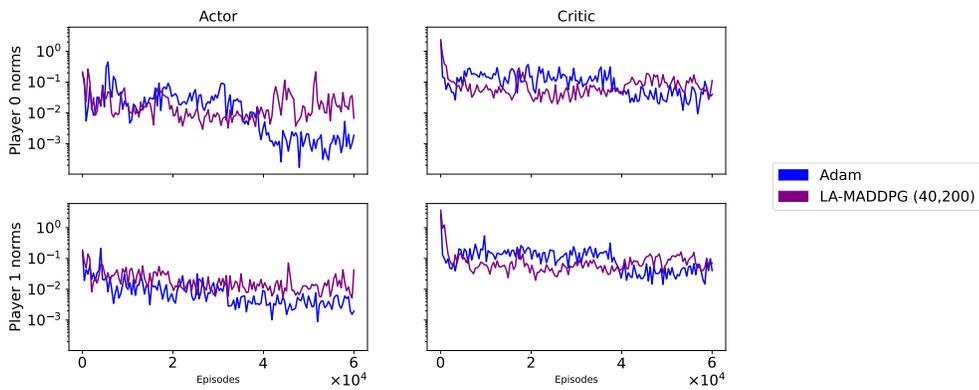


Figure 11: Gradient norms across training in the *rock-paper-scissors* game.