

---

# Variational Inequality Perspective and Optimizers for Multi-Agent Reinforcement Learning

---

**Baraah A. M. Sidahmed**

CISPA Helmholtz Center for Information Security  
Universität des Saarlandes  
baraah.adil@cispa.de

**Tatjana Chavdarova**

University of California, Berkeley  
CISPA Helmholtz Center for Information Security  
tatjana.chavdarova@berkeley.edu

## Abstract

Multi-agent reinforcement learning (MARL) presents unique challenges as agents learn strategies through trial and error. Gradient-based methods are often sensitive to hyperparameter selection and initial random seed variations. Recently, progress has been made in solving problems modeled by 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 performance of VI methods—namely, Nested-Lookahead VI (nLA-VI) and Extragradient (EG)—in enhancing 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 demonstrate 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 *MPE Predator-prey* environment [Lowe et al., 2017], where VI-based methods also foster more balanced participation among agents on the same team.

## 1 Introduction

We focus on multi-agent reinforcement learning (MARL) where multiple agents interact and learn simultaneously. MARL is often used to address complex, multi-agent problems across diverse domains, such as coordinating multi-robot and multi-drone systems for tasks like 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, agents aim to optimize a shared objective while acting on their own policies based on their observation of the environment. The interactions among agents can be competitive, cooperative, or a mix of both, leading to complex learning dynamics that differ significantly from single-agent reinforcement learning.

A significant challenge in single-agent reinforcement learning is the high sensitivity to hyperparameter selection and variations in initial random seeds. These methods require careful hyperparameter tuning, and performance can vary drastically with different random seeds used for sampling and model initialization [Wang et al., 2022, Eimer et al., 2023]. Consequently, the results are often not reproducible [Henderson et al., 2019], complicating research efforts and the deployment of these

algorithms. This issue propagates to multi-agent reinforcement learning (MARL) when single-agent methods are extended to the MARL framework. The work of [Gorsane et al. \[2022\]](#) highlights how a small change in the choice of hyperparameter values can lead to a significant difference in algorithm results. The same work also addresses that the results obtained from different seeds in popular MARL benchmarks like StarCraft multi-agent challenge [[Samvelyan et al., 2019](#)] exhibit a high variance. Furthermore, iterative gradient-based approaches for MARL have difficulty in exploring joint policy spaces of the multiple agents [[Li et al., 2023](#), [Christianos et al., 2021](#)] which may lead to sub-optimal solutions. Moreover, inherent cycling effects [[Zheng et al., 2021](#)] of some MARL structures is an added challenge for those approaches.

A concurrent line of works focuses on the Variational Inequality (VI) problem, a general class of problems that encompasses both equilibria and optima. VIs generalize standard constrained minimization problems, where  $F$  is a gradient field  $F \equiv \nabla f$ , and, by allowing  $F$  to be a general vector field, they also include problems such as finding equilibria in zero-sum and general-sum games [[Cottle and 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 simple problem instances of VIs. This failure is primarily due to the *rotational* component inherent in the gradient dynamics of these settings, which leads 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 and Shapley, 1996](#)] and a *Hamiltonian* 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 and Ba, 2015](#)] have no hope of converging for a more general problem class. 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 often diverges away [[Chavdarova et al., 2021](#)], unlike in standard minimization. Since practical GAN implementations are often not zero-sum, VIs allow for their modeling. 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.

In this paper, we pose the following question:

*Do MARL algorithms gain advantages from using VI 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 *nested-Lookahead-VI* (nLA-VI) [[Chavdarova et al., 2021](#)] and *Extragradient* (EG) [[Korpelevich, 1976](#)] methods for solving variational inequalities (VIs). In summary, our main contributions are as follows:

- 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 empirically compare our proposed methods to standard optimization methods in several two-player games and some MPE [[Lowe et al., 2017](#)] benchmarks.
- We provide further insights on using rewards as a metric within the MARL setting.

## 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, 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 and 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 and 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 and Pang, 2003a, 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 and Pang, 2003a, Daskalakis et al., 2018, Liang and Stokes, 2019, Gidel et al., 2019, 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 and Panageas, 2019, Cai et al., 2022, 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  $\pi_i$ , 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., 2019] 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*  $\mu_i$  which represents the policy to be updated, and (iv) a target actor network  $\mu'_i$  from which it selects its actions and is periodically updated with the learned policy  $\mu_i$ . 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 and Pang, 2003b] 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:  $F(\mathbf{z}) \equiv [\nabla_{\mathbf{z}_1} f_1(\mathbf{z}), \dots, \nabla_{\mathbf{z}_N} f_N(\mathbf{z})]^\top$ .

---

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

---

```

1: procedure NESTEDLOOKAHEAD:
2:   Input: Number of agents  $N$ , current episode  $e$ , current actor weights and snapshots
       $\{(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{i,s}, \boldsymbol{\theta}_{i,ss})\}_{i=1}^N$ , current critic weights and snapshots  $\{(\mathbf{w}_i, \mathbf{w}_{i,s}, \mathbf{w}_{i,ss})\}_{i=1}^N$ , lookahead
      hyperparameters  $k_s, k_{ss}$  (where  $k_{ss}$  can be  $\emptyset$ ) and  $\alpha_\theta, \alpha_w$ 
3:   Result: Updated actor and critic weights and snapshots for all agents
4:   if  $e \% k_s == 0$  then
5:     for all agent  $i \in 1, \dots, N$  do
6:        $\mathbf{w}_i \leftarrow \mathbf{w}_{i,s} + \alpha_w(\mathbf{w}_i - \mathbf{w}_{i,s})$  Apply lookahead (1st level)
7:        $\boldsymbol{\theta}_i \leftarrow \boldsymbol{\theta}_{i,s} + \alpha_\theta(\boldsymbol{\theta}_i - \boldsymbol{\theta}_{i,s})$ 
8:        $(\boldsymbol{\theta}_{i,s}, \mathbf{w}_{i,s}) \leftarrow (\boldsymbol{\theta}_i, \mathbf{w}_i)$  Update snapshots (1st level)
9:     end for
10:  end if
11:  if  $k_{ss}$  is not  $\emptyset$  and  $e \% k_{ss} == 0$  then
12:    for all agent  $i \in 1, \dots, N$  do
13:       $\mathbf{w}_i \leftarrow \mathbf{w}_{i,ss} + \alpha_w(\mathbf{w}_i - \mathbf{w}_{i,ss})$  Apply lookahead (2nd level)
14:       $\boldsymbol{\theta}_i \leftarrow \boldsymbol{\theta}_{i,ss} + \alpha_\theta(\boldsymbol{\theta}_i - \boldsymbol{\theta}_{i,ss})$ 
15:       $(\boldsymbol{\theta}_{i,s}, \boldsymbol{\theta}_{i,ss}, \mathbf{w}_{i,s}, \mathbf{w}_{i,ss}) \leftarrow (\boldsymbol{\theta}_i, \boldsymbol{\theta}_i, \mathbf{w}_i, \mathbf{w}_i)$  Update snapshots (1st & 2nd level)
16:    end for
17:  end if
18: end procedure

```

---

**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  $z_t$  as follows:

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

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].

## 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— $\mathbf{w}'_i, \theta'_i$  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:

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

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  $\theta_i$  and vice versa.

### 4.2 Proposed Methods

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

**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  $F_{\text{MADDPG}}$  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  $\theta, \mathbf{w}$ ), and their saved snapshots  $\theta_s, \mathbf{w}_s$ , as detailed in Algorithm 1. This method updates both the current networks and snapshots with

the  $\alpha$ -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**.

---

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

```

---

## 5 Experiments

### 5.1 Setup

We build<sup>1</sup> upon the open-source *PyTorch* implementation of MADDPG [Lowe et al., 2017]<sup>2</sup>. 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

<sup>1</sup>Code can be found at <https://github.com/badil96/VI-madpg.git>.

<sup>2</sup>Available at <https://github.com/Git-123-Hub/madpg-pettingzoo-pytorch>.



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

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

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

## 5.2 Results

**2-player games: rock–paper–scissors and matching pennies.** Figures 1a and 1b depict the average distance of the agents’ learned policies from the equilibrium policy. The baseline method eventually diverges. In contrast, LA-MADDPG consistently reduces the distance to the optimal policy, outperforming the baseline. While EG-MADDPG behaves similarly to the baseline, combining it with Lookahead stabilizes the performances. Additionally, Adam exhibits high variance across different seeds, while Lookahead significantly reduces variance, providing more stable and reliable results—an important factor in MARL experiments. We did not achieve full convergence to the Nash equilibrium with any of the algorithms, as we did not extensively tune the hyperparameters. For LA, we only used

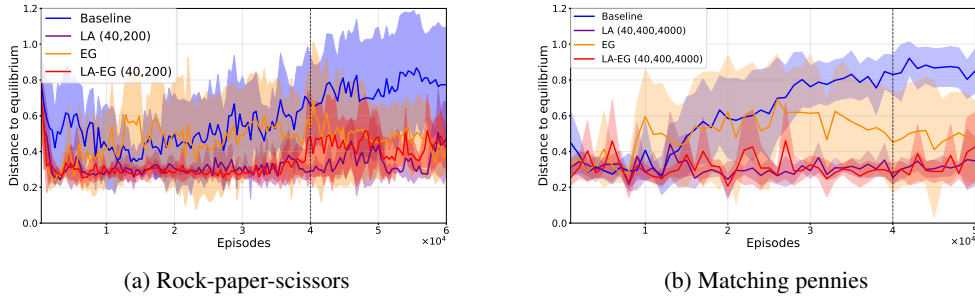


Figure 1: Comparison on the rock-paper-scissors game and matching pennies game between the *GD-MADDPG*, *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG* methods, denoted as *Baseline*, *LA*, *EG*, *LA-EG*, resp.  $x$ -axis: training episodes.  $y$ -axis: total distance of agents’ policies to the equilibrium 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.

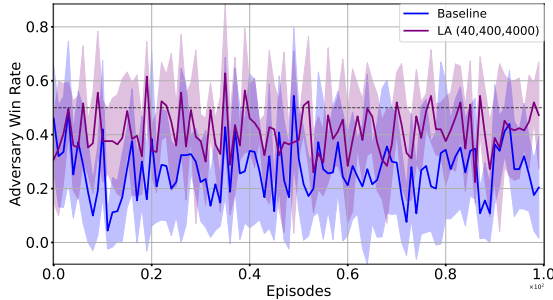


Figure 2: Comparison on the *MPE:Predator-prey* game between the *GD-MADDPG* and *LA-MADDPG*, optimization methods, denoted as *Baseline*, *LA*, resp.  $x$ -axis: evaluation episodes.  $y$ -axis: average win rate of adversary agents, averaged over 5 runs with different seeds. The dotted line depicts the start of the “shifting” (in first-in-first-out order) of the experiences in the replay buffer.

$\alpha = 0.5$  and randomly selected a few  $k$  values. Despite minimal tuning, the results provide strong indications.

**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.

**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 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

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.



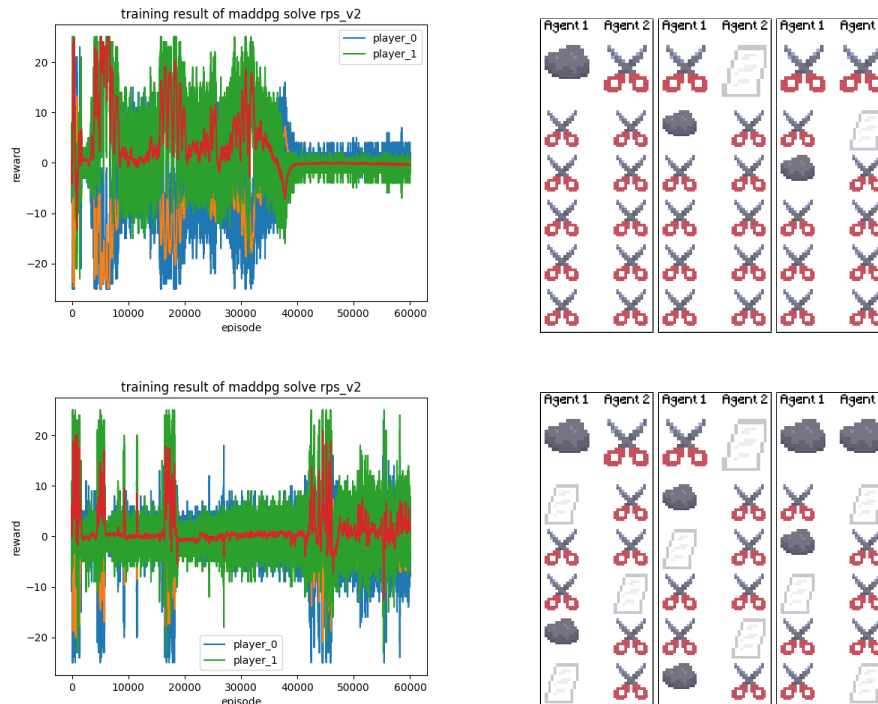


Figure 3: **Saturating rewards (left) versus actions of the learned policies at the end (right) in the rock–paper–scissors game. Top row: *GD-MADDPG*; bottom row: *LA-MADDPG*.** In the left column, green and blue represent actual rewards while red and orange show the running average through a window 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.

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 alternating their actions, which is the desired behavior. This underscores the need for stronger evaluation metrics in multi-agent reinforcement learning, particularly when the true equilibrium remains unknown. Refer to Appendix A.4 for additional discussion.

## 6 Conclusion

This paper tackles the inherent challenges in multi-agent reinforcement learning (MARL), where training is notoriously difficult due to high sensitivity to hyperparameters and random seed choices, making it challenging to reliably compare different methods. We explore whether Variational Inequality (VI) optimization techniques can improve 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.

While this work focuses on MADDPG, immediate future directions include investigating VI methods for other MARL algorithms and exploring additional existing VI techniques. These initial findings point toward promising opportunities for further development of VI-based methods in MARL, particularly in leveraging the structure of the optimization landscape.

## References

- W. Azizian, I. Mitliagkas, S. Lacoste-Julien, and G. Gidel. A tight and unified analysis of gradient-based methods for a whole spectrum of differentiable games. In *AISTATS*, pages 2863–2873, 2020.
- D. Balduzzi, S. Racaniere, J. Martens, J. Foerster, K. Tuyls, and T. Graepel. The mechanics of n-player differentiable games. In *ICML*, 2018.
- D. Bertsekas. *Rollout, Policy Iteration, and Distributed Reinforcement Learning*. Athena scientific optimization and computation series. Athena Scientific, 2021. ISBN 9781886529076.
- R. I. Bot, E. R. Csetnek, and P. T. Vuong. The forward-backward-forward method from continuous and discrete perspective for pseudo-monotone variational inequalities in Hilbert spaces. *arXiv:1808.08084*, 2020.
- R. I. Bot, E. R. Csetnek, and D.-K. Nguyen. Fast OGD in continuous and discrete time. *arXiv preprint arXiv:2203.10947*, 2022.
- M. Bowling. Convergence and no-regret in multiagent learning. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 17. MIT Press, 2004. URL [https://proceedings.neurips.cc/paper\\_files/paper/2004/file/88fee0421317424e4469f33a48f50cb0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2004/file/88fee0421317424e4469f33a48f50cb0-Paper.pdf).
- Y. Cai, A. Oikonomou, and W. Zheng. Tight last-iterate convergence of the extragradient method for constrained monotone variational inequalities. *arXiv:2204.09228*, 2022.
- T. Chavdarova, G. Gidel, F. Fleuret, and S. Lacoste-Julien. Reducing noise in GAN training with variance reduced extragradient. In *NeurIPS*, 2019.
- T. Chavdarova, M. Pagliardini, S. U. Stich, F. Fleuret, and M. Jaggi. Taming GANs with Lookahead-Minmax. In *ICLR*, 2021.
- T. Chavdarova, M. I. Jordan, and M. Zampetakis. Last-iterate convergence of saddle point optimizers via high-resolution differential equations. In *Minimax Theory and its Applications*, 2023.
- T. Chavdarova, T. Yang, M. Pagliardini, and M. I. Jordan. A primal-dual approach for solving variational inequalities with general-form constraints. In *ICLR*, 2024.
- F. Christianos, L. Schäfer, and S. V. Albrecht. Shared experience actor-critic for multi-agent reinforcement learning, 2021.
- R. W. Cottle and G. B. Dantzig. Complementary pivot theory of mathematical programming. *Linear Algebra and its Applications*, 1(1):103–125, 1968. ISSN 0024-3795.
- C. Daskalakis and I. Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. In *ITCS*, 2019.
- C. Daskalakis, A. Ilyas, V. Syrgkanis, and H. Zeng. Training GANs with optimism. In *ICLR*, 2018.
- J. Diakonikolas. Halpern iteration for near-optimal and parameter-free monotone inclusion and strong solutions to variational inequalities. In *COLT*, volume 125, 2020.
- T. Eimer, M. Lindauer, and R. Raileanu. Hyperparameters in reinforcement learning and how to tune them, 2023.
- F. Facchinei and J.-S. Pang. *Finite-dimensional Variational Inequalities and Complementarity Problems*. Springer, 2003a.
- F. Facchinei and J.-S. Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems Vol I*. Springer-Verlag, 2003b.
- W. Fan. A comprehensive analysis of game theory on multi-agent reinforcement. *Highlights in Science, Engineering and Technology*, 85:77–88, 03 2024. doi: 10.54097/gv6fpz53.

- G. Gidel, R. A. Hemmat, M. Pezeshki, R. L. Priol, G. Huang, S. Lacoste-Julien, and I. Mitliagkas. Negative momentum for improved game dynamics. In *AISTATS*, 2019.
- N. Golowich, S. Pattathil, and C. Daskalakis. Tight last-iterate convergence rates for no-regret learning in multi-player games. In *NeurIPS*, 2020a.
- N. Golowich, S. Pattathil, C. Daskalakis, and A. Ozdaglar. Last iterate is slower than averaged iterate in smooth convex-concave saddle point problems. In *COLT*, pages 1758–1784, 2020b.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- E. Gorbunov, N. Loizou, and G. Gidel. Extragradient method:  $\mathcal{O}(1/K)$  last-iterate convergence for monotone variational inequalities and connections with cocoercivity. In *AISTATS*, 2022.
- R. Gorsane, O. Mahjoub, R. de Kock, R. Dubb, S. Singh, and A. Pretorius. Towards a standardised performance evaluation protocol for cooperative marl, 2022.
- P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger. Deep reinforcement learning that matters, 2019.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- G. M. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 1976.
- A. Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Master’s thesis, 2009.
- J. G. Kuba, R. Chen, M. Wen, Y. Wen, F. Sun, J. Wang, and Y. Yang. Trust region policy optimisation in multi-agent reinforcement learning, 2022.
- M. Lanctot, J. Schultz, N. Burch, M. O. Smith, D. Hennes, T. Anthony, and J. Perolat. Population-based evaluation in repeated rock-paper-scissors as a benchmark for multiagent reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.03196>.
- P. Li, J. Hao, H. Tang, Y. Zheng, and X. Fu. Race: improve multi-agent reinforcement learning with representation asymmetry and collaborative evolution. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- T. Liang and J. Stokes. Interaction matters: A note on non-asymptotic local convergence of generative adversarial networks. *Artificial Intelligence and Statistics*, 2019.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning, 2019.
- R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Neural Information Processing Systems (NIPS)*, 2017.
- Y. Malitsky. Projected reflected gradient methods for monotone variational inequalities. *SIAM Journal on Optimization*, 25:502–520, 2015.
- L. Mescheder, A. Geiger, and S. Nowozin. Which Training Methods for GANs do actually Converge? In *ICML*, 2018.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015.
- D. Monderer and L. S. Shapley. Potential games. *Games and economic behavior*, 1996.
- S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian. Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In *ICML*, 2017.
- L. D. Popov. A modification of the arrow–hurwicz method for search of saddle points. *Mathematical Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.

- A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.
- T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. Foerster, and S. Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning, 2018.
- R. T. Rockafellar. Monotone operators associated with saddle-functions and minimax problems. *Nonlinear functional analysis*, 18(part 1):397–407, 1970.
- M. Rosca, Y. Wu, B. Dherin, and D. G. T. Barrett. Discretization drift in two-player games. In *ICML*, 2021.
- E. K. Ryu, K. Yuan, and W. Yin. Ode analysis of stochastic gradient methods with optimism and anchoring for minimax problems. *arXiv:1905.10899*, 2019.
- M. Samvelyan, T. Rashid, C. S. de Witt, G. Farquhar, N. Nardelli, T. G. J. Rudner, C.-M. Hung, P. H. S. Torr, J. Foerster, and S. Whiteson. The starcraft multi-agent challenge, 2019.
- K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning, 2019.
- R. Spica, D. Falanga, E. Cristofalo, E. Montijano, D. Scaramuzza, and M. Schwager. A real-time game theoretic planner for autonomous two-player drone racing. *arXiv:1801.02302*, 2018.
- G. Stampacchia. Formes bilinéaires coercitives sur les ensembles convexes. *Académie des Sciences de Paris*, 258:4413–4416, 1964.
- J. Terry, B. Black, N. Grammel, M. Jayakumar, A. Hari, R. Sullivan, L. S. Santos, C. Dieffendahl, C. Horsch, R. Perez-Vicente, et al. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:15032–15043, 2021.
- K. K. Thekumparampil, N. He, and S. Oh. Lifted primal-dual method for bilinearly coupled smooth minimax optimization. In *AISTATS*, 2022.
- P. Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60:237–252, 1995. ISSN 0377-0427.
- O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schrittwieser, J. Quan, S. Gaffney, S. Petersen, K. Simonyan, T. Schaul, H. van Hasselt, D. Silver, T. Lillicrap, K. Calderone, P. Keet, A. Brunasso, D. Lawrence, A. Ekeremo, J. Repp, and R. Tsing. Starcraft ii: A new challenge for reinforcement learning. *arXiv:1708.04782*, 2017.
- H. Wang, A. Sakhadeo, A. White, J. Bell, V. Liu, X. Zhao, P. Liu, T. Kozuno, A. Fyshe, and M. White. No more pesky hyperparameters: Offline hyperparameter tuning for rl, 2022.
- Z. Wang, B. Xu, and H.-J. Zhou. Social cycling and conditional responses in the rock-paper-scissors game. *Scientific Reports*, 4(1), July 2014. ISSN 2045-2322. doi: 10.1038/srep05830. URL <http://dx.doi.org/10.1038/srep05830>.
- Y. Yang and J. Wang. An overview of multi-agent reinforcement learning from game theoretical perspective, 2021.
- C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. The surprising effectiveness of ppo in cooperative, multi-agent games, 2022.
- M. Zhang, J. Lucas, J. Ba, and G. E. Hinton. Lookahead optimizer: k steps forward, 1 step back. In *NeurIPS*, 2019.
- L. Zheng, T. Fiez, Z. Alumbaugh, B. Chasnov, and L. J. Ratliff. Stackelberg actor-critic: Game-theoretic reinforcement learning algorithms, 2021.
- H.-J. Zhou. The rock–paper–scissors game. *Contemporary Physics*, 57(2):151–163, Mar. 2015. ISSN 1366-5812. doi: 10.1080/00107514.2015.1026556. URL <http://dx.doi.org/10.1080/00107514.2015.1026556>.
- M. Zhou, Y. Guan, M. Hayajneh, K. Niu, and C. Abdallah. Game theory and machine learning in uavs-assisted wireless communication networks: A survey. *arXiv:2108.03495*, 2021.

## A Appendix

### A.1 Additional Background

#### A.1.1 VI methods

In addition to those presented in the main part, we describe the following popular VI method, which we leave for future work.

**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  $\tau$ .



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

24:        Update actor policy using policy gradient formula and optimizer  $\mathcal{B}$ 
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  $a_i = \mu_i(o_i^j)$ 
26:        end for
27:        for all agent  $i \in 1, \dots, N$  do
28:           $\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$  (update target networks)
29:           $w_i' \leftarrow \tau w_i + (1 - \tau) w_i'$ 
30:        end for
31:      end if
32:       $\text{step} \leftarrow \text{step} + 1$ 
33:    until environment terminates
34: end for
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.

#### 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 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:      $(\theta_{i,s}, \theta_{i,ss}, w_{i,s}, w_{i,ss}) \leftarrow (\theta_i, \theta_i, w_i, w_i)$

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:          $y^j \leftarrow r_i^j + \gamma \mathbf{Q}^{\mu'}(\hat{\mathbf{x}}^j, a_1^j, \dots, a_N^j)$ , where  $\mathbf{a}_k^j = \{\mu_k^j(\sigma_k^j)\}$

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:**  $\theta, w$

---

---

**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:  $\varphi$ 

```

---

Table 2: Hyperparameters used for LA-MADDPG experiments.

Name	Description
Adam $lr$	0.01
Adam $\beta_1$	0.9
Adam $\beta_2$	0.999
Batch-size	1024
Update ratio $\tau$	0.01
Discount factor $\gamma$	0.95
Replay Buffer	$10^6$
learning step $T_{learn}$	100
$T_{rand}$	1024
Lookahead $\alpha$	0.5

## 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.

In all Rock-Paper-Scissors 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.

For the lookahead method, we experimented with different values of  $k$  and set  $\alpha = 0.5$ . We ran  $T = 60000$  training episodes, with a maximum of 25 environment steps ( $s$ ) per episode.

## 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.

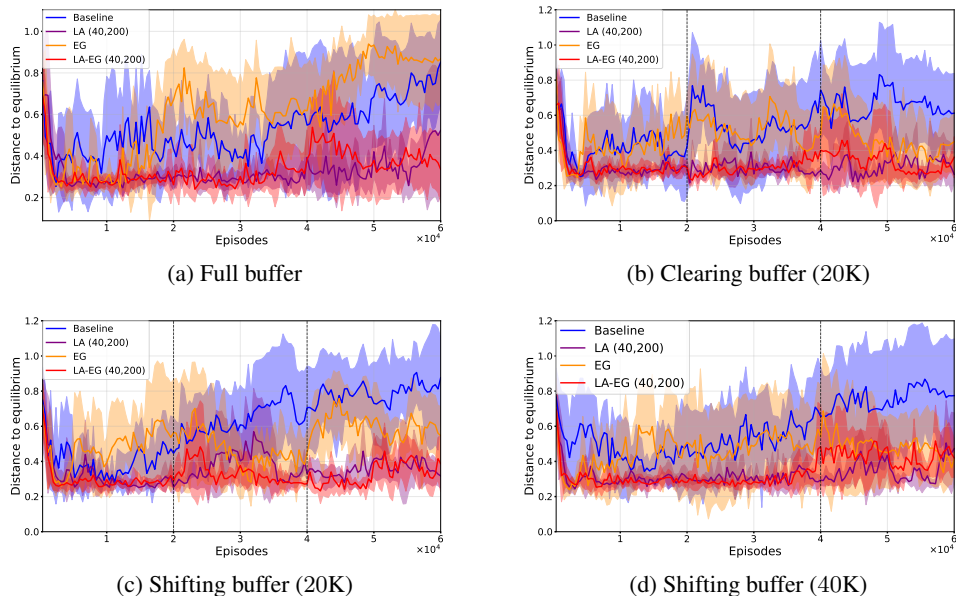


Figure 4: **Comparison of different buffer configurations (see Appendix A.3.1) and methods on Rock-paper-scissors game.**  $x$ -axis: training episodes.  $y$ -axis: 5-seed average norm between the two players’ policies and equilibrium policy  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$ . The dotted line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

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

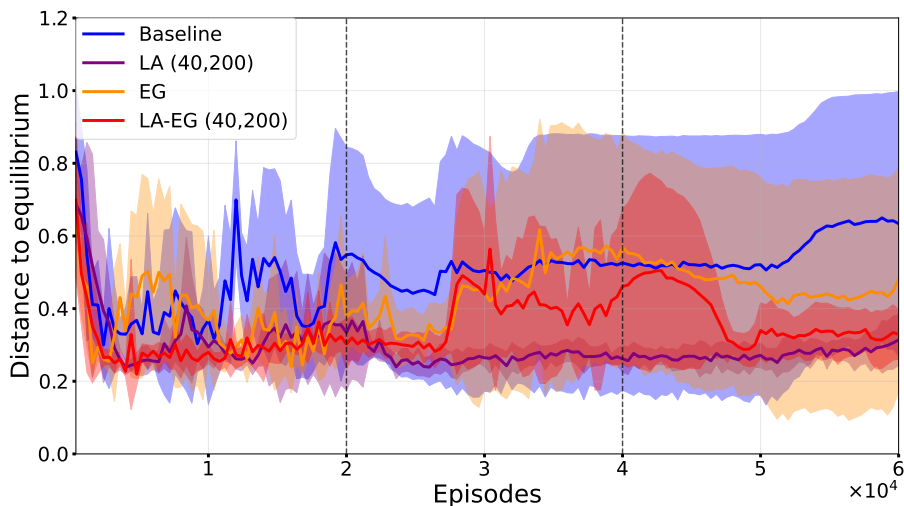


Figure 5: **Compares MADDPG with different LA-MADDPG configurations to the baseline MADDPG with (Adam) in rock-paper-scissors.**  $x$ -axis: training episodes.  $y$ -axis: 5-seed average norm between the two players’ policies and equilibrium policy  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$ . The dotted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

### 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.

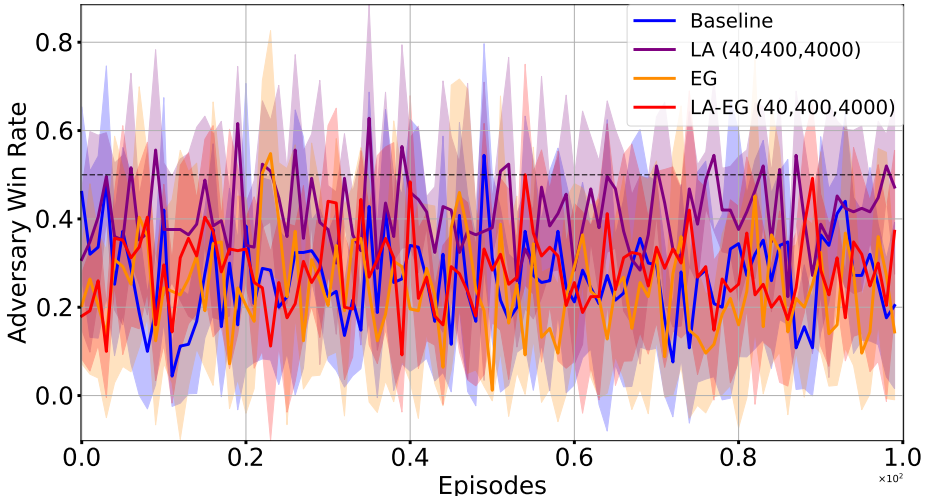


Figure 6: Comparison on the MPE–Predator-prey game between the *GD-MADDPG*, *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG* optimization methods, denoted as *Baseline*, *LA*, *EG*, *LA-EG*, resp. *x*-axis: evaluation episodes. *y*-axis: mean adversaries win rate, averaged over 5 runs with different seeds.

#### A.4 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 8, 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.



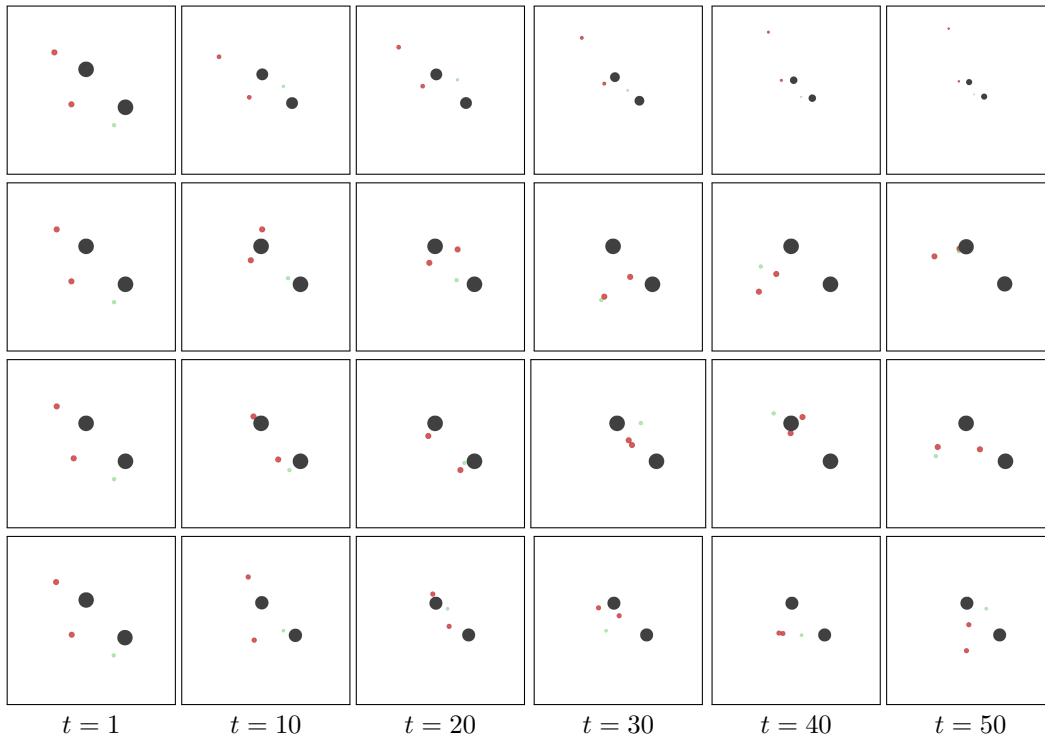


Figure 7: **Agents' trajectories of fully trained models with all considered optimization methods on the same environment seed of MPE: Predator-prey.** Snapshots show the progress of agents as time progresses in a 50 steps long environment. Each row contains snapshots of one method, from top to bottom: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

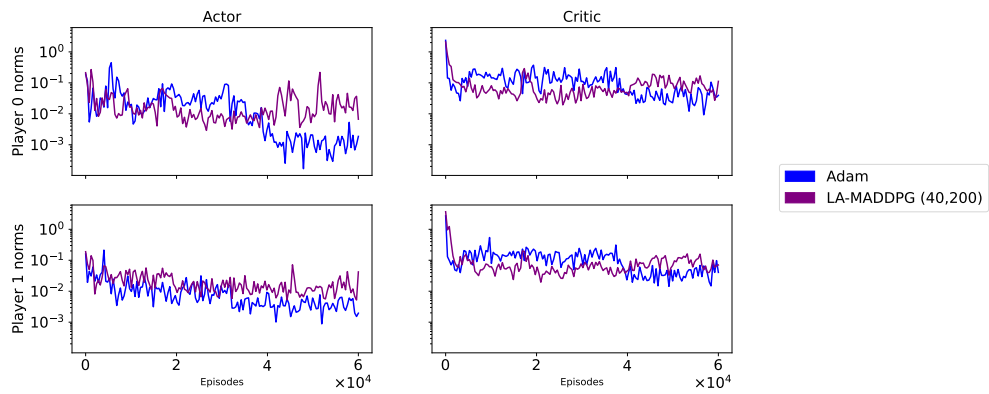


Figure 8: **Gradient norms across training in the rock-paper-scissors game.**