

LEARNING IN CIRCLES: ROTATIONAL DYNAMICS IN COMPETITIVE REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Optimization in competitive reinforcement learning (RL) differs from standard minimization. Actor-critic methods, in single- and multi-agent (MARL) settings, involve coupled objectives, so optimizing them jointly requires finding an equilibrium rather than performing independent descent. Through operator-theoretic viewpoint, we show that actor-critic models inherently exhibit rotational dynamics during learning, cycling around equilibria, thereby explaining in part the instability often observed in practice. Through the variational inequality (VI) framework for studying equilibrium seeking problems, we adopt the Lookahead method for VIs, which suppresses these rotations in actor-critic RL. Building on this, we introduce *Lookahead-(MA)RL (LA-(MA)RL)* to efficiently mitigate rotational dynamics. Across classical two-player games and multi-agent benchmarks, including *Rock-paper-scissors*, *Matching pennies*, and *Multi-Agent Particle environments*, LA-MARL consistently improves convergence and stability. Our results highlight optimization as a critical yet underexplored lever in RL: by rethinking the equilibrium-seeking dynamics, one can achieve substantial stability and performance gains.

1 INTRODUCTION

Competitive reinforcement learning (RL)—including single-agent actor-critic methods and multi-agent RL (MARL)—are widely used due to their ability to address complex challenges (see, for example, Omidshafiei et al., 2017; Vinyals et al., 2017; Spica et al., 2018; Zhou et al., 2021; Bertsekas, 2021). However, these methods are considerably harder to train than standard minimization problems. A key difficulty is the *interdependence of parameters*: for example, actors and critics update each other’s loss landscapes, so the optimization problem itself changes during training. Numerous algorithmic fixes have been proposed, such as target networks, regularizers, and alternative learning objectives, yet instabilities and reproducibility issues remain, especially in competitive multi-agent settings. Recent reports of large seed-to-seed performance variability in widely used MARL benchmarks highlight the scope of this challenge (Bettini et al., 2024b; Gorsane et al., 2022).

One mechanism rising from this interdependence in competitive RL algorithms is *rotational learning dynamics*, where iterates cycle around equilibria instead of converging (Mescheder et al., 2018; Balduzzi et al., 2018). For instance, the Gradient Descent (GD) method for the $\min_{z_1 \in \mathbb{R}^{d_1}} \max_{z_2 \in \mathbb{R}^{d_2}} z_1 \cdot z_2$ game, rotates around the equilibrium $(0, 0)$ for infinitesimally small learning rates, and diverges away for practical choices of its value. Consequently, GD—and its adaptive variants like *Adam* (Kingma & Ba, 2015)—fail to converge for a broad class of equilibrium-seeking problems.

Figure 1 illustrates this phenomenon in a simple two-player RL setting: under standard gradient descent the joint strategies rotate around the mixed-strategy equilibrium without convergence, while the *Lookahead* (Zhang et al., 2019b; Chavdarova et al., 2021) game optimization method (described in Section 2) contracts the trajectories and reaches equilibrium. Such learning dynamics may not be the only cause of instability in competitive RL, but offer a tractable lens through which we can study their cause.

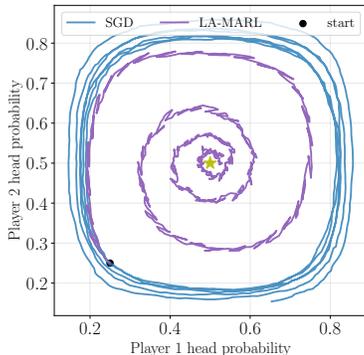


Figure 1: **Rotational dynamics in learning: Matching Pennies.** Joint probabilities of playing *Heads* for both players. x -axis: Player 1’s probability; y -axis: Player 2’s. SGD iterates (blue) circle around the mixed-strategy equilibrium (yellow star) without converging, while Lookahead trajectories contract and converge.

Hence, in this work, we ask: *to what extent are these instabilities due to the choice of optimization method, rather than the design of the algorithm itself?* Many RL algorithms implicitly aim to solve an equilibrium problem. Our goal is therefore to keep the high-level algorithm and objectives fixed, but to disentangle the effect of the optimization (learning) dynamics on training outcomes. To do so, we adopt insights from the *variational inequality* (VI) framework, which is a natural mathematical model for equilibrium problems. This motivates two central questions:

(i) *Do competitive RL methods such as actor–critic indeed exhibit rotational or game-like dynamics, and why?*

(ii) *When keeping all other choices fixed, how much can VI-based optimization improve stability and learning outcomes in competitive RL?*

Our objective is to provide systematic evidence that the choice of optimization method itself significantly impacts competitive RL outcomes. MARL serves as an especially relevant setting, since agents not only interact through coupled objectives but also through data collection, making optimization dynamics even more critical. Particularly, MARL problems are modeled with *stochastic games* (Littman, 1994); refer to Section 2. Three main MARL learning paradigms are commonly used:

- *value-based learning*—focuses on estimating so-called *value functions* (e.g., *Q*-learning, Deep *Q*-Networks (Mnih et al., 2015)) to learn action-values first and infer a policy based on it,
- *policy-based learning*—directly optimizes the *policy* (e.g., *REINFORCE* (Williams, 1992)) by adjusting action probabilities without explicitly learning the value functions, and
- *actor-critic methods*—combines value-based and policy-based approaches where an actor selects actions, and a critic evaluates them.

Furthermore, MARL can be broadly categorized into *centralized* and *independent* learning approaches. In centralized MARL, a global critic or shared value function leverages information from all agents to guide learning, improving coordination (Sunehag et al., 2017; Lowe et al., 2017; Yu et al., 2021). In contrast, independent MARL treats each agent as a separate learner, promoting scalability but introducing non-stationarity as agents continuously adapt to each other’s evolving policies (Matignon et al., 2012; Foerster et al., 2017). In this work, we focus on *Centralized Training Decentralized Execution* (CTDE) approaches, specifically the ones with centralized critics. Several of the new algorithms in MARL belong to this category such as (MADDPG, Lowe et al., 2017), (MATD3, Ackermann et al., 2019), (MAPPO, Yu et al., 2021), and (COMA, Foerster et al., 2018). In this work, we primarily focus on the CTDE actor-critic MARL learning paradigm, and adopt VI approaches leveraging a (combination of) *nested-Lookahead-VI* (Chavdarova et al., 2021) and *Extragradient* (Korpelevich, 1976) methods for iteratively solving variational inequalities (VIs).

Contributions. This paper:

- Analyzes standard actor–critic methods through their Jacobian and shows that, even in a single-agent setting, the joint actor–critic updates generically induce a non-symmetric operator with a rotational (game-like) component, which amplifies in centralized multi-agent extensions. To make this concrete, we introduce a minimal theoretical setting that isolates this phenomenon and validate it with numerical experiments (Section 3).
- Adopts the variational inequality (VI) viewpoint to cast common policy-gradient and centralized actor–critic MARL methods (e.g., MADDPG, MATD3) as explicit VI problems in the joint parameter space, providing concrete operators that match their practical implementations.
- We propose *LA-(MA)RL* (Algorithm 1), a scalable approach for neural network-based agents. While presented for actor-critic systems, the method generalizes to other MARL settings. *LA-MARL* is computationally efficient, making it well-suited for large-scale optimization tasks.
- We evaluate these methods on classic matrix games (rock–paper–scissors, matching pennies) and in benchmark multi-agent particle environments (Lowe et al., 2017), showing improved stability over standard optimization.

As a side contribution, we illustrate the limitations of reward-based performance metrics in competitive MARL and advocate distance-to-equilibrium and cross-play style evaluations as more informative alternatives in these settings.

Code implementation: <https://anonymous.4open.science/r/VI-marl-47A4>.

Brief related works discussion. In mathematics and numerical optimization, equilibrium-finding problems can be modeled using several frameworks, most notably the *Variational Inequality* (VIs, Stampacchia, 1964; Facchinei & Pang, 2003b) framework (see Section 2 for a formal definition). Recent advances in solving variational inequalities (VIs) have been heavily influenced by challenges observed in training generative adversarial networks (GANs, Goodfellow et al., 2014). This progress spans both theoretical developments—with new convergence guarantees (e.g., Golowich et al., 2020b; Daskalakis et al., 2020b; Gorbunov et al., 2022)—and practical algorithms for large-scale optimization (Diakonikolas, 2020; Chavdarova et al., 2021).

A large body of work has studied optimization in actor–critic (AC) and centralized MARL. In the single-agent case, numerous AC algorithms have been proposed (Lillicrap et al., 2015; Schulman et al., 2015; 2017) and their convergence properties analyzed (Mnih et al., 2016; Konda & Tsitsiklis, 1999; Holzleitner et al., 2020; Kumar et al., 2025; Chen & Zhao, 2025). Recent works emphasize their adversarial, game-like structure, connecting AC to GANs (Pfau & Vinyals, 2017) and showing that rotations can be interpreted as Stackelberg dynamics (Zheng et al., 2021). In MARL, research has focused on two-player zero-sum settings with regret guarantees (Bai & Jin, 2020; Xie et al., 2020), as well as extensions to linear quadratic games (Kalman, 1960; Fazel et al., 2018; Bu et al., 2019; Zhang et al., 2021; Mazumdar et al., 2020; Hambly et al., 2023), where convergence remains significantly harder in general-sum cases. Broader formulations such as polymatrix games further expose the limits of gradient descent and motivate alternative algorithms (Ma et al., 2021; Janovskaja, 1968). Despite progress, optimization and convergence in multi-agent actor–critic methods (Bettini et al., 2024a) remain open challenges, with non-stationarity and complex interactions continuing to drive instability.

We provide a detailed related works discussion in Appendix A.

2 PRELIMINARIES

Notation. We use bold lowercase for vectors, curly capitals for sets, lowercase for real-valued functions, and capitals for operators $\mathcal{Z} \mapsto \mathcal{Z}$ (e.g., F). We write $[n] = \{1, \dots, n\}$ and \succeq for positive semidefinite matrices. Let \mathcal{Z} be a convex, compact subset of Euclidean space equipped with inner product $\langle \cdot, \cdot \rangle$. We follow standard MARL notation in what follows.

MARL. *Markov games* (MGs, or *stochastic games*, Shapley, 1953; Littman, 1994) generalize *Markov Decision Processes* (MDPs Puterman, 1994) to a multi-agent setting. An MG is given by:

$$(n, \mathcal{S}, \{\mathcal{A}_i\}_{i=1}^n, p, \{r_i\}_{i=1}^n, \gamma), \quad (\text{MG})$$

with n agents sharing a state space \mathcal{S} . Each agent $i \in [n]$ receives observation $\mathbf{o}_i \in \mathcal{O}$ of the current state $\mathbf{s} \in \mathcal{S}$ of the environment. In the most general case, agent i 's observation $\mathbf{o}_i = f(\mathbf{s})$, where $f: \mathcal{S} \rightarrow \mathcal{O}_i$. For instance, f can be an identity or coordinate-selection map with $\mathcal{O}_i \subseteq \mathcal{S}$, or f can be a nontrivial mapping. Based on its policy $\pi_i: \mathcal{O}_i \rightarrow \mathcal{A}_i$, each agent $i \in [n]$ selects an action $a_i \in \mathcal{A}_i$, where \mathcal{A}_i is its finite action set. The joint actions of all agents are represented as $\mathbf{a} \triangleq (a_1, \dots, a_n)$, and the joint action space as $\mathcal{A} \triangleq \mathcal{A}_1 \times \dots \times \mathcal{A}_n$.

The environment transitions to a new state $\mathbf{s}' \in \mathcal{S}$ according to a *transition function* $p: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, where $\Delta(\mathcal{S})$ is the space of probability distributions over \mathcal{S} (non-negative $|\mathcal{S}|$ -dimensional vector summing to 1). The function p specifies the probability distribution of the next state \mathbf{s}' , given the current state \mathbf{s} and the joint action \mathbf{a} . Each agent $i \in [n]$ receives a reward r_i , where the reward function $r_i: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ depends on the current state, the joint action, and the resulting next state. The importance of future rewards is governed by the *discount factor* $\gamma \in [0, 1)$.

MGs generalize MDPs and *repeated games* (Aumann, 1995) by introducing non-stationary dynamics, where agents learn their policies jointly and adaptively. Each agent $i \in [n]$ aims to maximize its expected cumulative reward (return), defined as:

$$v_i^{\pi_i, \pi_{-i}}(\mathbf{s}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_i(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \sim \boldsymbol{\rho}, \mathbf{a}_t \sim \boldsymbol{\pi}(s_t) \right], \quad (\text{MA-Return})$$

with $\boldsymbol{\pi} \triangleq (\pi_1, \dots, \pi_n)$ the joint policy of all agents, π_{-i} the policies of all agents except agent i , and $\boldsymbol{\rho}$ the initial state distribution. Agents' interaction introduces non-stationarity (due to evolving policies) and reward interdependencies, leading to a distinct optimization landscape. *Nash equilibria* serve as solution concepts, where no agent can improve its return by unilaterally altering its policy.

MADDPG. *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, leveraging a centralized training with decentralized execution paradigm. Each agent $i \in [n]$ has: (i) *critic* network $Q_i: \mathcal{O}_1 \times \dots \times \mathcal{O}_n \times \mathcal{A} \rightarrow \mathbb{R}$, parametrized by $w_i \in \mathbb{R}^{d_i^Q}$: acts as a centralized action-value function, evaluating the expected return of joint actions \mathbf{a} in state \mathbf{s} and (ii) *actor* network $\mu_i: \mathcal{O}_i \rightarrow \Delta(\mathcal{A}_i)$, parametrized by $\theta_i \in \mathbb{R}^{d_i^\mu}$: represents the agent’s policy, mapping agents’ observation of states \mathbf{s} to a probability distribution over actions a_i .

MADDPG uses *target* networks for stability, which are delayed versions of the critic and actor networks: *target critic* \bar{Q}_i , is parametrized by $\bar{w}_i \in \mathbb{R}^{d_i^Q}$, and *target-actor* $\bar{\mu}_i$, parametrized by $\bar{\theta}_i$. The target networks are updated using a soft update, with $\tau \in (0, 1]$ as:

$$\bar{w}_i \leftarrow \tau w_i + (1 - \tau) \bar{w}_i, \quad (\text{Target-Critic}) \quad \bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i. \quad (\text{Target-Actor})$$

MATD3. *Multi-agent TD3* ([Ackermann et al., 2019](#)) extends MADDPG with: (i) *dual critics*: each agent uses two critics $Q_{i,1}, Q_{i,2}$ and the smaller target value to reduce overestimation bias; (ii) *delayed actor updates*: policies are updated every c critic steps for stability. Gaussian noise is also added to target actions for exploration. These modifications improve robustness over MADDPG in multi-agent settings.

Variational Inequality ([Stampacchia, 1964](#); [Facchinei & Pang, 2003b](#)). Variational Inequalities (VIs) extend beyond standard minimization problems to encompass a broad range of equilibrium-seeking problems. The connection to such general problems can be understood from the optimality condition for convex functions: a point z^* is an optimal solution if and only if $\langle z - z^*, \nabla f(z^*) \rangle \geq 0, \forall z \in \text{dom} f$. In the framework of VIs, this condition is generalized by replacing the gradient field ∇f with a more general vector field F , allowing for the modeling of a wider class of problems. Formally, the VI goal is to find an equilibrium z^* from the domain of continuous strategies \mathcal{Z} , such that:

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

where $F: \mathcal{Z} \rightarrow \mathbb{R}^d$, referred to as the *operator*, is continuous, and \mathcal{Z} is a subset of the Euclidean d -dimensional space \mathbb{R}^d . VIs are thus characterized by the tuple (F, \mathcal{Z}) , denoted herein as $\text{VI}(F, \mathcal{Z})$. For a more comprehensive introduction to VIs, including examples and applications, see [Appendix B.1](#).

VI methods. The *gradient descent* method straightforwardly extends for the VI problem as follows:

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

where t denotes the iteration count, and $\eta \in (0, 1)$ the step size or learning rate.

Extragradient ([Korpelevich, 1976](#)) is a modification of GD, which uses a “prediction” step to obtain an extrapolated point $z_{t+\frac{1}{2}}$ using GD: $z_{t+\frac{1}{2}} = z_t - \eta F(z_t)$, and the gradients at the *extrapolated* point are then applied to the *current* iterate z_t as follows:

$$z_{t+1} = z_t - \eta F\left(z_{t+\frac{1}{2}}\right). \quad (\text{EG})$$

Unlike gradient descent, EG converges in some simple game instances, such as in games linear in both players ([Korpelevich, 1976](#)).

The *nested-Lookahead-VI* (LA) algorithm for VI problems (Alg. 3, [Chavdarova et al., 2021](#)), is a general wrapper of a “base” optimizer $B: \mathbb{R}^n \rightarrow \mathbb{R}^n$ where, after every k iterations with B , $z_{t+1} = B(z_t)$ an averaging step is performed as follows:

$$z_{t+k} \leftarrow z_t + \alpha(z_{t+k} - z_t), \quad \alpha \in [0, 1]. \quad (\text{LA})$$

For this purpose a copy (snapshot) of the iterate after the averaging step is stored for the next LA update. See [Appendix B.1.1](#) for an alternative view.

This averaging can be applied recursively across multiple levels l , when using LA as base optimizer, typically with $l \in [1, 3]$. In [Algorithm 6](#), the parameter $k^{(j)}$ at level $j \in [l]$ is defined as the multiple of $k^{(j-1)}$ from the previous level $j - 1$, specifically $k^{(j)} = c_j \cdot k^{(j-1)}$. For $l = 1, k = 2$, LA has connections to EG ([Chavdarova et al., 2023](#)), however for higher values of k and l the resulting operator exhibits stronger contraction ([Chavdarova et al., 2021](#); [Ha & Kim, 2022](#)), which effectively addresses rotational learning dynamics.

3 ACTOR–CRITIC REINFORCEMENT LEARNING AS A GAME: BEYOND STANDARD MINIMIZATION

We begin with a *single-agent discounted MDP* and show that, even in this simplest case, the coupled actor–critic updates form a game in parameter space. This toy model isolates the interaction between actor and critic and will serve as the basis for our multi-agent extension in Section ??.

Single-agent setting. Consider the following setting with one unrolled state s' per trajectory as follows.

Setting 1 (off-policy deterministic AC setting). Action $a \in \mathbb{R}$, state $s \in \mathbb{R}^d$. Features $\phi(s, a) = f_\phi(s) + a m(s) \in \mathbb{R}^c$ with $m: \mathbb{R}^d \rightarrow \mathbb{R}^c$. Critic $Q_w(s, a) = \langle w, \phi(s, a) \rangle$, actor $\pi_\theta(s) = \langle \theta, s \rangle$, batch $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^B$, discount $\gamma \in [0, 1)$. Furthermore, let $\phi_i \triangleq \phi(s_i, a_i)$, $a'_i \triangleq \pi_\theta(s'_i)$, $\phi'_i \triangleq \phi(s'_i, a'_i)$, $\delta_i \triangleq \langle w, \phi_i \rangle - r_i - \gamma \langle w, \phi'_i \rangle$.

Lemma (informal) 1 (complex eigenvalues generically). Setting 1 with $\gamma > 0$, has operator $F = (F_w; F_\theta)^\top$ with $F_w = \frac{2}{B} \sum_i \delta_i (\phi_i - \gamma \phi'_i)$, $F_\theta = -\frac{1}{B} \sum_i s_i \langle w, m(s_i) \rangle$. Its Jacobian has block form $J \equiv (J_{ww}, J_{w\theta}; J_{\theta w}, \mathbf{0})^\top$ with:

$$J_{ww} = \frac{2}{B} \sum_i (\phi_i - \gamma \phi'_i) (\phi_i - \gamma \phi'_i)^\top \succeq 0,$$

$$J_{w\theta} = -\frac{2\gamma}{B} \sum_i [(\phi_i - \gamma \phi'_i) s_i^\top w^\top m(s'_i) + m(s'_i) s_i^\top \delta_i], \quad \text{and } J_{\theta, w} = -\frac{1}{B} \sum_i s_i m(s_i)^\top.$$

Although the underlying environment is a single-agent discounted MDP, the learning problem is an equilibrium-seeking problem between two coupled objectives: the critic seeks to minimize $F_w(w; \theta)$, while the actor seeks to maximize $F_\theta(w; \theta)$. We therefore interpret the pair (w, θ) as two players in a game in parameter space with joint dynamics governed by F . Furthermore, the Jacobian J is typically non-symmetric and, for generic data/parameters, has a non-real conjugate eigenpair (i.e., the linearized dynamics include a rotational component). All eigenvalues are real only under special symmetry/degeneracy (e.g., $\gamma = 0$ or $J_{w\theta} = J_{\theta w}^\top$, which in this model requires $\delta_i \equiv 0$ and $w^\top m(s'_i) \equiv 0$ for all i).

The above lemma shows setting 1 has *game structure*: except in special symmetry/degeneracy cases (e.g., $\gamma = 0$ or $J_{w\theta} = 0$), the Jacobian is non-symmetric and generically has complex conjugate eigenpair. This is a hallmark of game (rotation component), hence departs from standard minimization where in contrast, J —called Hessian in that context—is always symmetric. Full statement and proof deferred to Appendix C.1.

Since the operator of the problem is rotational, running the gradient descent method (GD) produces cyclic trajectories around an equilibrium. Without appropriate learning dynamics, the joint model w, θ may diverge away. Figure 2 (top) depicts (GD) on a concrete numerical instance of the problem class of setting 1. (GD) does not converge, but rather the joint last iterate (w, θ) cycles around the equilibrium point (w^*, θ^*) . Lookahead LA on the other hand, shown in Figure 2 (bottom), increases the contraction (see Appendix D.1) of the base (GD) and converges to the equilibrium. The example game details are given in appendix section C.1.1.

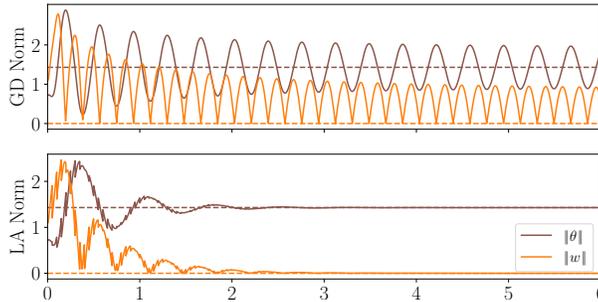


Figure 2: **Rotational learning dynamics**—on Setting 1 instance; see Appendix C.1.1. x -axis: training iterations, y -axis: $\|w\|$ and $\|\theta\|$ (times 10^3); the separate norms at equilibrium— $\|w^*\|, \|\theta^*\|$ —are depicted with dashed lines in orange and brown, resp. *Top*: gradient descent; *bottom*: Lookahead.

Extension to the multi-agent setting. With n actors and n centralized critics (each seeing the full state and joint action), the Jacobian keeps the single-agent saddle/block character but now aggregates n actor–critic couplings; the effective coupling (the product of the actor→critic and critic→actor blocks) sums over agents and is generically sign-indefinite, so a non-real conjugate eigenpair persists

and the spectral radius typically *increases* with n (at least matching the strongest single-agent term), while only degenerate triangular/symmetric cases (e.g., $\gamma = 0$ or vanishing cross-couplings) yield an all-real spectrum. Appendix C includes the corresponding lemma and proof for this setting.

While gradient descent might converge in problems with only a mild rotational component, our primary interest is the more competitive settings of (multi-agent) RL. This motivates the following Section 4, where we reinterpret existing MARL methods through the lens of variational inequalities (VIs) and introduce a general algorithmic family that incorporates the insights above.

4 COMPETITIVE RL THROUGH VARIATIONAL INEQUALITIES: EXISTING METHODS AND A GENERAL ALGORITHM

In this section, we reinterpret competitive RL methods through the lens of variational inequalities (VIs). We first unify existing approaches—including multi-agent policy gradient and actor–critic methods—under a common operator framework. Since single-agent actor–critic is a special case of the multi-agent setting, our focus is on the latter in the presentation, which captures a broader set of methods. Building on the insights from Section 3, we then introduce a general algorithmic family that augments these operators with Lookahead (LA) and/or Extragradient (EG) updates, applied directly in the joint parameter space of all agents. The pseudocode formulation naturally subsumes single-agent actor–critic and multi-agent policy gradient as special cases.

4.1 MARL OPERATORS

General MARL. Policy-based learning solves (**MA-Return**), where agents optimize their policies to maximize return. The associated operator F_{MAR} (*multi-agent-return*), with $\mathcal{Z} \equiv \mathcal{A}$, is:

$$F_{\text{MAR}} \left(\begin{bmatrix} \vdots \\ \boldsymbol{\pi}_i \\ \vdots \end{bmatrix} \right) \triangleq \begin{bmatrix} \vdots \\ \nabla_{\boldsymbol{\pi}_i} v_i^{\boldsymbol{\pi}_i, \boldsymbol{\pi}_{-i}} \\ \vdots \end{bmatrix} \equiv \begin{bmatrix} \vdots \\ \nabla_{\boldsymbol{\pi}_i} \left(\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \sim \boldsymbol{\rho}, \mathbf{a}_t \sim \boldsymbol{\pi}(\mathbf{s}_t) \right] \right) \\ \vdots \end{bmatrix}. \quad (F_{\text{MAR}})$$

Actor-critic (MA)RL. We denote by \mathbf{x} the full state representation, from which each agent’s observation \mathbf{o}_i is derived. Each agent $i \in [n]$ has a centralized critic $\mathbf{Q}_i^\mu(\mathbf{x}_t, \mathbf{a}_t; \mathbf{w}_i)$ and a policy $\mu_i(\mathbf{o}_i; \boldsymbol{\theta}_i)$. Given a batch $\mathcal{B} = \{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}'^j)\}_{j=1}^{|\mathcal{B}|}$ from the replay buffer \mathcal{D} , equilibrium is sought by solving

(VI) with operator:

$$F_{\text{MAAC}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \ell_i^w(\cdot; \mathbf{w}_i, \boldsymbol{\theta}_i) \right) \\ \nabla_{\boldsymbol{\theta}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \ell_i^\theta(\cdot; \mathbf{w}_i, \boldsymbol{\theta}_i) \right) \\ \vdots \end{bmatrix}, \quad (F_{\text{MAAC}})$$

with parameter space $\mathcal{Z} \equiv \mathbb{R}^d$, $d = \sum_{i=1}^n (d_i^Q + d_i^\mu)$. MAAC stands for *multi-agent actor–critic*. Even for $n=1$, the actor–critic interaction forms a game: \mathbf{w}_i updates depend on $\boldsymbol{\theta}_i$ and vice versa, a coupling that drives the dynamics in MARL, see Section 3.

MADDPG (Lowe et al., 2017) As an example, we expand (F_{MAAC}) for MADDPG. The critic and actor losses are:

$$\ell_i^w(\cdot; \mathbf{w}_i, \boldsymbol{\theta}_i) = (y_i - \mathbf{Q}_i^\mu(\mathbf{x}^j, \mathbf{a}^j; \mathbf{w}_i))^2, y_i = r_i^j + \gamma \mathbf{Q}_i^\mu(\mathbf{x}'^j, \mathbf{a}'; \mathbf{w}_i) \Big|_{\mathbf{a}' = \bar{\mu}(\sigma'^j)}. \quad (\ell_{\text{MADDPG}}^w)$$

$$\ell_i^\theta(\cdot; \mathbf{w}_i, \boldsymbol{\theta}_i) = \mu_i(\sigma_i^j; \boldsymbol{\theta}_i) \nabla_{\mathbf{a}_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, \mathbf{a}_1^j, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n^j; \mathbf{w}_i) \Big|_{\mathbf{a}_i = \mu_i(\sigma_i^j)}. \quad (\ell_{\text{MADDPG}}^\theta)$$

Here, the critic minimizes the Bellman error while the actor maximizes the Q -value, highlighting their coupled dynamics. Operators for other algorithms are deferred to Appendix D.

4.2 LA-(MA)RL PSEUDOCODE

To solve the VI problem with MARL operators such as (F_{MAR}) or (F_{MAAC}), we propose two methods: *LA-MARL* and *EG-MARL*.

Algorithm 1 LA-(MA)RL Pseudocode

```

324 1: Input: Env  $\mathcal{E}$ ; agents  $n$ ; episodes  $T$ ; actor nets  $\{\mu_i\}_{i=1}^n$  (params  $\theta$ ); critic nets  $\{Q_i\}_{i=1}^n$  (params
325  $w$ ); targets  $\{\bar{\mu}_i, \bar{Q}_i\}$ ; VI operator  $F$ ; base optimizer  $B$ ; discount  $\gamma$ ; soft update  $\tau$ ; random steps
326  $t_{\text{rand}}$ ; learn interval  $t_{\text{learn}}$ ; LA hyperparams  $\mathcal{L} = (l, \{k^{(j)}\}, \alpha_\theta, \alpha_w)$ .
327
328 2: Init: Replay buffer  $\mathcal{D} \leftarrow \emptyset$ ; set targets  $\bar{\theta} \leftarrow \theta, \bar{w} \leftarrow w$ ; LA state  $\phi$ .
329
330 3: for  $e = 1$  to  $T$  do
331 4:   reset state  $x$  from  $\mathcal{E}$ .
332 5:   while episode not done do
333 6:     Act: for each agent  $i$ , choose  $a_i \sim \text{Uniform}(\mathcal{A}_i)$  if  $\text{step} \leq t_{\text{rand}}$ , else  $a_i = \mu_i(o_i)$ .
334 7:     Step: execute  $a$ , observe  $(r, x')$ ; store  $(x, a, r, x')$  in  $\mathcal{D}$ ; set  $x \leftarrow x'$ .
335 8:     if  $\text{step} \bmod t_{\text{learn}} = 0$  and  $\mathcal{D}$  sufficiently full then
336 9:       Learn: sample batch  $\mathcal{B} \subset \mathcal{D}$ .
337 10:       $(\theta, w) \leftarrow \text{VI-UPDATE}(F, \mathcal{B}, \theta, w; B, \gamma)$ 
338 11:      Targets:  $\bar{\theta} \leftarrow \tau\theta + (1 - \tau)\bar{\theta}, \bar{w} \leftarrow \tau w + (1 - \tau)\bar{w}$ .
339 12:    end if
340 13:  end while
341 14:  Lookahead:  $(\theta, w) \leftarrow \text{NESTEDLOOKAHEAD}((\theta, w), \mathcal{L})$ .
342 15: end for
343 16: Output:  $\theta^{(l-1)}, w^{(l-1)}$ .

```

LA-MARL, Algorithm 1. LA-MARL periodically (at $\text{step} = k$) saves snapshots of all agents’ actor and critic networks and averages them with the current parameters, using α -averaging (Algorithm 6). It can be nested at multiple levels, where higher levels update less frequently. All agents run the LA step synchronously to ensure contraction. Variants of the algorithm specialized for MADDPG and MATD3 are provided in the appendix (Algorithms 7, 8).

Adaptability. Although designed for off-policy actor–critic, LA-(MA)RL extends naturally to other MARL settings: (i) policy gradient methods, by removing the critic; and (ii) on-policy methods, by omitting target networks. Importantly, lookahead must be applied in the joint parameter space, since adversarial objectives induce rotational dynamics, and per-agent averaging would be inconsistent.

(LA-)EG-MARL. EG-MARL applies the extragradient update rule (EG) to all agents’ actor and critic networks (Algorithm 2). LA-EG-MARL further combines extragradient with nested lookahead by using (EG) as the base optimizer.

Convergence. For monotone operators (see Appendix B.1.1 for the definition), gradient descent is known to diverge (Korpelevich, 1976). In contrast, EG-MARL and its LA variants converge under this assumption (Korpelevich, 1976; Chavdarova et al., 2021; Gorbunov et al., 2022; Pethick et al., 2023). Lookahead improves the contractiveness of the fixed-point operator, and nesting further strengthens stability, mitigating divergence in competitive (rotational) settings. See Appendix D for details.

5 EXPERIMENTS

5.1 SETUP

We build on the open-source *PyTorch* MADDPG implementation and extend it to MATD3, using the original papers’ hyperparameters (Appendix E). Experiments span two zero-sum games—rock–paper–scissors and matching pennies—and a competitive MPE scenario, using *Petting-Zoo* (Terry et al., 2021) implementations.

Rock–paper–scissors (RPS) and Matching Pennies (MP). These canonical games have closed-form mixed Nash equilibria (MNE) and expose inherent cycling (Zhou, 2015; Wang et al., 2014; Srinivasan et al., 2018). In RPS, two players ($n = 2$) choose among three actions ($m = 3$) with MNE $\pi_{\text{RPS}}^* (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$; players observe the opponent’s previous action and receive $+1/0/-1$ for win/tie/loss

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

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

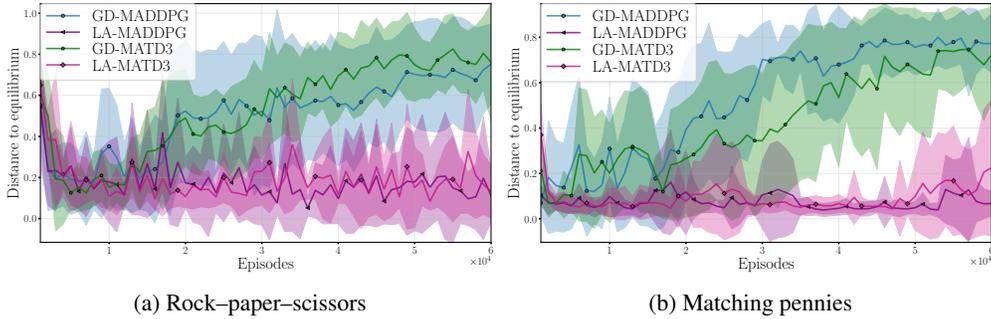


Figure 3: **Comparison between GD-(MADDPG/MATD3) and LA-(MADDPG/MATD3), on Rock-paper-scissors and Matching pennies.** x -axis: training episodes. y -axis: shows the distance between the agents’ policies and the equilibrium policy. Curves are averaged over 10 random seeds, and the shaded regions show ± 1 standard deviation over these 10 seeds.

#Players	GD vs. GD	GD vs. LA	LA vs. LA	LA vs. GD
$n = 3$	2.99 ± 1.73	$2.14 \pm .91 \downarrow$	5.44 ± 1.27	$7.41 \pm 1.75 \uparrow$
$n = 5$	15.69 ± 7.18	$15.5 \pm 5.32 \downarrow$	14.58 ± 5.45	$22.58 \pm 8.97 \uparrow$

Table 1: **Competition between agents trained with different algorithms.** Means and standard deviations (over 5 seeds) of adversary reward in MPE: Predator Prey, on 100 test environments. When GD’s opponent is switched to LA, its reward decreases, and vice versa. See Section 5.2.

over $t = 25$ steps. MP is a two-player, two-action game ($m = 2$) with MNE $\pi_{MP}^* = (\frac{1}{2}, \frac{1}{2})$; Even wins on a match, Odd on a mismatch. We report $\|\pi - \pi^*\|^2$ for learned policies in both games.

MPE: Predator-Prey and Physical Deception. We evaluate one competitive environment from the Multi-Agent Particle Environments (MPE) benchmark (Lowe et al., 2017). *Predator-Prey*, has p good agents, m adversaries, and l landmarks, where good agents are faster but penalized for being caught or going out of bounds, while adversaries collaborate to capture them. We use $n \equiv (p + m) \in [3, 5]$, and $l = 2$.

Methods. We evaluate our methods against the baseline, GD-MARL (GD): MADDPG/MATD3 with Adam (Kingma & Ba, 2015) as base optimizer B . When referring to LA-based methods, we will indicate the k values for each lookahead level in brackets. For instance, LA (10, 1000) denotes a two-level lookahead where $k^{(1)} = 10$ and $k^{(2)} = 1000$. EG denotes the EG method, and refer to it analogously. Further details on hyperparameters are provided in Appendix E.

5.2 RESULTS

2-player games: RPS & MP. Figures 3a–3b plot verage distance between learned and equilibrium policies. We show that GD-MARL drifts/diverges, while LA-MARL reliably contracts to near-equilibrium and outperforms GD. MADDPG and MATD3 perform similarly, with MATD3 showing lower seed variance. Figure 1 shows how Lookahead dampens rotations in MP, reducing the distance to equilibrium relative to GD.

MPE: Predator–Prey. Table 1 reports mean adversary rewards versus good agents. We train GD-MATD3 (baseline) and LA-MATD3 with five seeds and cross-play them. The results show that LA agents perform well against LA and consistently outperform GD in head-to-head matchups.

Summary. Overall, our results indicate the following: (i) VI-based methods consistently outperform their respective baselines, by effectively handling the rotational dynamics. (ii) LA-VI outperforms the other methods.

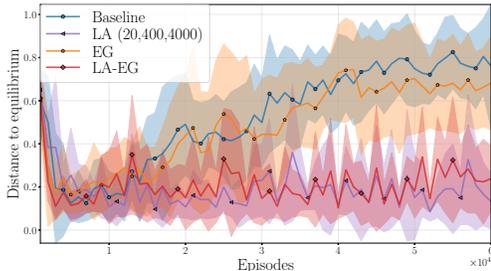


Figure 4: **Comparison between GD, LA, EG, and LA-EG optimization methods on the Rock-paper-scissors game.** x -axis: training episodes. y -axis: squared norm of the learned policy probabilities relative to the equilibrium. Shaded regions indicate one standard deviation across the random seeds used. EG uses solely one extrapolation, and thus, as a method, is very close to GD; refer to Section 5.2 for a discussion.

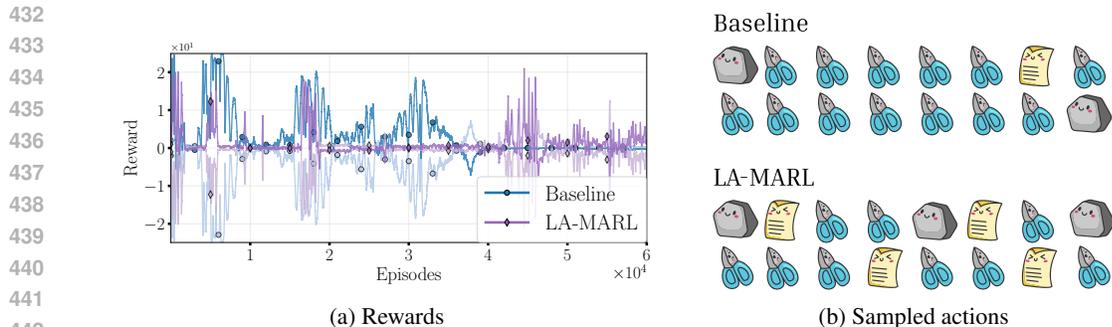


Figure 5: **Rewards (left) vs. sampled actions from learned policies (right), of (LA-)MADDPG in the Rock-paper-scissors game.** The baseline has saturating rewards (in the last part), however, that is not indicative of the agents’ performances. Refer to Section 5.2 for a discussion, and Figure 16 for more detailed plots and larger action samples.

Comparison among VI methods & insights from GANs. EG ((EG))—while convergent for monotone VIs—yields only marginal gains over GD here, reflecting its small local correction. In contrast, Lookahead (LA) imposes stronger contraction, markedly improving stability and convergence; deeper nesting helps—especially preventing last-iterate divergence (Fig. 6)—but too many levels slow learning. Three nested levels offered the best trade-off (Fig. 4). This mirrors GAN results (Chavdarova et al., 2021): EG gives modest benefits, whereas more contractive VI updates consistently perform better. Overall, the evidence indicates strong rotational dynamics; for highly competitive tasks, prefer more contractive VI methods (e.g., multi-level LA).

Limitations of Rewards as a Metric in MARL. While saturating rewards are commonly used in MARL, few works challenge their reliability (e.g., Bowling, 2004) as reward convergence doesn’t equal optimality. Rewards can saturate with suboptimal policies—e.g., in Fig. 5, baseline agents tie by repeating a subset of actions, masking failure to reach equilibrium. Conversely, LA-MADDPG shows no saturation yet learns near-equilibrium via randomization. For this, stronger metrics are needed, especially when the true equilibrium is unknown (Appendix F.9).

6 DISCUSSION

This work studied competitive reinforcement learning through the lens of variational inequalities (VIs). Through second derivative analysis, we showed why even single-agent actor-critic methods inherit game-like dynamics with rotational components, and that these effects amplify in multi-agent settings. Building on this operator view, we unified common multi-agent policy-gradient and actor-critic methods under a VI formulation and instantiated VI-inspired optimizers as *drop-in* updates. In particular, we introduced a general algorithmic framework (Algorithm 1) that augments existing MARL operators with Lookahead (LA) and/or Extragradient (EG) steps, yielding computationally efficient methods tailored to practical CTDE MARL. We provided consistent empirical evidence that integrating VI-based methods improves stability and last-iterate behavior. Conceptually, our LA-MARL and EG-MARL variants act as lightweight, general wrappers that transfer tools from VI theory into practical actor-critic and multi-agent policy-gradient frameworks.

Limitations. While rotations offer a useful and tractable lens on instability, they are not the sole driver of failure modes in RL. Function-approximation error, exploration strategies, and environment stochasticity all interact with the rotational dynamics we study. Understanding this interplay more systematically, and isolating when rotational effects dominate versus when other factors are critical, is an important direction for follow-up work.

Future directions. Our results suggest several additional promising avenues: further developing VI-inspired optimizers tailored specifically to stochastic and off-policy settings; extending analysis beyond centralized critics to decentralized settings; and designing new evaluation protocols that go beyond reward-based metrics to better capture model performance. This work highlights the importance of optimization dynamics in competitive RL and aims to bridge VI theory and deep MARL practice, opening up further research at the intersection of game theory, optimization, and reinforcement learning.

REFERENCES

- 486
487
488 Johann J.H. Ackermann, Volker Gabler, Takayuki Osa, and Masashi Sugiyama. Reducing overesti-
489 mation bias in multi-agent domains using double centralized critics. *ArXiv:1910.01465*, 2019.
- 490
491 Ioannis Anagnostides, Ioannis Panageas, Gabriele Farina, and Tuomas Sandholm. Optimistic policy
492 gradient in multi-player markov games with a single controller: Convergence beyond the minty
493 property. In *AAAI Conference on Artificial Intelligence*, 2023.
- 494
495 Gürdal Arslan and Serdar Yüksel. Decentralized q-learning for stochastic teams and games. *IEEE*
496 *Transactions on Automatic Control*, 62:1545–1558, 2015.
- 497
498 Robert J. Aumann. *Repeated Games with Incomplete Information*, volume 1 of *MIT Press Books*.
499 The MIT Press, December 1995. ISBN ARRAY(0x6d381400).
- 500
501 Wäiss Azizian, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. A tight and unified
502 analysis of gradient-based methods for a whole spectrum of differentiable games. In *AISTATS*, pp.
503 2863–2873, 2020.
- 504
505 Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In *ICML*,
506 pp. 528–537, 2020.
- 507
508 David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel.
509 The mechanics of n-player differentiable games. In *ICML*, 2018.
- 510
511 Alain Bensoussan and Jacques Louis Lions. Contrôle impulsif et inéquations quasi variationnelles
512 stationnaires. In *Comptes Rendus Academie Sciences Paris 276*, pp. 1279–1284, 1973a.
- 513
514 Alain Bensoussan and Jacques Louis Lions. Nouvelle formulation de problèmes de contrôle im-
515 pulsif et applications. In *Comptes Rendus Academie Sciences Paris 276*, pp. 1189–1192,
516 1973b.
- 517
518 Alain Bensoussan and Jacques Louis Lions. Nouvelle méthodes en contrôle impulsif. In *Applied*
519 *Mathematics and Optimization*, pp. 289–312, 1974.
- 520
521 D. Bertsekas. *Rollout, Policy Iteration, and Distributed Reinforcement Learning*. Athena scientific
522 optimization and computation series. Athena Scientific, 2021. ISBN 9781886529076.
- 523
524 Matteo Bettini, Amanda Prorok, and Vincent Moens. Benchmark: Benchmarking multi-agent
525 reinforcement learning. *Journal of Machine Learning Research*, 25(217):1–10, 2024a.
- 526
527 Matteo Bettini, Amanda Prorok, and Vincent Moens. Benchmark: Benchmarking multi-agent
528 reinforcement learning. *arXiv:2312.01472*, 2024b.
- 529
530 Radu Ioan Bot, Ernő Robert Csetnek, and Phan Tu Vuong. The forward-backward-forward method
531 from continuous and discrete perspective for pseudo-monotone variational inequalities in Hilbert
532 spaces. *arXiv:1808.08084*, 2020.
- 533
534 Radu Ioan Bot, Ernő Robert Csetnek, and Dang-Khoa Nguyen. Fast OGDA in continuous and
535 discrete time. *arXiv:2203.10947*, 2022.
- 536
537 Michael Bowling. Convergence and no-regret in multiagent learning. In *NIPS*, volume 17. MIT Press,
538 2004.
- 539
540 Jingjing Bu, Lillian J. Ratliff, and Mehran Mesbahi. Global convergence of policy gradient for
541 sequential zero-sum linear quadratic dynamic games. *arXiv:1911.04672*, 2019.
- 542
543 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Tight last-iterate convergence of the extragrad-
544 ient method for constrained monotone variational inequalities. *arXiv:2204.09228*, 2022.
- 545
546 Tatjana Chavdarova, Gauthier Gidel, François Fleuret, and Simon Lacoste-Julien. Reducing noise in
547 GAN training with variance reduced extragradient. In *NeurIPS*, 2019.
- 548
549 Tatjana Chavdarova, Matteo Pagliardini, Sebastian U Stich, François Fleuret, and Martin Jaggi.
550 Taming GANs with Lookahead-Minmax. In *ICLR*, 2021.

- 540 Tatjana Chavdarova, Michael I. Jordan, and Manolis Zampetakis. Last-iterate convergence of saddle
541 point optimizers via high-resolution differential equations. In *Minimax Theory and its Applications*,
542 2023.
- 543 Tatjana Chavdarova, Tong Yang, Matteo Pagliardini, and Michael I. Jordan. A primal-dual approach
544 for solving variational inequalities with general-form constraints. In *ICLR*, 2024.
- 545 Xuyang Chen and Lin Zhao. On the convergence of continuous single-timescale actor-critic. In *Forty-
546 second International Conference on Machine Learning*, 2025. URL [https://openreview.
547 net/forum?id=pV7hSmGJXP](https://openreview.net/forum?id=pV7hSmGJXP).
- 548 Richard W. Cottle and George B. Dantzig. Complementary pivot theory of mathematical programming.
549 *Linear Algebra and its Applications*, 1(1):103–125, 1968. ISSN 0024-3795.
- 550 Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and
551 constrained min-max optimization. In *ITCS*, 2019.
- 552 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with
553 optimism. In *ICLR*, 2018.
- 554 Constantinos Daskalakis, Dylan J. Foster, and Noah Golowich. Independent policy gradient methods
555 for competitive reinforcement learning. In *NeurIPS*. Curran Associates Inc., 2020a. ISBN
556 9781713829546.
- 557 Constantinos Daskalakis, Stratis Skoulakis, and Manolis Zampetakis. The complexity of constrained
558 min-max optimization. *arXiv:2009.09623*, 2020b.
- 559 Jelena Diakonikolas. Halpern iteration for near-optimal and parameter-free monotone inclusion and
560 strong solutions to variational inequalities. In *COLT*, volume 125, 2020.
- 561 Francisco Facchinei and J. S. Pang. Finite-dimensional variational inequalities and comple-
562 mentarity problems. 2003a. URL [https://api.semanticscholar.org/CorpusID:
563 118457067](https://api.semanticscholar.org/CorpusID:118457067).
- 564 Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional Variational Inequalities and Complemen-
565 tarity Problems*. Springer, 2003b.
- 566 Wentao Fan. A comprehensive analysis of game theory on multi-agent reinforcement. *Highlights in
567 Science, Engineering and Technology*, 85:77–88, 03 2024. doi: 10.54097/gv6fpz53.
- 568 M. Fazel, R. Ge, S. M. Kakade, and M. Mesbahi. Global convergence of policy gradient methods for
569 the linear quadratic regulator. In *International Conference on Machine Learning*, 2018.
- 570 Jakob Foerster, Nantas Nardelli, Gregory Farquhar, Triantafyllos Afouras, Philip H. S. Torr, Pushmeet
571 Kohli, and Shimon Whiteson. Stabilising experience replay for deep multi-agent reinforcement
572 learning. In *ICML*, pp. 1146–1155. JMLR, 2017.
- 573 Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.
574 Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial
575 intelligence*, volume 32, 2018.
- 576 Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel Huang,
577 Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics.
578 In *AISTATS*, 2019.
- 579 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates
580 for no-regret learning in multi-player games. In *NeurIPS*, 2020a.
- 581 Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is
582 slower than averaged iterate in smooth convex-concave saddle point problems. In *COLT*, pp.
583 1758–1784, 2020b.
- 584 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
585 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.

- 594 Eduard Gorbunov, Nicolas Loizou, and Gauthier Gidel. Extragradient method: $\mathcal{O}(1/K)$ last-iterate
595 convergence for monotone variational inequalities and connections with cocoercivity. In *AISTATS*,
596 2022.
- 597
598 Rihab Gorsane, Omayma Mahjoub, Ruan de Kock, Roland Dubb, Siddarth Singh, and Arnu Pretorius.
599 Towards a standardised performance evaluation protocol for cooperative marl. *arXiv:2209.10485*,
600 2022.
- 601 Yuxiang Guan, Giulio Salizzoni, Maryam Kamgarpour, and Tyler H. Summers. A policy iteration
602 algorithm for n-player general-sum linear quadratic dynamic games. *arXiv:2410.03106*, 2024.
603
- 604 Junsoo Ha and Gunhee Kim. On convergence of lookahead in smooth games. In *Proceedings of The*
605 *25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings*
606 *of Machine Learning Research*, pp. 4659–4684. PMLR, 28–30 Mar 2022.
- 607 Ben Hambly, Renyuan Xu, and Huining Yang. Policy gradient methods find the nash equilibrium in
608 n-player general-sum linear-quadratic games. *Journal of Machine Learning Research*, 24(139):
609 1–56, 2023.
- 611 Markus Holzleitner, Lukas Gruber, José Arjona-Medina, Johannes Brandstetter, and Sepp Hochreiter.
612 Convergence proof for actor-critic methods applied to ppo and rudder, 2020.
- 613 Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *ICML*,
614 2018.
- 616 E. Janovskaja. Equilibrium points in polymatrix games. *Lithuanian Mathematical Journal*, 8(2):
617 381–384, Apr. 1968. doi: 10.15388/LMJ.1968.20224.
- 618 Jiechuan Jiang and Zongqing Lu. I2q: A fully decentralized q-learning algorithm. In *NeurIPS*,
619 volume 35, pp. 20469–20481, 2022.
- 621 R. E. Kalman. Contributions to the theory of optimal control. *Boletin de la Sociedad Matematica*
622 *Mexicana*, 5, 1960.
- 624 Fivos Kalogiannis, Jingming Yan, and Ioannis Panageas. Learning equilibria in adversarial team
625 markov games: a nonconvex-hidden-concave min-max optimization problem. In *Proceedings of*
626 *the 38th International Conference on Neural Information Processing Systems*, 2024.
- 627 Tosio Kato. *Perturbation Theory for Linear Operators*. 1966.
- 628
629 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- 630
631 Vijay Konda and John Tsitsiklis. Actor-critic algorithms. In S. Solla, T. Leen, and K. Müller (eds.),
632 *NIPS*, volume 12. MIT Press, 1999.
- 633 Galina Michailovna Korpelevich. The extragradient method for finding saddle points and other
634 problems. *Matecon*, 1976.
- 636 Georgios Kotsalis, Guanghui Lan, and Tianjiao Li. Simple and optimal methods for stochastic
637 variational inequalities, ii: Markovian noise and policy evaluation in reinforcement learning, 2021.
638
- 639 Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Master’s thesis, 2009.
- 640
641 Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong
642 Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv:2109.11251*,
643 2022.
- 644 Navdeep Kumar, Priyank Agrawal, Giorgia Ramponi, Kfir Yehuda Levy, and Shie Mannor. On the
645 convergence of single-timescale actor-critic, 2025.
- 646
647 P. Lancaster and L. Rodman. *Algebraic Riccati Equations*. Oxford Mathematical Monographs.
Clarendon Press, 1995. ISBN 9780198537953.

- 648 Marc Lanctot, John Schultz, Neil Burch, Max Olan Smith, Daniel Hennes, Thomas Anthony, and
649 Julien Perolat. Population-based evaluation in repeated rock-paper-scissors as a benchmark for
650 multiagent reinforcement learning. *arXiv:2303.03196*, 2023.
- 651
652 Hepeng Li and Haibo He. Multiagent trust region policy optimization. *IEEE transactions on neural
653 networks and learning systems*, PP, 04 2023.
- 654 Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. Robust multi-agent
655 reinforcement learning via minimax deep deterministic policy gradient. In *AAAI*, volume 33, pp.
656 4213–4220, 2019.
- 657
658 Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence
659 of generative adversarial networks. *AISTATS*, 2019.
- 660
661 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Manfred Otto Heess, Tom Erez,
662 Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement
663 learning. *CoRR*, abs/1509.02971, 2015.
- 664
665 Michael L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *ICML*,
666 *ICML’94*, pp. 157–163. Morgan Kaufmann Publishers Inc., 1994. ISBN 1558603352.
- 667
668 Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-
669 critic for mixed cooperative-competitive environments. *NIPS*, 2017.
- 670
671 Songtao Lu, Kaiqing Zhang, Tianyi Chen, Tamer Başar, and Lior Horesh. Decentralized policy
672 gradient descent ascent for safe multi-agent reinforcement learning. *AAAI*, 35(10):8767–8775,
673 2021. doi: 10.1609/aaai.v35i10.17062.
- 674
675 Jeffrey Ma, Alistair Letcher, Florian Schäfer, Yuanyuan Shi, and Anima Anandkumar. Polymatrix
676 competitive gradient descent. *arXiv:2111.08565*, 2021.
- 677
678 Yu. Malitsky. Projected reflected gradient methods for monotone variational inequalities. *SIAM
679 Journal on Optimization*, 25:502–520, 2015.
- 680
681 Laëtitia Matignon, Guillaume J. Laurent, and Nadine Le Fort-Piat. Independent reinforcement
682 learners in cooperative Markov games: a survey regarding coordination problems. *Knowledge
683 Engineering Review*, 27(1):1–31, March 2012. doi: 10.1017/S026988891200057.
- 684
685 Eric Mazumdar, Lillian J. Ratliff, Michael I. Jordan, and S. Shankar Sastry. Policy-gradient algorithms
686 have no guarantees of convergence in linear quadratic games. In *International Conference on
687 Autonomous Agents and Multiagent Systems*, 2020.
- 688
689 Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which Training Methods for GANs do
690 actually Converge? In *ICML*, 2018.
- 691
692 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-
693 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,
694 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
695 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.
696 *Nature*, 518(7540):529–533, 2015. ISSN 00280836.
- 697
698 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim
699 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
700 learning, 2016.
- 701
702 Dov Monderer and Lloyd S Shapley. Potential games. *Games and economic behavior*, 1996.
- 703
704 Shayegan Omidshafiei, Jason Pazis, Chris Amato, Jonathan P. How, and John Vian. Deep de-
705 centralized multi-task multi-agent reinforcement learning under partial observability. In *ICML*,
706 2017.
- 707
708 Thomas Pethick, Wanyun Xie, and Volkan Cevher. Stable nonconvex-nonconcave training via linear
709 interpolation. In *NeurIPS*, 2023.

- 702 David Pfau and Oriol Vinyals. Connecting generative adversarial networks and actor-critic methods,
703 2017.
- 704 Leonid Denisovich Popov. A modification of the arrow–hurwicz method for search of saddle points.
705 *Mathematical Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- 706
707 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
708 Wiley & Sons, Inc., USA, 1st edition, 1994. ISBN 0471619779.
- 709
710 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder de Witt, Gregory Farquhar, Jakob Foerster,
711 and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent
712 reinforcement learning. *arXiv:1803.11485*, 2018.
- 713
714 Ralph Tyrrell Rockafellar. Monotone operators associated with saddle-functions and minimax
715 problems. *Nonlinear functional analysis*, 18(part 1):397–407, 1970.
- 716
717 Mihaela Rosca, Yan Wu, Benoit Dherin, and David G. T. Barrett. Discretization drift in two-player
718 games. In *ICML*, 2021.
- 719
720 Masoud Roudneshin, Jalal Arabneydi, and Amir G. Aghdam. Reinforcement learning in nonzero-sum
721 linear quadratic deep structured games: Global convergence of policy optimization. In *IEEE
722 Conference on Decision and Control (CDC)*, pp. 512–517. IEEE Press, 2020. doi: 10.1109/
723 CDC42340.2020.9303950.
- 724
725 Ernest K. Ryu, Kun Yuan, and Wotao Yin. Ode analysis of stochastic gradient methods with optimism
726 and anchoring for minimax problems. *arXiv:1905.10899*, 2019.
- 727
728 Muhammed O. Sayin, Kaiqing Zhang, David S. Leslie, Tamer Basar, and Asuman E. Ozdaglar.
729 Decentralized q-learning in zero-sum markov games. In *NeurIPS*, 2021.
- 730
731 John Schulman, Sergey Levine, P. Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region policy
732 optimization. *ArXiv:1502.05477*, 2015.
- 733
734 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
735 optimization algorithms. *ArXiv:1707.06347*, 2017.
- 736
737 Lloyd S. Shapley. Stochastic games*. *Proceedings of the National Academy of Sciences*, 39(10):
738 1095–1100, 1953. doi: 10.1073/pnas.39.10.1095.
- 739
740 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
741 Deterministic policy gradient algorithms. In *Proceedings of the 31st International Conference on
742 International Conference on Machine Learning - Volume 32*, pp. 1–387–1–395. JMLR.org, 2014.
- 743
744 Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran:
745 Learning to factorize with transformation for cooperative multi-agent reinforcement learning.
746 *arXiv:1905.05408*, 2019.
- 747
748 Riccardo Spica, Davide Falanga, Eric Cristofalo, Eduardo Montijano, Davide Scaramuzza, and
749 Mac Schwager. A real-time game theoretic planner for autonomous two-player drone racing.
750 *arXiv:1801.02302*, 2018.
- 751
752 Sriram Srinivasan, Marc Lanctot, Vinicius Zambaldi, Julien Perolat, Karl Tuyls, Remi Munos, and
753 Michael Bowling. Actor-critic policy optimization in partially observable multiagent environments.
754 In *NeurIPS*, volume 31, 2018.
- 755
756 Guido Stampacchia. Formes bilineaires coercitives sur les ensembles convexes. *Académie des
757 Sciences de Paris*, 258:4413–4416, 1964.
- 758
759 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech M. Czarnecki, Vinícius Flores Zambaldi,
760 Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel.
761 Value-decomposition networks for cooperative multi-agent learning. *ArXiv:1706.05296*, 2017.
- 762
763 J Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S
764 Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym for
765 multi-agent reinforcement learning. *NeurIPS*, 34:15032–15043, 2021.

- 756 Kiran Koshy Thekumparampil, Niao He, and Sewoong Oh. Lifted primal-dual method for bilinearly
757 coupled smooth minimax optimization. In *AISTATS*, 2022.
- 758
- 759 Haoxing Tian, Alex Olshevsky, and Ioannis Paschalidis. Convergence of actor-critic with multi-layer
760 neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
761 URL <https://openreview.net/forum?id=QlfGOVD5PO>.
- 762 Paul Tseng. On linear convergence of iterative methods for the variational inequality problem.
763 *Journal of Computational and Applied Mathematics*, 60:237–252, 1995. ISSN 0377-0427.
- 764
- 765 Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle
766 Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, John Quan, Stephen
767 Gaffney, Stig Petersen, Karen Simonyan, Tom Schaul, Hado van Hasselt, David Silver, Tim-
768 othy Lillicrap, Kevin Calderone, Paul Keet, Anthony Brunasso, David Lawrence, Anders Ek-
769 ermo, Jacob Repp, and Rodney Tsing. Starcraft ii: A new challenge for reinforcement learning.
770 *arXiv:1708.04782*, 2017.
- 771 Zhijian Wang, Bin Xu, and Hai-Jun Zhou. Social cycling and conditional responses in the rock-paper-
772 scissors game. *Scientific Reports*, 4(1), 2014. ISSN 2045-2322. doi: 10.1038/srep05830.
- 773
- 774 Chen-Yu Wei, Yi-Te Hong, and Chi-Jen Lu. Online reinforcement learning in stochastic games. In
775 *NIPS*, volume 30. Curran Associates, Inc., 2017.
- 776 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
777 learning. *Machine Learning*, 8(3–4):229–256, 1992. ISSN 0885-6125. doi: 10.1007/BF00992696.
- 778
- 779 Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-
780 move markov games using function approximation and correlated equilibrium. In *COLT*, volume
781 125, pp. 3674–3682. PMLR, 09–12 Jul 2020.
- 782 Tong Yang, Michael I. Jordan, and Tatjana Chavdarova. Solving constrained variational inequalities
783 via an interior point method. In *ICLR*, 2023.
- 784
- 785 Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game
786 theoretical perspective. *arXiv:2011.00583*, 2021.
- 787
- 788 Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre M. Bayen, and Yi Wu. The surprising
789 effectiveness of mappo in cooperative, multi-agent games. *ArXiv:2103.01955*, 2021.
- 790
- 791 Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The
792 surprising effectiveness of ppo in cooperative, multi-agent games. *arXiv:2103.01955*, 2022.
- 793
- 794 Kaiqing Zhang, Zhuoran Yang, and Tamer Basar. Policy optimization provably converges to nash
795 equilibria in zero-sum linear quadratic games. In *NeurIPS*, volume 32, 2019a.
- 796
- 797 Kaiqing Zhang, Xiangyuan Zhang, Bin Hu, and Tamer Basar. Derivative-free policy optimization for
798 linear risk-sensitive and robust control design: Implicit regularization and sample complexity. In
799 *NeurIPS*, 2021.
- 800
- 801 Michael Zhang, James Lucas, Jimmy Ba, and Geoffrey E Hinton. Lookahead optimizer: k steps
802 forward, 1 step back. In *NeurIPS*, 2019b.
- 803
- 804 Liyuan Zheng, Tanner Fiez, Zane Alumbaugh, Benjamin Chasnov, and Lillian J. Ratliff. Stackelberg
805 actor-critic: Game-theoretic reinforcement learning algorithms. *arXiv:2109.12286*, 2021.
- 806
- 807 Hai-Jun Zhou. The rock–paper–scissors game. *Contemporary Physics*, 57(2):151–163, March 2015.
808 ISSN 1366-5812. doi: 10.1080/00107514.2015.1026556.
- 809
- 805 M. Zhou, Y. Guan, M. Hayajneh, K. Niu, and C. Abdallah. Game theory and machine learning in
806 uavs-assisted wireless communication networks: A survey. *arXiv:2108.03495*, 2021.

810 A EXTENDED RELATED WORKS DISCUSSION

811
812
813 Our work is primarily grounded in three key areas: Actor critic (AC) methods, Multi-Agent Rein-
814forcement Learning (MARL) and Variational Inequalities (VIs), which we discuss next. Additionally,
815we extend our discussion on related works on Linear-Quadratic (LQ) games and discuss relevant
816literature on independent MARL.

817 **Actor critic.** Many AC algorithms have been proposed for single-agent reinforcement learning (e.g.,
818Lillicrap et al., 2015; Schulman et al., 2015; 2017). In this work, we focus on the rotational, game-like
819dynamics that can arise in AC. Pfau & Vinyals (2017) formalize the connection between AC and GANs
820as a bilevel optimization problem, highlighting the adversarial structure. Zheng et al. (2021) make
821the cycling phenomenon explicit by casting actor–critic updates as a two-player Stackelberg game
822and show that leader–follower updates can mitigate such rotations. Actor–critic (AC) convergence is
823well understood in classic two–time-scale settings with linear function approximation (critic faster
824than actor), where the iterates converge to a stationary point under standard conditions (Konda &
825Tsitsiklis, 1999). More recent results push beyond linearity and tabular assumptions: Holzeitner et al.
826(2020) prove (local) convergence of deep actor–critic schemes applied to PPO and RUDDER using
827a two–time-scale stochastic-approximation framework that accommodates episodic sampling and
828policies that become greedier over time. On the single-time-scale front, Kumar et al. (2025) analyze
829the coupled actor/critic recursions and establish global convergence to an ϵ -optimal policy with
830sample complexity) (finite-state discounted MDPs), refining prior bounds and specifying step-size
831schedules for both actor and critic. In continuous state–action spaces, Chen & Zhao (2025) give a
832finite-time convergence analysis for the canonical single-time-scale AC with Markovian sampling
833via a unified Lyapunov framework, closing a gap between practice and theory. Complementing
834these, Tian et al. (2023) prove convergence guarantees for AC with deep multi-layer neural networks,
835quantifying how network width and critic approximation quality control the stationarity gap.

836
837 **Multi-Agent Reinforcement Learning (MARL).** Various MARL algorithms have been devel-
838oped (Lowe et al., 2017; Iqbal & Sha, 2018; Ackermann et al., 2019; Yu et al., 2021), with some
839extending existing single-agent reinforcement learning (RL) methods (Rashid et al., 2018; Son et al.,
8402019; Yu et al., 2022; Kuba et al., 2022). Lowe et al. (2017) extend an actor-critic algorithm to the
841MARL setting using the *centralized training decentralized execution* framework. In the proposed
842algorithm, named *multi-agent deep deterministic policy gradient (MADDPG)*, each agent in the game
843consists of two components: an *actor* and a *critic*. The actor is a policy network that has access
844only to the local observations of the corresponding agent and is trained to output appropriate actions.
845The critic is a value network that receives additional information about the policies of other agents
846and learns to output the Q-value; see Section 2. After a phase of experience collection, a batch is
847sampled from a replay buffer and used for training the agents. To our knowledge, all deep MARL
848implementations rely on either stochastic gradient descent or *Adam* optimizer (Kingma & Ba, 2015)
849to train all networks. Game theory and MARL share many foundational concepts, and several studies
850explore the relationships between the two fields (Yang & Wang, 2021; Fan, 2024), with some using
851game-theoretic approaches to model MARL problems (Zheng et al., 2021). This work proposes
852incorporating game-theoretic techniques into the optimization process of existing MARL methods to
853determine if these techniques can enhance MARL optimization.

854 Li et al. (2019) introduced an algorithm called *M3DDPG*, aimed at enhancing the robustness of
855learned policies. Specifically, it focuses on making policies resilient to worst-case adversarial
856perturbations, as well as uncertainties in the environment or the behaviors of other agents. Several
857works rely on two-player zero-sum Markov games to study the regret of an agent relative to a perfect
858adversary. For instance, Bai & Jin (2020) introduces self-play algorithms for online learning—
859the *Value Iteration with Upper/Lower Confidence Bound (VI-ULCB)* and an explore-then-exploit
860algorithm—and show the respective regret bounds. In addition to the online setting, Xie et al. (2020)
861also consider the offline setting where they propose using Coarse Correlated Equilibria (CCE) instead
862of Nash Equilibrium (NE) and derive concentration bounds for CCEs. Convergence in MARL is
863challenging due to complex interactions and non-stationarity among agents. While multi-agent
actor-critic methods are widely used (Bettini et al., 2024a), their optimization and convergence
properties remain underexplored, making this an open problem.

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, 2003b, 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, 2003b; Daskalakis et al., 2018; Liang & 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 & 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 2 and Appendix B.1.1. In this work, we primarily focus on the nested Lookahead-VI (LA) method, which has achieved state-of-the-art results on the CIFAR-10 (Krizhevsky, 2009) benchmark for generative adversarial networks (Goodfellow et al., 2014).

VIs and Markov games. Variational inequalities (VIs) offer a convenient framework for expressing equilibrium-seeking problems in games and multi-agent learning. Recently, this perspective has been used to analyze learning dynamics in multi-player Markov games. In that line of work, Anagnostides et al. (2023) focus on single-controller Markov games and design a new optimistic policy-gradient algorithm that converges to stationary ϵ -Nash equilibria under a generalized Minty condition (Facchinei & Pang, 2003a). Similarly, Kalogiannis et al. (2024) study two-player adversarial team Markov games, reformulating equilibrium computation as a nonconvex-concave min-max problem and exploiting hidden monotone-operator structure to obtain algorithms to compute the Nash equilibrium in a Markov game with polynomial complexity. In both cases, the main contribution is *algorithmic*: they introduce new learning rules tailored to specific structured classes of Markov games and analyze their convergence. By contrast, our work keeps standard CTDE actor-critic algorithms (e.g., MADDPG, MATD3) fixed and instead changes only the *optimization layer*, using the VI framework to (i) reveal and quantify rotational dynamics in existing actor-critic updates and (ii) plug in VI-inspired optimizers such as Lookahead and Extragradient as drop-in replacements, thereby isolating the impact of optimization dynamics on stability and convergence behavior.

Another line of work studies simple and optimal methods for stochastic VIs under Markovian noise, with reinforcement-learning policy evaluation as a motivating application. Kotsalis et al. (2021) design several TD-like algorithms—one similar to an extragradient type—for solving VIs with inexact information and under the existence of a stochastic oracle and provide non-asymptotic convergence rates in this setting. In contrast, we work in practical deep multi-agent reinforcement learning setup: we cast practical actor-critic and centralized MARL methods as VI problems in joint parameter space and use VI-based optimization as a generic wrapper to improve their empirical behavior, rather than proposing new TD schemes for policy evaluation.

General-sum linear quadratic (LQ) games. In LQ games, each agent’s action linearly impacts the state process, and their goal is to minimize a quadratic cost function dependent on the state and control actions of both themselves and their opponents. LQ games are widely studied as they admit global Nash equilibria (NE), which can be analytically computed using coupled algebraic Riccati equations (Lancaster & Rodman, 1995).

918 Several works establish global convergence for policy gradient methods in zero-sum settings. Zhang
919 et al. (2019a) propose an alternating policy update with projection for deterministic infinite-horizon
920 settings, proving sublinear convergence. Bu et al. (2019) study leader-follower policy gradient in
921 a deterministic setup, and showing sublinear convergence. Zhang et al. (2021) study the sample
922 complexity of policy gradient with alternating policy updates.

923 For the deterministic n -agent setting, Mazumdar et al. (2020) showed that policy gradient methods
924 fail to guarantee even local convergence. Roudneshin et al. (2020) prove global convergence for
925 policy gradient in a *mean-field* LQ game with infinite horizon and stochastic dynamics. Hambly et al.
926 (2023) show that the *natural policy gradient* method achieves global convergence in finite-horizon
927 general-sum LQ games, provided that a certain condition on an added noise to the system is satisfied.
928 Recently, Guan et al. (2024) proposed a policy iteration method for the infinite horizon setting.

929
930 **Independent MARL.** In independent MARL, each agent learns its policy independently, without
931 direct access to the observations, actions, or rewards of other agents (Matignon et al., 2012; Foerster
932 et al., 2017). Each agent treats the environment as stationary and ignores the presence of other agents,
933 effectively treating them as part of the environment.

934 (Daskalakis et al., 2020a) study two-agent zero-sum MARL setting of independent learning algorithms.
935 The authors show that if both players run policy gradient methods jointly, their policies will converge
936 to a min-max equilibrium of the game, as long as their learning rates follow a two-timescale rule.
937 (Arslan & Yüksel, 2015) propose a decentralized Q -learning algorithm for MARL setting where
938 agents have limited information and access solely of their local observations and rewards. Jiang
939 & Lu (2022) proposes a decentralized algorithm. Sayin et al. (2021) explore a decentralized Q -
940 learning algorithm for zero-sum Markov games, where two competing agents learn optimal policies
941 without direct coordination or knowledge of each other’s strategies. Each agent relies solely on local
942 observations and rewards, updating their Q -values independently while interacting in a stochastic
943 environment. (Lu et al., 2021) study a decentralized cooperative multi-agent setting with coupled
944 safety constraints.

945 Wei et al. (2017) rely on the framework of average-reward stochastic games to model single player
946 with a perfect adversary, yielding a two-player zero-sum game, in a Markov environment, and study
947 the regret bound.

948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

B EXTENDED BACKGROUND

In this section, we further discuss VIs, and provide additional background and relevant algorithms.

B.1 VI DISCUSSION

Variational Inequality. We first recall the definition of VIs. A $\mathbf{VI}(F, \mathcal{Z})$ is defined as:

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

where $F: \mathcal{Z} \rightarrow \mathbb{R}^d$, referred to as the *operator*, is continuous, and \mathcal{Z} is a subset of the Euclidean d -dimensional space \mathbb{R}^d .

When $F \equiv \nabla f$ and f is a real-valued function $f: \mathcal{Z} \rightarrow \mathbb{R}$, the problem **VI** is equivalent to standard minimization. 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). We refer the reader to (Facchinei & 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

$$\mathbf{z} \equiv \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_n \end{bmatrix} \in \mathbb{R}^d, \quad \text{with} \quad d = \sum_{i=1}^n d_i.$$

Each agent $i \in [n]$ aims to optimize its objective $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$, which, in the general case, depends on all players' strategies. Then, finding an equilibrium in this game is equivalent to solving a **VI** where the operator F corresponds to:

$$F_{n\text{-agents}}(\mathbf{z}) \equiv \begin{bmatrix} \nabla_{\mathbf{z}_1} f_1(\mathbf{z}) \\ \vdots \\ \nabla_{\mathbf{z}_n} f_n(\mathbf{z}) \end{bmatrix}. \quad (F_{n\text{-agents}})$$

An instructive way to understand the difference between non-rotational and rotational learning dynamics is to consider the second-derivative matrix $J: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ referred herein as the *Jacobian*. For the above ($F_{n\text{-agents}}$) problem the Jacobian is as follows:

$$J_{n\text{-agents}}(\mathbf{z}) \equiv \begin{bmatrix} \nabla_{\mathbf{z}_1}^2 f_1(\mathbf{z}) & \nabla_{\mathbf{z}_1 \mathbf{z}_2}^2 f_1(\mathbf{z}) & \dots & \nabla_{\mathbf{z}_1 \mathbf{z}_n}^2 f_1(\mathbf{z}) \\ \vdots & \vdots & \dots & \vdots \\ \nabla_{\mathbf{z}_n \mathbf{z}_1}^2 f_n(\mathbf{z}) & \nabla_{\mathbf{z}_n \mathbf{z}_2}^2 f_n(\mathbf{z}) & \dots & \nabla_{\mathbf{z}_n}^2 f_n(\mathbf{z}) \end{bmatrix}. \quad (J_{n\text{-agents}})$$

Notably, unlike in minimization, where the second-derivative matrix—the so-called *Hessian*—is always symmetric, the Jacobian is not necessarily symmetric. Hence, its eigenvalues may belong to the complex plane. In some cases, the Jacobian of the associated vector field 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.

In Section D we will also rely on a more general problem, referred to as the *Quasi Variational Inequality*.

Quasi Variational Inequality. Given a map $F: \mathcal{X} \rightarrow \mathbb{R}^n$ —herein referred as an *operator*—the goal is to:

$$\text{find } \mathbf{x}^* \quad \text{s.t.} \quad \langle \mathbf{x} - \mathbf{x}^*, F(\mathbf{x}^*) \rangle \geq 0, \quad \forall \mathbf{x} \in \mathcal{K}(\mathbf{x}^*), \quad (\text{QVI})$$

where $\mathcal{K}(\mathbf{x}) \subseteq \mathbb{R}^d$ is a point-to-set mapping from \mathbb{R}^d into subsets of \mathbb{R}^d such that for every $\mathbf{x} \in \mathcal{X}$, $\mathcal{K}(\mathbf{x}) \subseteq \mathbb{R}^d$ which can be possibly empty.

In other words, the constraint set for QVIs depends on the variable \mathbf{x} . This contrasts with a standard variational inequality (**VI**), where the constraint set \mathcal{K} is fixed and does not depend on \mathbf{x} . QVIs were introduced in a series of works by Bensoussan & Lions (1973a;b; 1974).

B.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 $z, z' \in \mathbb{R}^d$.

The following provides an alternative but equivalent formulation of LA. LA was originally proposed for minimization problems (Zhang et al., 2019b).

LA equivalent formulation. We can equivalently write (LA) as follows. At a step t : (i) a copy of the current iterate \tilde{z}_t is made: $\tilde{z}_t \leftarrow z_t$, (ii) \tilde{z}_t is updated $k \geq 1$ times using B , yielding \tilde{z}_{t+k} , and finally (iii) the actual update z_{t+1} is obtained as a *point that lies on a line between* the current z_t iterate and the predicted one \tilde{z}_{t+k} :

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

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) (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.

B.1.2 PSEUDOCODE FOR EXTRAGRADIENT

In Algorithm 2 outlines the *Extragradient* optimizer (Korpelevich, 1976), which we employ in EG-MARL. 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 2 Extragradient optimizer; Can be used as B in algorithm 1.

```

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

```

B.1.3 PSEUDOCODE FOR NESTED LOOKAHEAD FOR A TWO-PLAYER GAME

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

In the given algorithm, the actor and critic parameters are first updated using a gradient-based optimizer. At interval $k^{(1)}$, backtracking is done between the current parameters and first-level copies (slow weights) and they get updated. At interval $k^{(2)} = c_j k^{(1)}$ backtracking is performed again with second-level copies (slower weights), updating both first- and second-level copies with the averaged version.

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

```

1: Input: number of episodes  $t$ , learning rates  $\eta_\theta, \eta_w$ , initial weights  $\{\theta, \theta^{(1)}, \theta^{(2)}\}$  and
    $\{w, w^{(1)}, w^{(2)}\}$ , LA hyperparameters: levels  $l = 2, (k^{(1)}, k^{(2)})$  and  $(\alpha_\theta, \alpha_w)$ , losses  $\ell^\theta, \ell^w$ ,
   real–data distribution  $p_d$ , noise–data distribution  $p_z$ .
2: for  $r \in 1, \dots, t$  do
3:    $x \sim p_d, z \sim p_z$ 
4:    $w \leftarrow w - \eta_w \nabla_w \ell^w(w, x, z)$  (update  $w$ )
5:    $\theta \leftarrow \theta - \eta_\theta \nabla_\theta \ell^\theta(\theta, x, z)$  (update  $\theta$ )
6:   if  $r \% k^{(1)} == 0$  then
7:      $w \leftarrow w^{(1)} + \alpha_w (w - w^{(1)})$  (backtracking on interpolated line  $w^{(1)}, w$ )
8:      $\theta \leftarrow \theta^{(1)} + \alpha_\theta (\theta - \theta^{(1)})$  (backtracking on interpolated line  $\theta^{(1)}, \theta$ )
9:      $(\theta^{(1)}, w^{(1)}) \leftarrow (\theta, w)$  (update slow checkpoints)
10:  end if
11:  if  $r \% k^{(2)} == 0$  then
12:     $w \leftarrow w^{(2)} + \alpha_w (w - w^{(2)})$  (backtracking on interpolated line  $w^{(2)}, w$ )
13:     $\theta \leftarrow \theta^{(2)} + \alpha_\theta (\theta - \theta^{(2)})$  (backtracking on interpolated line  $\theta^{(2)}, \theta$ )
14:     $(\theta^{(2)}, w^{(2)}) \leftarrow (\theta, w)$  (update super-slow checkpoints)
15:     $(\theta_{(1)}, w_{(1)}) \leftarrow (\theta, w)$  (update slow checkpoints)
16:  end if
17: end for
18: Output:  $\theta^{(2)}, w^{(2)}$ 

```

B.2 MARL ALGORITHMS

B.2.1 DETAILS ON THE MADDPG ALGORITHM

The MADDPG algorithm (Lowe et al., 2017) is outlined in Algorithm 4. An empty replay buffer \mathcal{D} is initialized to store experiences. In each episode, the environment is reset and agents choose actions to perform accordingly. After, 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 -value of agent i . Deep Q-learning (Mnih et al., 2015) is then used to update the critic network; lines 20–21. Then, the agents’ policy network is optimized using policy gradient; refer to 23. Finally, following each learning iteration, the target networks are updated towards current actor and critic networks using a fraction τ . Then the process repeats until the end of training.

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.

B.2.2 MATD3 ALGORITHM

We provide a psuedo code for MATD3 algorithm from (Ackermann et al., 2019) in algorithm 5. As discussed in the main section, MATD3 was introduced as an improvement to MADDPG and follows a similar structure, except for the learning steps. After sampling a batch from the replay buffer \mathcal{D} , both critics of each agent are updated using Deep Q-learning, with the target computed using the minimum of the two critic values (notice the difference in lines 20 and 20of the two algorithms). The actor networks are then updated via policy gradient, using only the Q -value from the first critic; see line 24.

B.2.3 COUNTERFACTUAL MULTI-AGENT POLICY GRADIENTS (COMA, (FOERSTER ET AL., 2018))

COMA is an actor-critic multi-agent algorithm based on the CTDE paradigm, with one centralized critic and n decentralized actors. Additionally, COMA directly addresses the credit assignment problem in multi-agent settings by: (*i*) computing a counterfactual baseline for each agent $b_i(s, \mathbf{a}_{-i})$,

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

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_{rand} before learning, learning interval t_{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 η_θ, η_w , optimizer B (e.g., Adam), discount factor γ , soft update parameter τ .

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 $a_i = \mu_i(\sigma_i) + \mathcal{N}_t$ using current policy and exploration
 noise
 12: **end if**
 13: Execute actions $\mathbf{a} = (a_1, \dots, a_n)$, observe rewards \mathbf{r} and new state \mathbf{x}' *(apply actions and
record results)*
 14: replay buffer $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \mathbf{x}')$
 15: $\mathbf{x} \leftarrow \mathbf{x}'$
 16: *(apply learning step if applicable)*
 17: **if** $\text{step} \% t_{\text{learn}} = 0$ **then**
 18: **for all** agent $i \in 1, \dots, n$ **do**
 19: sample batch $\mathcal{B} : \{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}'^j)\}_{j=1}^{|\mathcal{B}|}$ from \mathcal{D}
 20: $y^j \leftarrow r^j + \gamma \mathbf{Q}^\mu(\mathbf{x}'^j, a'_1, \dots, a'_n)$, where $a'_k = \bar{\mu}_k(\sigma_k^j)$
 21: Update critic by minimizing the loss (using optimizer B):

$$\ell(w_i) = \frac{1}{|\mathcal{B}|} \sum_j \left(y^j - \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j) \right)^2$$

 22: Update actor policy using policy gradient formula and optimizer B

$$\nabla_{\theta_i} J \approx \frac{1}{|\mathcal{B}|} \sum_j \nabla_{\theta_i} \mu_i(\sigma_i^j) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j)$$
, where $a_i = \mu_i(\sigma_i^j)$
 24: **end for**
 25: **for all** agent $i \in [n]$ **do**
 26: $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$ *(update target networks)*
 27: $\bar{w}_i \leftarrow \tau w_i + (1 - \tau) \bar{w}_i$
 28: **end for**
 29: **end if**
 30: $\text{step} \leftarrow \text{step} + 1$
 31: **until** environment terminates
 32: **end for**
 33: **Output:** θ, w

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Algorithm 5 Pseudocode for MATD3 (Ackermann et al., 2019).

```

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  $\theta \equiv \{\theta_i\}_{i=1}^n$ , both critic networks,  $\{Q_{i,1}, Q_{i,2}\}_{i=1}^n$  with initial weights
    $w \equiv \{w_{i,1}, w_{i,2}\}_{i=1}^n$ , learning rates  $\eta_\theta, \eta_w$ , optimizer  $B$  (e.g., Adam), discount factor  $\gamma$ , soft
   update parameter  $\tau$ , policy update frequency  $p$ .
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  $a_i = \mu_i(o_i) + \epsilon$  using current policy with some exploration
        noise
12:    end if
13:    Execute actions  $\mathbf{a} = (a_1, \dots, a_n)$ , observe rewards  $\mathbf{r}$  and new state  $\mathbf{x}'$  (apply actions and
        record results)
14:    replay buffer  $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \mathbf{x}')$ 
15:     $\mathbf{x} \leftarrow \mathbf{x}'$ 
16:    (apply learning step if applicable)
17:    if  $\text{step} \% t_{\text{learn}} = 0$  then
18:      for all agent  $i \in [n]$  do
19:        sample batch  $\{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}'^j)\}_{j=1}^{|\mathcal{B}|}$  from  $\mathcal{D}$ 
20:         $y^j \leftarrow r_i^j + \gamma \min_{m=1,2} Q_{i,m}^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j)$ , where  $a_k^j = \bar{\mu}_k(o_k^j) + \epsilon$ 
21:        Update both critics,  $m = 1, 2$  by minimizing the loss (using optimizer  $B$ ):
22:          
$$\ell(w_{i,m}) = \frac{1}{|\mathcal{B}|} \sum_j \left( y^j - Q_{i,m}^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j) \right)^2$$

23:        if  $\text{step} \% p = 0$  then
24:          Update actor policy using policy gradient formula and optimizer  $B$ 
25:           $\nabla_{\theta_i} J \approx \frac{1}{|\mathcal{B}|} \sum_j \nabla_{\theta_i} \mu_i(o_i^j) \nabla_{a_i} Q_{i,1}^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j)$ , where  $a_i = \mu_i(o_i^j)$ 
26:           $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$  (update target networks)
27:           $\bar{w}_{i,m} \leftarrow \tau w_{i,m} + (1 - \tau) \bar{w}_{i,m}$ 
28:        end if
29:      end for
30:    end if
31:     $\text{step} \leftarrow \text{step} + 1$ 
32:  until environment terminates
33: end for
Output:  $\theta, w$ 

```

1242 (ii) using this baseline to estimate the advantage A_i of the chosen action over all others in \mathcal{A}_i , and
1243 (iii) leveraging this advantage to update individual policies. This ensures that policy updates reflect
1244 each agent’s true contribution to the overall reward.

1245 1246 B.2.4 MULTI-AGENT TRUST REGION POLICY OPTIMIZATION (MATRPO, (LI & HE, 2023)) 1247

1248 Trust Region Policy Optimization (TRPO, Schulman et al., 2015) is a policy optimization method
1249 that ensures stable updates by constraining policy changes within a trust region. This constraint is
1250 enforced using the KL-divergence, and the update step is computed using natural gradient descent.

1251 Extending TRPO to the cooperative multi-agent setting introduces challenges due to non-stationarity.
1252 To address this, MATRPO employs a centralized critic, represented by a central value function $V(\mathbf{s})$,
1253 which leverages shared information among agents to estimate the Generalized Advantage Estimator
1254 (GAE) A_i . The advantage function is then used in the policy gradient update, while ensuring that the
1255 KL-divergence constraint is respected, maintaining stable and coordinated learning across agents.

1256 1257 B.2.5 MULTI-AGENT PROXIMAL POLICY OPTIMIZATION (MAPPO, YU ET AL., 2021)

1258 One of the widely used algorithms in practice is MAPPO, an extension of Proximal Policy Opti-
1259 mization (PPO, Schulman et al., 2017) to the multi-agent setting. Similar to TRPO, PPO ensures
1260 that policy updates remain within a small, stable region, but instead of enforcing a KL-divergence
1261 constraint, it uses clipping. This clipping mechanism simplifies the optimization process, allowing
1262 updates to be performed efficiently using standard gradient ascent methods.

1263 MAPPO is an on-policy algorithm that employs a centralized critic while maintaining decentralized
1264 actor networks for each agent. Its critic update follows the same rule as MATRPO, but for the policy
1265 update, it optimizes a clipped surrogate objective, which restricts the policy update step size, ensuring
1266 stable and efficient learning.

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

C OFF-POLICY DETERMINISTIC ACTOR-CRITIC METHODS ANALYSIS & MISSING PROOFS

C.1 SINGLE-AGENT SETTING

The following describes the considered setting in Lemma 1 in more detail. For simplicity, we consider a single state s' for each trajectory.

Setting 1 (off-policy deterministic AC setting). *Consider the following off-policy deterministic actor-critic setting (Silver et al., 2014):*

1. Action $a \in \mathbb{R}$
2. State $\mathbf{s} \in \mathbb{R}^d$
3. State-action features $\phi \in \mathbb{R}^c$ defined as:

$$\phi(\mathbf{s}, a) \triangleq f_\phi(\mathbf{s}) + a m(\mathbf{s}) \in \mathbb{R}^c, \quad \text{with } f_\phi: \mathbb{R}^d \rightarrow \mathbb{R}^c, m: \mathbb{R}^d \rightarrow \mathbb{R}^c.$$

Notice that:

$$\frac{\partial \phi(\mathbf{s}, a)}{\partial a} = m(\mathbf{s}), \quad \left. \frac{\partial \phi(\mathbf{s}, a)}{\partial \theta} \right|_{a \text{ fixed}} = 0.$$

4. Critic parameters $\mathbf{w} \in \mathbb{R}^c$
5. Actor parameters $\theta \in \mathbb{R}^d$
6. Critic (linear Q -value): $Q_{\mathbf{w}}(\mathbf{s}, a) = \langle \mathbf{w}, \phi(\mathbf{s}, a) \rangle$
7. Actor (deterministic linear policy): $\pi_{\theta}(\mathbf{s}) \triangleq \langle \theta, \mathbf{s} \rangle$
8. Batch of experiences: $\{\mathbf{s}_i, a_i, r_i, \mathbf{s}'_i\}_{i=1}^B$

Taking $\mathbf{s} \in \mathbb{R}^d$, of the same dimension as θ , is without loss of generality for the following analysis, since one can use a different dimension and a feature map. We also introduce the following shorthands:

$$\phi_i \triangleq \phi(\mathbf{s}_i, a_i), \tag{\phi_i}$$

$$a'_i \triangleq \pi_{\theta}(\mathbf{s}'_i) = \langle \theta, \mathbf{s}'_i \rangle, \tag{a'_i}$$

$$\phi'_i \triangleq \phi(\mathbf{s}'_i, a'_i) = f_\phi(\mathbf{s}'_i) + a'_i m(\mathbf{s}'_i), \tag{\phi'_i}$$

$$\delta_i \triangleq \langle \mathbf{w}, \phi_i \rangle - r_i - \gamma \langle \mathbf{w}, \phi'_i \rangle. \tag{\delta_i}$$

We next state the complete statement of the informal Lemma 1 and provide its proof.

Lemma 1 (Complex eigenvalues in off-policy deterministic AC). *Consider the off-policy deterministic AC model of setting 1 with $\gamma > 0$. Its associated (simultaneous) gradient operator $F: \mathbb{R}^{d+c} \rightarrow \mathbb{R}^{d+c}$ is:*

$$F\left(\begin{bmatrix} \mathbf{w} \\ \theta \end{bmatrix}\right) = \begin{bmatrix} F_{\mathbf{w}} \\ F_{\theta} \end{bmatrix} = \begin{bmatrix} \frac{2}{B} \sum_{i=1}^B \delta_i (\phi_i - \gamma \phi'_i) \\ -\frac{1}{B} \sum_{i=1}^B \mathbf{s}_i \langle \mathbf{w}, m(\mathbf{s}_i) \rangle \end{bmatrix}. \tag{1}$$

Its associated Jacobian $J: \mathbb{R}^{c+d} \rightarrow \mathbb{R}^{(c+d) \times (c+d)}$ has the following block form:

$$J = \begin{bmatrix} J_{\mathbf{w}\mathbf{w}} & J_{\mathbf{w}\theta} \\ J_{\theta\mathbf{w}} & J_{\theta\theta} \end{bmatrix}, \quad \text{with } J_{\mathbf{w}\mathbf{w}} \succeq 0,$$

with:

$$\begin{aligned}
 (1) \quad J_{\mathbf{w}\mathbf{w}} &= \frac{\partial F_{\mathbf{w}}}{\partial \mathbf{w}} = \frac{2}{B} \sum_{i=1}^B (\phi_i - \gamma \phi'_i) (\phi_i - \gamma \phi'_i)^\top \in \mathbb{R}^{c \times c} \\
 (2) \quad J_{\mathbf{w}\theta} &= \frac{\partial F_{\mathbf{w}}}{\partial \theta} = -\frac{2\gamma}{B} \sum_{i=1}^B \left[(\phi_i - \gamma \phi'_i) \mathbf{s}'_i{}^\top \mathbf{w}^\top m(\mathbf{s}'_i) + m(\mathbf{s}'_i) \mathbf{s}'_i{}^\top \delta_i \right] \in \mathbb{R}^{c \times d} \\
 (3) \quad J_{\theta\mathbf{w}} &= \frac{\partial F_{\theta}}{\partial \mathbf{w}} = -\frac{1}{B} \sum_{i=1}^B \mathbf{s}_i m(\mathbf{s}_i)^\top \in \mathbb{R}^{d \times c} \\
 (4) \quad J_{\theta\theta} &= \frac{\partial F_{\theta}}{\partial \theta} = \mathbf{0}_{d \times d}.
 \end{aligned}$$

Then:

- (i) (Pure cross-term case) If $J_{\mathbf{w}\mathbf{w}} = \mathbf{0}$, then $\text{Spec}(J) = \{\pm \sqrt{\lambda_k(J_{\mathbf{w}\theta} J_{\theta\mathbf{w}})}\}_k$. In particular, if $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ has a negative eigenvalue, J has a purely imaginary conjugate eigenpair.
- (ii) (Persistence) If $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ has a negative eigenvalue, there exists $\varepsilon > 0$ such that for all $\|J_{\mathbf{w}\mathbf{w}}\| < \varepsilon$ (operator norm), J has a non-real conjugate eigenpair.
- (iii) (Symmetry case) If $J_{\mathbf{w}\theta} = J_{\theta\mathbf{w}}^\top$ and $J_{\mathbf{w}\mathbf{w}} = J_{\mathbf{w}\mathbf{w}}^\top \succeq 0$, then J is symmetric and all eigenvalues are real.

Moreover, $J_{\mathbf{w}\theta} \neq J_{\theta\mathbf{w}}^\top$ unless $\gamma = 0$ or simultaneously $\delta_i \equiv 0$ and $\mathbf{w}^\top m(\mathbf{s}'_i) \equiv 0$ for all i . Hence, for $\gamma > 0$ and generic data/parameters, $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ is sign-indefinite and J exhibits a non-real eigenpair.

Proof of Lemma 1. STEP 1: CRITIC LOSS AND DERIVATIVES

Loss. Given a batch $\{\mathbf{s}_i, a_i, r_i, \mathbf{s}'_i\}_{i=1}^B$, consider the *temporal difference* (TD) loss for the critic, parametrized by \mathbf{w} , defined as:

$$\begin{aligned}
 \mathcal{L}(\mathbf{w}) &= \frac{1}{B} \sum_{i=1}^B \left(\underbrace{\langle \mathbf{w}, \phi(\mathbf{s}_i, a_i) \rangle}_{\hat{y}_i} - r_i - \gamma \underbrace{\langle \mathbf{w}, \phi(\mathbf{s}'_i, \pi_\theta(\mathbf{s}'_i)) \rangle}_{y_i} \right)^2 \\
 &= \frac{1}{B} \sum_{i=1}^B \left(\langle \mathbf{w}, \phi_i \rangle - r_i - \gamma \langle \mathbf{w}, \phi'_i \rangle \right)^2.
 \end{aligned} \tag{2}$$

Notice that the θ -dependence in the critic loss is implicit, via a'_i in ϕ'_i .

First derivative. The first derivative of the critic loss with respect to \mathbf{w} is then:

$$F_{\mathbf{w}}(\mathbf{w}, \theta) \equiv \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{2}{B} \sum_{i=1}^B \delta_i (\phi_i - \gamma \phi'_i) \in \mathbb{R}^c. \tag{F_{\mathbf{w}}}$$

Second derivative. The second derivative matrix for \mathbf{w} is then:

$$J_{\mathbf{w}\mathbf{w}} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}) = \frac{2}{B} \sum_{i=1}^B (\phi_i - \gamma \phi'_i) (\phi_i - \gamma \phi'_i)^\top \in \mathbb{R}^{c \times c}. \tag{J_{\mathbf{w}\mathbf{w}}}$$

$J_{\mathbf{w}\mathbf{w}}$ is symmetric, and due to the outer product structure of $J_{\mathbf{w}\mathbf{w}}$, it is positive semi-definite (PSD); thus, it has non-negative eigenvalues.

Mixed second derivative: derivative of critic's first derivative w.r.t. actor. Since:

$$\frac{\partial \phi(\mathbf{s}'_i, \pi_\theta(\mathbf{s}'_i))}{\partial \theta} = \frac{\partial \phi(\mathbf{s}'_i, a'_i)}{\partial a} \mathbf{s}'_i{}^\top = m(\mathbf{s}'_i) \mathbf{s}'_i{}^\top \in \mathbb{R}^{c \times d},$$

differentiating $(F_{\mathbf{w}})$ w.r.t. θ , and applying the product and chain rules, gives:

$$J_{\mathbf{w}\theta} = \nabla_{\theta} F_{\mathbf{w}} = -\frac{2\gamma}{B} \sum_{i=1}^B \left[(\phi_i - \gamma \phi'_i) \mathbf{s}'_i{}^\top \mathbf{w}^\top m(\mathbf{s}'_i) + m(\mathbf{s}'_i) \mathbf{s}'_i{}^\top \delta_i \right] \in \mathbb{R}^{c \times d}. \tag{J_{\mathbf{w}\theta}}$$

STEP 2: ACTOR DERIVATIVES

First derivative. We adopt the convention that both actor and critic run descent step; thus, for consistency, we use a minus sign to the derivative for θ ; i.e., the actor block is the negative policy gradient. Due to the policy gradient theorem for the deterministic setting, it holds:

$$F_{\theta}(\theta) = -\frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \pi_{\theta}(s_i) \nabla_a Q_{\mathbf{w}}(s_i, a)|_{a=\pi_{\theta}(s_i)}. \quad (3)$$

Since in the considered setting $Q_{\mathbf{w}}(s, a) = \langle \mathbf{w}, \phi(s, a) \rangle = \langle \mathbf{w}, f_{\phi}(s) \rangle + a \langle \mathbf{w}, m(s) \rangle$, we have: $\nabla_a Q_{\mathbf{w}}(s, a) = \mathbf{w}^{\top} m(s)$. Thus, for the first derivative w.r.t. θ for the considered setting we have:

$$\begin{aligned} F_{\theta}(\mathbf{w}, \theta) &= -\frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \pi_{\theta}(s_i) \nabla_a Q_{\mathbf{w}}(s_i, a)|_{a=\pi_{\theta}(s_i)} \\ &= -\frac{1}{B} \sum_{i=1}^B s_i \langle \mathbf{w}, m(s_i) \rangle. \end{aligned} \quad (F_{\theta})$$

Second derivative. Since the actor gradient in equation (F_{θ}) doesn't have θ terms:

$$J_{\theta\theta} = \nabla_{\theta} F_{\theta} = \mathbf{0}_{d \times d} \in \mathbb{R}^{d \times d}. \quad (J_{\theta\theta})$$

Mixed second derivative: derivative of actor's first derivative w.r.t. critic. Differentiating (F_{θ}) with respect to \mathbf{w} yields:

$$J_{\theta\mathbf{w}} = \nabla_{\mathbf{w}} F_{\theta} = -\frac{1}{B} \sum_{i=1}^B s_i m(s_i)^{\top} \in \mathbb{R}^{d \times c}. \quad (J_{\theta\mathbf{w}})$$

STEP 3: SPECTRAL ANALYSIS OF J

(i) (*Pure cross-term case*) If $J_{\mathbf{w}\mathbf{w}} = \mathbf{0}$, then $J^2 = \text{diag}(J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}, J_{\theta\mathbf{w}} J_{\mathbf{w}\theta})$, so the nonzero eigenvalues of J are the signed square roots of those of $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ and $J_{\theta\mathbf{w}} J_{\mathbf{w}\theta}$ (where $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ and $J_{\theta\mathbf{w}} J_{\mathbf{w}\theta}$ share the same nonzero spectrum). Negative eigenvalues of $J_{\mathbf{w}\theta} J_{\theta\mathbf{w}}$ yield purely imaginary eigenvalues of J .

(ii) (*Persistence under small symmetric $J_{\mathbf{w}\mathbf{w}} \succeq 0$*) Let

$$J(\epsilon) = \begin{bmatrix} \epsilon J_{\mathbf{w}\mathbf{w}} & J_{\mathbf{w}\theta} \\ J_{\theta\mathbf{w}} & \mathbf{0} \end{bmatrix},$$

where $\epsilon J_{\mathbf{w}\mathbf{w}}$ is symmetric semidefinite and $\epsilon = \|J_{\mathbf{w}\mathbf{w}}\|$ (operator norm) is small. By classical matrix perturbation theory, eigenvalues depend continuously (analytically) on its matrix entries (e.g., [Kato, 1966](#)); a simple imaginary pair at $J_{\mathbf{w}\mathbf{w}} = \mathbf{0}$ cannot become real under a sufficiently small symmetric perturbation $J_{\mathbf{w}\mathbf{w}} \succeq 0$ without crossing the real axis. Thus the non-real pair persists for small $\|J_{\mathbf{w}\mathbf{w}}\|$, that is for all $0 \leq \epsilon \leq \epsilon_0$ small enough.

(iii) If $J_{\mathbf{w}\theta} = J_{\theta\mathbf{w}}^{\top}$ and $J_{\mathbf{w}\mathbf{w}} = J_{\mathbf{w}\mathbf{w}}^{\top}$, the Jacobian J is symmetric, so its spectrum is real.

The following shows that (iii) forces degenerate conditions. In our setting,

$$J_{\mathbf{w}\theta} = -\frac{2\gamma}{B} \sum_{i=1}^B \left[(\phi_i - \gamma \phi'_i) s_i^{\top} \mathbf{w}^{\top} m(s'_i) + m(s'_i) s_i^{\top} \delta_i \right], \quad J_{\theta\mathbf{w}} = -\frac{1}{B} \sum_{i=1}^B s_i m(s_i)^{\top},$$

and

$$J_{\theta\mathbf{w}}^{\top} = -\frac{1}{B} \sum_{i=1}^B m(s_i) s_i^{\top}.$$

Thus, $J_{\mathbf{w}\theta} = J_{\theta\mathbf{w}}^\top$ would require (i) $\gamma = 0$ (which kills the dependence on the next-state) and that $J_{\theta\mathbf{w}} = \mathbf{0}$, or (ii) the vanishing of both the factors δ_i and $\mathbf{w}^\top m(\mathbf{s}'_i)$ for all i ; these are stringent/degenerate conditions. So generically $J_{\mathbf{w}\theta} \neq J_{\theta\mathbf{w}}^\top$.

□

Remark 1 (Sign convention). Above, we used the sign convention for descending in F , which, due to the negative sign in the derivative of θ , implies that the actor follows the ascent direction. Even if the alternative sign convention is used, the above lemma is invariant under the actor sign flip.

C.1.1 AN INSTANCE OF SINGLE ACTOR-CRITIC SETTING 1: FULL DESCRIPTION AND STATIONARY POINT

This section describes an instance setting 1 in Appendix C.1, which is used for the numerical experiment in Figure 2. The specific example uses $d=c=2$, and a fixed off-policy batch of size $B=2$. The purpose of the example is to give an intuition of the rotational dynamics in actor-critic and that LA can indeed help with mitigating them.

Setup. Suppose we have state and action spaces

$$\mathbf{s} \in \mathbb{R}^2, \quad a \in \mathbb{R},$$

and consider a deterministic linear policy

$$\pi_\theta(\mathbf{s}) = \langle \theta, \mathbf{s} \rangle, \quad \theta \in \mathbb{R}^2,$$

together with a linear critic

$$Q_{\mathbf{w}}(\mathbf{s}, a) = \langle \mathbf{w}, \phi(\mathbf{s}, a) \rangle, \quad \mathbf{w} \in \mathbb{R}^2.$$

The state-action features are defined as in Setting 1 with

$$\phi(\mathbf{s}, a) = f_\phi(\mathbf{s}) + a m(\mathbf{s}), \quad f_\phi(\mathbf{s}) \equiv \mathbf{0}, \quad m(\mathbf{s}) = M \mathbf{s}, \quad M = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}.$$

Thus, m maps any vector to a $+90^\circ$ rotation.

We fix a small off-policy batch of size $B=2$:

$$\{(\mathbf{s}_i, a_i, r_i, \mathbf{s}'_i)\}_{i=1}^2 = \{(e_1, 1, r_1, e_2), (e_2, 1, r_2, e_1)\}, \quad r_1 > 0, r_2 > 0, \quad \gamma \in [0, 1),$$

where $e_1 = [1, 0]^\top$ and $e_2 = [0, 1]^\top$, hence, the future states are *swapped*: $\mathbf{s}'_1 = e_2$ and $\mathbf{s}'_2 = e_1$. We use the shorthand notations from (ϕ_i) – (δ_i) defined previously.

Deriving the operator. We first compute the feature maps induced by the chosen batch. Using $m(e_1) = M e_1 = e_2$, $m(e_2) = M e_2 = -e_1$, and $f_\phi \equiv \mathbf{0}$, the features of the sampled state-action pairs are

$$\phi_1 = \phi(e_1, 1) = m(e_1) = e_2, \quad \phi_2 = \phi(e_2, 1) = m(e_2) = -e_1.$$

Next, we compute the actions at the next states under the current policy:

$$a'_1 = \pi_\theta(e_2) = \langle \theta, e_2 \rangle = \theta_2, \quad a'_2 = \pi_\theta(e_1) = \langle \theta, e_1 \rangle = \theta_1.$$

Using these, the features of the next states and actions are

$$\phi'_1 = \phi(e_2, a'_1) = \theta_2 m(e_2) = -\theta_2 e_1, \quad \phi'_2 = \phi(e_1, a'_2) = \theta_1 m(e_1) = \theta_1 e_2.$$

Substituting those formulas and using direct computations, the difference terms $\phi_i - \gamma \phi'_i$, can be written as:

$$\phi_1 - \gamma \phi'_1 = e_2 - \gamma(-\theta_2 e_1) = \begin{bmatrix} \gamma \theta_2 \\ 1 \end{bmatrix}, \quad \phi_2 - \gamma \phi'_2 = -e_1 - \gamma(\theta_1 e_2) = \begin{bmatrix} -1 \\ -\gamma \theta_1 \end{bmatrix}.$$

We previously defined the TD-error in (δ_i) :

$$\delta_i = \langle \mathbf{w}, \phi_i \rangle - r_i - \gamma \langle \mathbf{w}, \phi'_i \rangle.$$

Writing $\mathbf{w} = (w_1, w_2)^\top$, we obtain

$$\begin{aligned} \delta_1 &= \langle \mathbf{w}, \phi_1 \rangle - r_1 - \gamma \langle \mathbf{w}, \phi'_1 \rangle = \langle \mathbf{w}, e_2 \rangle - r_1 - \gamma \langle \mathbf{w}, -\theta_2 e_1 \rangle \\ &= w_2 - r_1 + \gamma \theta_2 w_1, \end{aligned}$$

$$\begin{aligned} \delta_2 &= \langle \mathbf{w}, \phi_2 \rangle - r_2 - \gamma \langle \mathbf{w}, \phi'_2 \rangle = \langle \mathbf{w}, -e_1 \rangle - r_2 - \gamma \langle \mathbf{w}, \theta_1 e_2 \rangle \\ &= -w_1 - r_2 - \gamma \theta_1 w_2. \end{aligned}$$

Recall from the proof of Lemma 1 setting that:

$$F_{\mathbf{w}}(\mathbf{w}, \boldsymbol{\theta}) = \frac{2}{B} \sum_{i=1}^B \delta_i (\phi_i - \gamma \phi'_i), \quad F_{\boldsymbol{\theta}}(\mathbf{w}, \boldsymbol{\theta}) = -\frac{1}{B} \sum_{i=1}^B \mathbf{s}_i \langle \mathbf{w}, m(\mathbf{s}_i) \rangle.$$

For our toy example, we have $B = 2$, so the factor $2/B$ in $F_{\mathbf{w}}$ equals 1.

Substituting the expressions above, we get the operator of our toy example:

$$\begin{aligned} F_{\mathbf{w}}(\mathbf{w}, \boldsymbol{\theta}) &= \delta_1 (\phi_1 - \gamma \phi'_1) + \delta_2 (\phi_2 - \gamma \phi'_2) \\ &= \delta_1 \begin{bmatrix} \gamma \theta_2 \\ 1 \end{bmatrix} + \delta_2 \begin{bmatrix} -1 \\ -\gamma \theta_1 \end{bmatrix} = \begin{bmatrix} \delta_1 (\gamma \theta_2) - \delta_2 \\ \delta_1 - \delta_2 (\gamma \theta_1) \end{bmatrix}, \end{aligned} \quad (4)$$

$$\begin{aligned} F_{\boldsymbol{\theta}}(\mathbf{w}, \boldsymbol{\theta}) &= -\frac{1}{2} \left(e_1 \langle \mathbf{w}, m(e_1) \rangle + e_2 \langle \mathbf{w}, m(e_2) \rangle \right) \\ &= -\frac{1}{2} \left(e_1 \langle \mathbf{w}, e_2 \rangle + e_2 \langle \mathbf{w}, -e_1 \rangle \right) = -\frac{1}{2} \begin{bmatrix} w_2 \\ -w_1 \end{bmatrix} = -\frac{1}{2} M^\top \mathbf{w}. \end{aligned} \quad (5)$$

Closed-form stationary point $(\mathbf{w}^*, \boldsymbol{\theta}^*)$. We now solve for a stationary point $(\mathbf{w}^*, \boldsymbol{\theta}^*)$ satisfying

$$F_{\boldsymbol{\theta}}(\mathbf{w}^*, \boldsymbol{\theta}^*) = \mathbf{0}, \quad F_{\mathbf{w}}(\mathbf{w}^*, \boldsymbol{\theta}^*) = \mathbf{0}.$$

Actor block. From (5), the condition $F_{\boldsymbol{\theta}}(\mathbf{w}^*, \boldsymbol{\theta}^*) = \mathbf{0}$ reads

$$-\frac{1}{2} \begin{bmatrix} w_2^* \\ -w_1^* \end{bmatrix} = \mathbf{0} \implies \mathbf{w}^* = (0, 0)^\top.$$

Thus, at any stationary point of F , the critic parameters must be zero in this toy setup.

Critic block. With $\mathbf{w}^* = \mathbf{0}$, the TD errors simplify to

$$\delta_1 = -r_1, \quad \delta_2 = -r_2.$$

Plugging these values into (4) gives

$$F_{\mathbf{w}}(\mathbf{0}, \boldsymbol{\theta}) = \begin{bmatrix} -r_1 (\gamma \theta_2) + r_2 \\ -r_1 + r_2 (\gamma \theta_1) \end{bmatrix}.$$

Setting $F_{\mathbf{w}}(\mathbf{0}, \boldsymbol{\theta}^*) = \mathbf{0}$ yields the linear system

$$-r_1 (\gamma \theta_2^*) + r_2 = 0, \quad -r_1 + r_2 (\gamma \theta_1^*) = 0,$$

whose solution is

$$\theta_1^* = \frac{r_1}{\gamma r_2}, \quad \theta_2^* = \frac{r_2}{\gamma r_1}.$$

Collecting the two blocks, we obtain the unique stationary point of F in this toy example:

$$\boxed{\mathbf{w}^* = (0, 0)^\top, \quad \boldsymbol{\theta}^* = \left(\frac{r_1}{\gamma r_2}, \frac{r_2}{\gamma r_1} \right)}.$$

In the symmetric case $r_1 = r_2 = r$, this reduces to

$$\boldsymbol{\theta}^* = (1/\gamma, 1/\gamma).$$

In the numerical example used for Figure 2, we chose $r_1 = r_2 = 3$ and $\gamma = 0.99$ which gives the stationary points marked by the dashed lines in the figure.

C.2 MULTI-AGENT SETTING WITH A CENTRALIZED CRITIC

We extend the single-agent linear setting to n agents with a shared state and scalar actions. Each agent i has its own *centralized* critic with access to the full state and the joint action, and a linear deterministic actor. We keep one unrolled next-state s' per transition, as in the single-agent analysis. We rely on similar shorthand definitions as in Appendix C.1.

Setting 2 (off-policy deterministic n -agent AC; full-access critics). *Consider the following:*

1. **State and joint action:** $\mathbf{s} \in \mathbb{R}^d$, $\mathbf{a} = (a^{(1)}, \dots, a^{(n)}) \in \mathbb{R}^n$ with scalar $a^{(i)} \in \mathbb{R}$.
2. **Per-agent centralized critic (linear in \mathbf{a}):**

$$Q_i(\mathbf{s}, \mathbf{a}; \mathbf{w}_i) = \langle \mathbf{w}_i, \phi_i(\mathbf{s}, \mathbf{a}) \rangle, \quad \phi_i(\mathbf{s}, \mathbf{a}) = f_{\phi_i}(\mathbf{s}) + \sum_{p=1}^n a^{(p)} m_i^{(p)}(\mathbf{s}) \in \mathbb{R}^{c_i},$$

where $m_i^{(p)}: \mathbb{R}^d \rightarrow \mathbb{R}^{c_i}$ and $\mathbf{w}_i \in \mathbb{R}^{c_i}$. Then $\frac{\partial \phi_i(\mathbf{s}, \mathbf{a})}{\partial a^{(p)}} = m_i^{(p)}(\mathbf{s})$ and $\left. \frac{\partial \phi_i(\mathbf{s}, \mathbf{a})}{\partial \theta_p} \right|_{\mathbf{a} \text{ fixed}} = \mathbf{0}$.

3. **Per-agent actor (linear deterministic):**

$$a^{(i)} = \pi_i(\mathbf{s}; \boldsymbol{\theta}_i) \triangleq \langle \boldsymbol{\theta}_i, \mathbf{s} \rangle, \quad \boldsymbol{\theta}_i \in \mathbb{R}^d.$$

4. **Batch and unrolling:** $\mathcal{B} = \{(\mathbf{s}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{s}'^j)\}_{j=1}^{|\mathcal{B}|}$ with $\mathbf{r}^j = (r_1^j, \dots, r_n^j)$; next actions $a^{(p)j} = \langle \boldsymbol{\theta}_p, \mathbf{s}'^j \rangle$, joint $\mathbf{a}^j = (a^{(p)j})_{p=1}^n$. Define the shorthands

$$\phi_i^j = \phi_i(\mathbf{s}^j, \mathbf{a}^j), \quad \phi_i'^j = \phi_i(\mathbf{s}'^j, \mathbf{a}'^j).$$

5. **TD residual (per critic):**

$$\delta_i^j \triangleq \langle \mathbf{w}_i, \phi_i^j \rangle - r_i^j - \gamma \langle \mathbf{w}_i, \phi_i'^j \rangle.$$

Per-critic i TD loss and derivatives.

$$\begin{aligned} \mathcal{L}_i(\mathbf{w}_i) &= \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(\underbrace{\langle \mathbf{w}_i, \phi_i(\mathbf{s}^j, \mathbf{a}^j) \rangle}_{\hat{y}_i^j} - \underbrace{(r_i^j + \gamma \langle \mathbf{w}_i, \phi_i(\mathbf{s}'^j, \mathbf{a}'^j) \rangle)}_{y_i^j} \right)^2 \\ &= \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(\langle \mathbf{w}_i, \phi_i^j \rangle - r_i^j - \gamma \langle \mathbf{w}_i, \phi_i'^j \rangle \right)^2. \end{aligned} \quad (6)$$

$$F_{\mathbf{w}_i}(\mathbf{w}_i, \boldsymbol{\theta}_{1:n}) \equiv \nabla_{\mathbf{w}_i} \mathcal{L}_i(\mathbf{w}_i) = \frac{2}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \delta_i^j (\phi_i^j - \gamma \phi_i'^j) \in \mathbb{R}^{c_i}, \quad (7)$$

$$J_{\mathbf{w}_i \mathbf{w}_i} = \nabla_{\mathbf{w}_i}^2 \mathcal{L}_i(\mathbf{w}_i) = \frac{2}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} (\phi_i^j - \gamma \phi_i'^j) (\phi_i^j - \gamma \phi_i'^j)^\top \succeq \mathbf{0} \in \mathbb{R}^{c_i \times c_i}. \quad (8)$$

Using

$$\frac{\partial \phi_i(\mathbf{s}'^j, \mathbf{a}'^j)}{\partial \theta_p} = \left. \frac{\partial \phi_i(\mathbf{s}'^j, \mathbf{a})}{\partial a^{(p)}} \right|_{\mathbf{a}=\mathbf{a}'^j} \frac{\partial a^{(p)j}}{\partial \theta_p} = m_i^{(p)}(\mathbf{s}'^j) (\mathbf{s}'^j)^\top \in \mathbb{R}^{c_i \times d},$$

the mixed critic-actor block is

$$J_{\mathbf{w}_i \boldsymbol{\theta}_p} = \nabla_{\boldsymbol{\theta}_p} F_{\mathbf{w}_i} = -\frac{2\gamma}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left[(\phi_i^j - \gamma \phi_i'^j) (\mathbf{s}'^j)^\top \underbrace{\mathbf{w}_i^\top m_i^{(p)}(\mathbf{s}'^j)}_{\alpha_{i,p}^j} + m_i^{(p)}(\mathbf{s}'^j) (\mathbf{s}'^j)^\top \delta_i^j \right] \in \mathbb{R}^{c_i \times d}. \quad (9)$$

Actor derivatives (descent-style operator). Each actor p follows the negative policy gradient using its *own* critic Q_p :

$$\nabla_{\mathbf{a}^{(p)}} Q_p(\mathbf{s}, \mathbf{a}; \mathbf{w}_p) = \mathbf{w}_p^\top m_p^{(p)}(\mathbf{s}), \quad \nabla_{\boldsymbol{\theta}_p} \pi_p(\mathbf{s}; \boldsymbol{\theta}_p) = \mathbf{s}.$$

Thus

$$F_{\boldsymbol{\theta}_p}(\mathbf{w}_p, \boldsymbol{\theta}_p) \equiv -\nabla_{\boldsymbol{\theta}_p} J^{(p)} = -\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \mathbf{s}^j \langle \mathbf{w}_p, m_p^{(p)}(\mathbf{s}^j) \rangle \in \mathbb{R}^d, \quad (10)$$

$$J_{\boldsymbol{\theta}_p \boldsymbol{\theta}_q} = \nabla_{\boldsymbol{\theta}_q} F_{\boldsymbol{\theta}_p} = \mathbf{0}_{d \times d} \quad (\text{for all } p, q), \quad (11)$$

$$J_{\boldsymbol{\theta}_p \mathbf{w}_i} = \nabla_{\mathbf{w}_i} F_{\boldsymbol{\theta}_p} = \begin{cases} -\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \mathbf{s}^j m_p^{(p)}(\mathbf{s}^j)^\top \in \mathbb{R}^{d \times c_p}, & \text{if } i = p, \\ \mathbf{0}_{d \times c_i}, & \text{if } i \neq p. \end{cases} \quad (12)$$

Hence, the actor-critic block $J_{\boldsymbol{\theta}_w}$ is block-diagonal across agents.

Operator and Jacobian. Stack parameters as $[\mathbf{w}_1; \dots; \mathbf{w}_n; \boldsymbol{\theta}_1; \dots; \boldsymbol{\theta}_n] \in \mathbb{R}^{c_{\text{tot}} + nd}$ with $c_{\text{tot}} = \sum_i c_i$. Stack the operator as

$$F = \begin{bmatrix} F_{\mathbf{w}_1} \\ \vdots \\ F_{\mathbf{w}_n} \\ F_{\boldsymbol{\theta}_1} \\ \vdots \\ F_{\boldsymbol{\theta}_n} \end{bmatrix}.$$

The Jacobian has the block structure

$$J = \begin{bmatrix} J_{\mathbf{w}\mathbf{w}} & J_{\mathbf{w}\boldsymbol{\theta}} \\ J_{\boldsymbol{\theta}\mathbf{w}} & \mathbf{0} \end{bmatrix}, \quad J_{\mathbf{w}\mathbf{w}} = \text{diag}(J_{\mathbf{w}_1 \mathbf{w}_1}, \dots, J_{\mathbf{w}_n \mathbf{w}_n}),$$

with $J_{\mathbf{w}_i \mathbf{w}_i}$ from (8); $J_{\mathbf{w}\boldsymbol{\theta}}$ the $c_{\text{tot}} \times nd$ block matrix whose (i, p) block is (9); and $J_{\boldsymbol{\theta}\mathbf{w}}$ the $nd \times c_{\text{tot}}$ block matrix whose (p, i) block is (12).

We now state the spectral result and relate it to the single-agent case (same one-step unrolling).

Lemma 2 (Complex eigenvalues and scaling vs. single agent). *Consider Setting 2 with $\gamma > 0$ and the Jacobian*

$$J = \begin{bmatrix} J_{\mathbf{w}\mathbf{w}} & J_{\mathbf{w}\boldsymbol{\theta}} \\ J_{\boldsymbol{\theta}\mathbf{w}} & \mathbf{0} \end{bmatrix}.$$

Let $B_{i,p} \triangleq J_{\mathbf{w}_i \boldsymbol{\theta}_p} \in \mathbb{R}^{c_i \times d}$ and $C_{p,i} \triangleq J_{\boldsymbol{\theta}_p \mathbf{w}_i} \in \mathbb{R}^{d \times c_i}$; define

$$\mathcal{M}_n \triangleq J_{\mathbf{w}\boldsymbol{\theta}} J_{\boldsymbol{\theta}\mathbf{w}} = \sum_{p=1}^n \begin{bmatrix} B_{1,p} \\ \vdots \\ B_{n,p} \end{bmatrix} [C_{p,1} \quad \dots \quad C_{p,n}] \in \mathbb{R}^{c_{\text{tot}} \times c_{\text{tot}}}.$$

Then:

(i) (Pure cross-term case) If $J_{\mathbf{w}\mathbf{w}} = \mathbf{0}$, then

$$\text{spec}(J) = \left\{ \pm \sqrt{\lambda_k(\mathcal{M}_n)} \right\}_k.$$

In particular, if \mathcal{M}_n has a negative eigenvalue, J has a purely imaginary conjugate pair. Because $B_{i,p}$ couples critic i with every actor p via the next-state joint action, while $C_{p,i}$ is nonzero only for $i = p$ (actor p uses its own critic p), \mathcal{M}_n is generically sign-indefinite, hence complex pairs are typical.

(ii) (Persistence) If \mathcal{M}_n has a negative eigenvalue, there exists $t_0 > 0$ such that for all $0 \leq t < t_0$

$$J(t) = \begin{bmatrix} t J_{\mathbf{w}\mathbf{w}} & J_{\mathbf{w}\boldsymbol{\theta}} \\ J_{\boldsymbol{\theta}\mathbf{w}} & \mathbf{0} \end{bmatrix}$$

has a non-real conjugate eigenpair.

(iii) (Relation to single agent) Let \mathcal{M}_1 denote the single-agent product (with one s' unroll). Then

$$\mathcal{M}_n = \mathcal{M}_1 + \sum_{p=2}^n \begin{bmatrix} B_{1,p} \\ \vdots \\ B_{n,p} \end{bmatrix} [C_{p,1} \quad \cdots \quad C_{p,n}].$$

Hence

$$\max_p \rho \left(\begin{bmatrix} B_{1,p} \\ \vdots \\ B_{n,p} \end{bmatrix} [C_{p,1} \quad \cdots \quad C_{p,n}] \right) \leq \rho(\mathcal{M}_n) \leq \sum_{p=1}^n \left\| \begin{bmatrix} B_{1,p} \\ \vdots \\ B_{n,p} \end{bmatrix} \right\| \| [C_{p,1} \quad \cdots \quad C_{p,n}] \|,$$

so in the pure cross-term case

$$\rho(J) = \sqrt{\rho(\mathcal{M}_n)} \geq \max_p \sqrt{\rho \left(\begin{bmatrix} B_{1,p} \\ \vdots \\ B_{n,p} \end{bmatrix} [C_{p,1} \quad \cdots \quad C_{p,n}] \right)}.$$

Under alignment (shared invariant directions / commuting contributions / PSD terms) the spectral radius grows with n (at least as the square root of the sum of aligned eigenvalues).

(iv) (Triangular/degenerate real-spectrum cases) If $J_{w_i \theta_p} = \mathbf{0}$ for all i, p (e.g., $\gamma = 0$, or $\delta_i^j \equiv 0$ and $w_i^\top m_i^{(p)}(s^j) \equiv 0$), J is block upper triangular with real spectrum $\text{spec}(J_{ww}) \cup \{0\}$. If instead $J_{\theta_p w_i} = \mathbf{0}$ for all p, i (e.g., $m_p^{(p)} \equiv 0$), J is block lower triangular, again yielding a real spectrum. If $J_{ww} = J_{ww}^\top \succeq 0$ and $J_{w_i \theta_p} = (J_{\theta_p w_i})^\top$ for all i, p , then J is symmetric and all eigenvalues are real (stringent/degenerate off-policy conditions).

Proof sketch. (i) With $J_{ww} = 0$, $J^2 = \text{diag}(\mathcal{M}_n, J_{\theta w} J_{w \theta})$, so eigenvalues of J are $\pm \sqrt{\lambda_k(\mathcal{M}_n)}$. (ii) Analytic perturbation theory implies a simple imaginary pair at $t = 0$ cannot become real for small $t > 0$ without crossing the real axis. (iii) Bounds follow from subadditivity of the operator norm and $\rho(\sum_p X_p) \geq \max_p \rho(X_p)$ for square matrices $\{X_p\}$ of common size; alignment/commutation yields additive growth of dominant eigenvalues, hence $\rho(\mathcal{M}_n)$ (and thus $\rho(J)$) increases with n . (iv) Triangular and symmetric cases yield real spectra as stated. \square

Remark 2 (How adding agents changes the Jacobian). Relative to the single-agent Jacobian: (a) J_{ww} becomes block-diagonal across critics; (b) $J_{w \theta}$ gains n columns of blocks $J_{w_i \theta_p}$ since each critic i depends on the *full* next-state joint action; (c) $J_{\theta w}$ becomes block-diagonal because actor p depends only on its own critic p . These changes enlarge $\mathcal{M}_n = J_{w \theta} J_{\theta w}$ by a sum of n player-specific terms, making negative eigenvalues (and thus complex pairs) more likely and typically *increasing* the spectral radius compared to the single-agent case with one s' unroll.

Remark 3 (Sign convention). We use a descent-style operator for both critic and actors, matching the single-agent section. Flipping all actor signs multiplies $J_{\theta w}$ by -1 and leaves conclusions about non-symmetry and complex pairs unchanged, since they depend on $J_{w \theta} J_{\theta w}$.

D VI MARL CONVERGENCE, PERSPECTIVES & DETAILS ON THE PROPOSED ALGORITHMS

In this section we extend our discussion on the convergence on VI-MARL operator, then we present the VI operators of additional MARL algorithms within the centralized critic CTDE paradigm. After, we provide detailed versions of Algorithm 1 for MADDPG and MATD3, outlining the full training process when incorporating LA or LA-EG.

D.1 VI MARL CONVERGENCE

We first recall the abstract multi-player operator definition from Appendix B.1. Each agent $i \in [n]$ aims to optimize its objective $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$, which, in the general case, depends on all players' strategies. Then, we have the following operator F :

$$F_{n\text{-agents}}(\mathbf{z}) \equiv \begin{bmatrix} \nabla_{\mathbf{z}_1} f_1(\mathbf{z}) \\ \vdots \\ \nabla_{\mathbf{z}_n} f_n(\mathbf{z}) \end{bmatrix}, \quad (F_{n\text{-agents}})$$

with the game Jacobian as follows:

$$J_{n\text{-agents}}(\mathbf{z}) \equiv \begin{bmatrix} \nabla_{\mathbf{z}_1^2}^2 f_1(\mathbf{z}) & \nabla_{\mathbf{z}_1 \mathbf{z}_2}^2 f_1(\mathbf{z}) & \dots & \nabla_{\mathbf{z}_1 \mathbf{z}_n}^2 f_1(\mathbf{z}) \\ \vdots & \vdots & \dots & \vdots \\ \nabla_{\mathbf{z}_n \mathbf{z}_1}^2 f_n(\mathbf{z}) & \nabla_{\mathbf{z}_n \mathbf{z}_2}^2 f_n(\mathbf{z}) & \dots & \nabla_{\mathbf{z}_n^2}^2 f_n(\mathbf{z}) \end{bmatrix}. \quad (J_{n\text{-agents}})$$

More precisely, for multi-agent actor-critic RL we have the following operator:

$$F_{\text{MAAC}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \ell_i^{\mathbf{w}}(\cdot; \mathbf{w}_i, \boldsymbol{\theta}) \right) \\ \nabla_{\boldsymbol{\theta}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \ell_i^{\boldsymbol{\theta}}(\cdot; \mathbf{w}_i, \boldsymbol{\theta}_i) \right) \\ \vdots \end{bmatrix}, \quad (F_{\text{MAAC}})$$

where the parameter space is $\mathcal{Z} \equiv \mathbb{R}^d$, with $d = \sum_{i=1}^n (d_i^Q + d_i^\mu)$; and MAAC stands for *multi-agent-actor-critic*.

Then, we can notice by computing the Jacobian of the above operator that the eigenvalues are in the complex plane. Applying lookahead results in interpolating the largest eigenvalue (in magnitude) with the point (1,0) in the complex plane, thus reducing the spectral radius of the Jacobian. Furthermore, applying this recursively (nested Lookahead) leads to larger contraction.

To make this more precise, consider the gradient descent operator as a base optimizer

$$T_{GD} \equiv I - \alpha F,$$

where α is the step size vector.

Let λ denote the eigenvalue of $J^{base} \triangleq \nabla T_{GD}(\cdot)$ with largest modulus, i.e. $\rho(J^{base}(\cdot)) = |\lambda|$, let \mathbf{u} be its associated eigenvector: $J^{base} \mathbf{u} = \lambda \mathbf{u}$.

The Jacobian of Lookahead is then:

$$J^{LA} = \nabla F^{LA}(\cdot) = (1 - \alpha)I + \alpha(J^{base})^k.$$

The power k rotates the eigenvector in the complex plane; see (Chavdarova et al., 2021). By noticing that:

$$\begin{aligned} J^{LA} \mathbf{u} &= ((1 - \alpha)I + \alpha(J^{base})^k) \mathbf{u} \\ &= ((1 - \alpha) + \alpha \lambda^k) \mathbf{u}, \end{aligned}$$

we deduce \mathbf{u} is an eigenvector of J^{LA} with eigenvalue $1 - \alpha + \alpha \lambda^k$. Thus, this is strictly closer to the unit ball in the complex plane, increasing the contractiveness.

D.2 VI MARL PERSPECTIVES

In the main text, we introduced the general VI operator for multi-agent actor-critic algorithms (F_{MAAC}) and provided the specific equations for MADDPG in (ℓ_{MADDPG}^w & $\ell_{\text{MADDPG}}^\theta$), with the operator corresponding to:

$$F_{\text{MADDPG}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(r_i^j + \gamma \mathbf{Q}_i^{\bar{\mu}}(\mathbf{x}^j, \mathbf{a}'; \mathbf{w}_i) \Big|_{\mathbf{a}'=\bar{\mu}(\sigma^j)} - \mathbf{Q}_i^\mu(\mathbf{x}^j, \mathbf{a}^j; \mathbf{w}_i) \right)^2 \right) \\ \nabla_{\boldsymbol{\theta}_i} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \mu_i(\sigma_i^j; \boldsymbol{\theta}_i) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j; \mathbf{w}_i) \Big|_{a_i=\mu_i(\sigma_i^j)} \right) \\ \vdots \end{bmatrix}, \quad (F_{\text{MADDPG}})$$

where the parameter space is $\mathcal{Z} \equiv \mathbb{R}^d$, with $d = \sum_{i=1}^n (d_i^Q + d_i^\mu)$.

We now show how update equations for several well-known MARL algorithms—that follow the CTDE paradigm with a centralized critic—can be written as a VI. Our VI-based methods can also be applied to these algorithms using the operators below.

For a more general notation, for each agent $i \in [n]$ we assume: (i) central critic network (one or multiple) that estimates either action value Q -Network(s, \mathbf{a}): $\mathbf{Q}_i(\mathbf{x}_t, \mathbf{a}_t; \mathbf{w}_i)$, or state value V -network(s): $\mathbf{V}_i(\mathbf{x}_t; \mathbf{w}_i)$, and (ii) a decentralized policy network that can be deterministic $\mu_i(\sigma_i; \boldsymbol{\theta}_i)$ or stochastic $\pi_i(\sigma_i; \boldsymbol{\theta}_i)$, depending on the algorithm. Given a batch of experiences \mathcal{B} : $(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}^j)$, sampled from a replay buffer (\mathcal{D}), we provide the necessary equations and the final operator (F) for each of the following popular MARL algorithms.

D.2.1 MATD3

The VI formulation for MATD3 is very similar to MADDPG, except here, for each agent, we have two critic networks; we write: $\mathbf{w}_i \equiv \{\mathbf{w}_{i,1}, \mathbf{w}_{i,2}\}$. Accordingly, target computation for the critic ($Q_{i,m}$) is calculated by taking the minimum of both critic networks, but only the value of critic 1 is used for the actor (policy network) update. We have:

$$F_{\text{MATD3}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_{i,1} \\ \mathbf{w}_{i,2} \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_{i,1}} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(\underbrace{r_i^j + \gamma \min_{m \in \{1,2\}} \mathbf{Q}_{i,m}^{\bar{\mu}}(\mathbf{x}^j, a'_1, \dots, a'_n) \Big|_{\mathbf{a}'=\bar{\mu}(\sigma^j)}}_{\text{target } y_i} - \mathbf{Q}_{i,1}^\mu(\mathbf{x}^j, \mathbf{a}^j; \mathbf{w}_{i,1}) \right)^2 \right) \\ \nabla_{\mathbf{w}_{i,2}} \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(\underbrace{r_i^j + \gamma \min_{m \in \{1,2\}} \mathbf{Q}_{i,m}^{\bar{\mu}}(\mathbf{x}^j, a'_1, \dots, a'_n) \Big|_{\mathbf{a}'=\bar{\mu}(\sigma^j)}}_{\text{target } y_i} - \mathbf{Q}_{i,2}^\mu(\mathbf{x}^j, \mathbf{a}^j; \mathbf{w}_{i,2}) \right)^2 \right) \\ \left(\frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \nabla_{\boldsymbol{\theta}_i} \mu_i(\sigma_i^j; \boldsymbol{\theta}_i) \nabla_{a_i} \mathbf{Q}_{i,1}^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j) \Big|_{a_i=\mu_i(\sigma_i^j)} \right) \\ \vdots \end{bmatrix}. \quad (F_{\text{MATD3}})$$

1836 D.2.2 COMA

1837
1838 In COMA, critic is trained using a $TD(\lambda)$ target (y^λ) computed from a target network parameterized
1839 by $\bar{\mathbf{w}}$ that get updated to main network weights every couple iterations. Given the following
1840 Advantage A_i calculations:

$$1841 A_i(\mathbf{x}, \mathbf{a}) = Q(\mathbf{x}, \mathbf{a}) - b_i(\mathbf{x}, \mathbf{a}_{-i})$$

$$1842 b_i(\mathbf{x}, \mathbf{a}_{-i}) = \sum_{a_i} \pi_i(a_i|\mathbf{o}_i) Q(\mathbf{x}, (a_i, \mathbf{a}_{-i})),$$

1843 the operator for COMA corresponds to:

$$1844 F_{\text{COMA}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \mathbb{E} \left[(y_i^\lambda - Q_i(\mathbf{x}^j, \mathbf{a}; \mathbf{w}_i))^2 \right] \\ \mathbb{E} \left[\nabla_{\boldsymbol{\theta}_i} \sum_i A_i(\mathbf{x}, \mathbf{a}) \log \pi_{\boldsymbol{\theta}_i}(a_i|\mathbf{o}_i) \right] \\ \vdots \end{bmatrix}. \quad (F_{\text{COMA}})$$

1853 D.2.3 MAPPO

1854
1855 As previously noted, MAPPO can be seen as a simplified version of MATRPO. It shares a similar
1856 critic loss with MATRPO but simplifies the actor loss by using a clipped objective instead of a KL
1857 constraint, making the optimization problem more tractable. This allows it to be formulated as a VI,
1858 as shown below:

$$1859 \hat{V}_t = (1 - \lambda) \sum_{n=1}^T \lambda^{n-1} \left(\sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n \mathbf{V}(\mathbf{o}'_t) \right),$$

$$1860 F_{\text{MAPPO}} \left(\begin{bmatrix} \vdots \\ \mathbf{w}_i \\ \boldsymbol{\theta}_i \\ \vdots \end{bmatrix} \right) \equiv \begin{bmatrix} \vdots \\ \nabla_{\mathbf{w}_i} \mathbb{E} \left[\left(\mathbf{V}(\mathbf{x}; \mathbf{w}_i) - \hat{V}_t \right)^2 \right] \\ \nabla_{\boldsymbol{\theta}_i} \mathbb{E} \left[\min \left\{ \frac{\pi_{\boldsymbol{\theta}_i}(a_i|\mathbf{o}_i)}{\pi_{\boldsymbol{\theta}_i^{\text{old}}}(a_i|\mathbf{o}_i)} A_i^{\boldsymbol{\theta}_i^{\text{old}}}, \text{clip} \left(\frac{\pi_{\boldsymbol{\theta}_i}(a_i|\mathbf{o}_i)}{\pi_{\boldsymbol{\theta}_i^{\text{old}}}(a_i|\mathbf{o}_i)}, 1 - \epsilon, 1 + \epsilon \right) A_i^{\boldsymbol{\theta}_i^{\text{old}}} \right\} \right] \\ \vdots \end{bmatrix}. \quad (F_{\text{MAPPO}})$$

1873 D.3 DETAILED ALGORITHMS

1874
1875 Herein we provide procedure NestedLookahead called from algorithm 1 to compute the extrapolations
1876 and after present two pseudocodes considered as extended versions of the main algorithm in algorithm
1877 1; in which we detail how the lookahead approach can be integrated in the training process of
1878 MADDPG and MATD3.

1880 D.3.1 NESTED LOOKAHEAD ALGORITHM

1881
1882 In algorithm 6 below we share a detailed version of Nested lookahead procedure called from
1883 algorithms 1, 7 and 8.

1884 D.3.2 EXTENDED VERSION OF LA-MADDPG PSEUDOCODE

1885
1886 We include an extended version for the LA-MADDPG algorithm without VI notations in algorithm 7.

1888 D.3.3 EXTENDED VERSION OF LA-MATD3 PSEUDOCODE

1889
1889 We include an extended version for the LA-MATD3 algorithm without VI notations in algorithm 8.

Algorithm 6 Pseudocode for LA-VI, called from Algorithm 1. Updates the parameters in-place.

```

1890 1: procedure NESTEDLOOKAHEAD:
1891 2:   Input: #agents  $n$ , episode counter  $e$ , actor and critic weights and snapshots:
1892    $\{(\theta_i, \theta_i^{(1)}, \dots, \theta_i^{(l)})\}_{i=1}^n$  and  $\{(\mathbf{w}_i, \mathbf{w}_i^{(1)}, \dots, \mathbf{w}_i^{(l)})\}_{i=1}^n$ , LA hyperparameters: levels  $l$ ,  $(k^{(1)},$ 
1893    $\dots, k^{(l)})$  and  $(\alpha_\theta, \alpha_w)$ .
1894 3:   for all  $j \in [l]$  do
1895 4:     if  $e \% k^{(j)} == 0$  then
1896 5:       for all agent  $i \in [n]$  do
1897 6:          $\mathbf{w}_i \leftarrow \mathbf{w}_i^{(j)} + \alpha_w(\mathbf{w}_i - \mathbf{w}_i^{(j)})$  LA  $j^{\text{th}}$  level
1898 7:          $\theta_i \leftarrow \theta_i^{(j)} + \alpha_\theta(\theta_i - \theta_i^{(j)})$ 
1899 8:          $(\theta_i^{(1)}, \dots, \theta_i^{(j)}, \mathbf{w}_i^{(1)}, \dots, \mathbf{w}_i^{(j)}) \leftarrow (\{\theta_i\}_{\times j}, \{\mathbf{w}_i\}_{\times j})$  Update copies up to  $j^{\text{th}}$ 
1900 9:       end for
1901 10:    end if
1902 11:  end for
1903 12: end procedure

```

E DETAILS ON THE IMPLEMENTATION

We used the configurations and hyperparameters from the original MADDPG paper for our implementation. For completeness, these are listed in Table 2. We ran $t = 60000$ training episodes for all environments, with a maximum of 25 environment steps (*step*) per episode.

In all experiments, we used a 2-layer MLP with 64 units per layer. ReLU activation was applied between layers for both the policy and value networks of all agents.

E.1 HYPERPARAMETER SELECTION FOR LOOKAHEAD

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

Summary.

- We observed two- or three-level of Lookahead outperform single-level Lookahead (figure 6).
- Each level $j \in [l]$ has different k , denoted here with $k^{(j)}$. These should be selected as multiple of the selected k for the level before, that is, $k^{(j)} = c_j \cdot k^{(j-1)}$, where c_j is positive integer.
- We observed that for the innermost lookahead, small values for $k^{(1)}$, such as smaller than or equal to 50, perform better than using large values. For the outer $k^{(j)}$, $j > 1$ large values work well, such as in the range between 5 – 10 for the c_j ,
- We typically used $\alpha = 0.5$, and we observed lower values, such as $\alpha = 0.3$, give better performances then $\alpha > 0.5$.

Discussion.

- To give an intuition regarding the above-listed conclusions, small values for $k^{(1)}$ 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 $k^{(1)}$ 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.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

Algorithm 7 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  $B$  (e.g., Adam), discount factor  $\gamma$ , lookahead hyperparameters
 $\mathcal{L} \equiv (l, \{k^{(j)}\}_{j=1}^l, \alpha_\theta, \alpha_w)$ , soft update parameter  $\tau$ .
2: Initialize:
3:   Replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
4:   LA parameters:  $\phi \leftarrow \{\theta\}_{\times l}, \{w\}_{\times l}$  (store snapshots for LA)
5: for all episode  $e \in 1, \dots, t$  do
6:    $x \leftarrow \text{Sample}(\mathcal{E})$  (sample from environment  $\mathcal{E}$ )
7:    $step \leftarrow 1$ 
8:   repeat
9:     if  $e \leq t_{\text{rand}}$  then
10:      for each agent  $i, a_i \sim \mathcal{A}_i$  (sample actions randomly)
11:     else
12:      for each agent  $i$ , select action  $a_i$  using current policy and exploration
13:     end if
14:      (apply actions and record results)
15:     Execute actions  $\mathbf{a} = (a_1, \dots, a_n)$ , observe rewards  $\mathbf{r}$  and new state  $\mathbf{x}'$ 
16:     replay buffer  $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \mathbf{x}')$ 
17:      $\mathbf{x} \leftarrow \mathbf{x}'$ 
18:     (apply learning step if applicable)
19:     if  $step \% t_{\text{learn}} = 0$  then
20:       for all agents  $i \in 1, \dots, n$  do
21:         sample batch  $\{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}'^j)\}_{j=1}^{|\mathcal{B}|}$  from  $\mathcal{D}$ 
22:          $y^j \leftarrow r_i^j + \gamma \mathbf{Q}^\mu(\mathbf{x}'^j, a'_1, \dots, a'_n)$ , where  $a'_k = \bar{\mu}_k(\sigma_k^{j'})$ 
23:         Update critic by minimizing the loss  $\ell(w_i) = \frac{1}{|\mathcal{B}|} \sum_j (y^j - \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j))^2$ 
using  $B$ 
24:         Update actor policy using policy gradient formula  $B$ 
25:          $\nabla_{\theta_i} J \approx \frac{1}{|\mathcal{B}|} \sum_j \nabla_{\theta_i} \mu_i(\sigma_i^j) \nabla_{a_i} \mathbf{Q}_i^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j)$ , where  $a_i = \mu_i(\sigma_i^j)$ 
26:       end for
27:       for all agents  $i \in [n]$  do
28:          $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$  (update target networks)
29:          $\bar{w}_i \leftarrow \tau w_i + (1 - \tau) \bar{w}_i$ 
30:       end for
31:     end if
32:      $step \leftarrow step + 1$ 
33:   until environment terminates
34:   NESTEDLOOKAHEAD( $n, e, \phi, \mathcal{L}$ )
35: end for
36: Output:  $\theta, w$ 

```

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Algorithm 8 Pseudocode for LA–MATD3: MATD3 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  $\theta \equiv \{\theta_i\}_{i=1}^n$ , both critic networks,  $\{Q_{i,1}, Q_{i,2}\}_{i=1}^n$  with initial weights
    $w \equiv \{w_{i,1}, w_{i,2}\}_{i=1}^n$ , learning rates  $\eta_\theta, \eta_w$ , base optimizer  $B$  (e.g., Adam), discount factor  $\gamma$ ,
   lookahead hyperparameters  $\mathcal{L} \equiv (l, \{k^{(j)}\}_{j=1}^l, \alpha_\theta, \alpha_w)$ , soft update parameter  $\tau$ , policy update
   frequency  $p$ .
2: Initialize:
3:   Replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
4:   LA parameters:  $\phi \leftarrow \{\theta\}_{\times l}, \{w\}_{\times l}$  (store snapshots for LA)
5: for all episode  $e \in 1, \dots, t$  do
6:    $x \leftarrow \text{Sample}(\mathcal{E})$  (sample from environment  $\mathcal{E}$ )
7:    $step \leftarrow 1$ 
8:   repeat
9:     if  $e \leq t_{\text{rand}}$  then
10:      for each agent  $i, a_i \sim \mathcal{A}_i$  (sample actions randomly)
11:     else
12:      for each agent  $i$ , select action  $a_i$  using current policy and exploration
13:     end if
14:      (apply actions and record results)
15:     Execute actions  $\mathbf{a} = (a_1, \dots, a_n)$ , observe rewards  $\mathbf{r}$  and new state  $\mathbf{x}'$ 
16:     replay buffer  $\mathcal{D} \leftarrow (\mathbf{x}, \mathbf{a}, \mathbf{r}, \mathbf{x}')$ 
17:      $x \leftarrow \mathbf{x}'$ 
18:      (apply learning step if applicable)
19:   if  $step \% t_{\text{learn}} = 0$  then
20:     for all agent  $i \in [n]$  do
21:       sample batch  $\{(\mathbf{x}^j, \mathbf{a}^j, \mathbf{r}^j, \mathbf{x}'^j)\}_{j=1}^{|\mathcal{B}|}$  from  $\mathcal{D}$ 
22:        $y^j \leftarrow r_i^j + \gamma \min_{m=1,2} Q_{i,m}^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j)$ , where  $a_k^j = \bar{\mu}_k(\sigma_k^j) + \epsilon$ 
23:       Update both critics,  $m = 1, 2$  by minimizing the loss (using optimizer  $B$ ):
24:         
$$\ell(w_{i,m}) = \frac{1}{|\mathcal{B}|} \sum_j \left( y^j - Q_{i,m}^\mu(\mathbf{x}^j, a_1^j, \dots, a_n^j) \right)^2$$

25:       if  $step \% p = 0$  then
26:         Update actor policy using policy gradient formula and optimizer  $B$ 
27:          $\nabla_{\theta_i} J \approx \frac{1}{|\mathcal{B}|} \sum_j \nabla_{\theta_i} \mu_i(\sigma_i^j) \nabla_{a_i} Q_{i,1}^\mu(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_n^j)$ , where  $a_i = \mu_i(\sigma_i^j)$ 
28:          $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$  (update target networks)
29:          $\bar{w}_{i,m} \leftarrow \tau w_{i,m} + (1 - \tau) \bar{w}_{i,m}$ 
30:       end if
31:     end for
32:   end if
33:    $step \leftarrow step + 1$ 
34: until environment terminates
35:   NESTEDLOOKAHEAD( $n, e, \phi, \mathcal{L}$ )
36: end for
Output:  $\theta, w$ 

```

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	1.5×10^6
learning step t_{learn}	100
t_{rand}	1024
Policy update ratio (MATD3) p	2
Noise std (MATD3)	0.2
Noise clip (MATD3)	0.5
Lookahead α	0.5

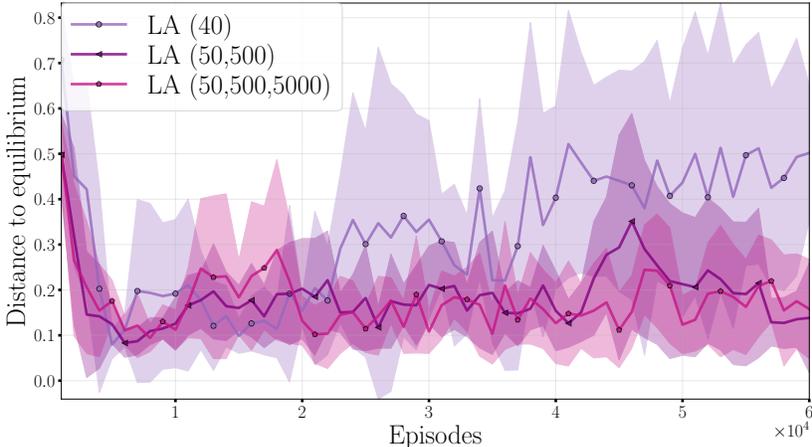


Figure 6: Comparison of LA-MATD3 with different levels in Rock-paper-scissors. x -axis: training episodes. y -axis: 5-seed average norm between the two players’ policies and equilibrium policy $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$.

E.2 COMPUTE RESOURCES

We ran the multi-agent experiments (RPS, MP, MPE) on Google Colab enterprise using an e2-standard-8 type machine with 100 GB Standard disk (pd-standard).

F ADDITIONAL EMPIRICAL RESULTS

F.1 MPE: PHYSICAL DECEPTION: COOPERATIVE-COMPETITIVE ENVIRONMENT

In *Physical deception*, we have p good agents, one adversary, and p landmarks, with one landmark designated as the *target*. Good agents aim to get close to the target landmark while misleading the adversary, which must infer the target’s location. Unlike Predator-Prey, this environment does not

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

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

Table 3: *Equilibrium reached?* Average adversary win rate for MPE: Physical deception on 100 test environments. The *win rate* is the fraction of times the adversary was closer to the target at the end of episode. *Closer to 0.5 is better*. Refer to Section 5.2.

involve direct competition for the adversary—its reward depends solely on its own policy. In our experiments, we set $p = 2$.

Table F.1 presents the mean and standard deviation of the adversary’s win rate, measuring how often it successfully reaches the target. In this setting, equilibrium is achieved when both teams win with equal probability across multiple instances. Given the *cooperative* nature of the game, the baseline performs relatively well, with EG-MADDPG achieving similar performance. However, both LA-MADDPG and LA-EG-MADDPG outperform their respective base optimizers (MADDPG and EG-MADDPG), demonstrating improved stability and effectiveness.

F.2 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 7 depicts the results when using different buffer options for the RPS game.

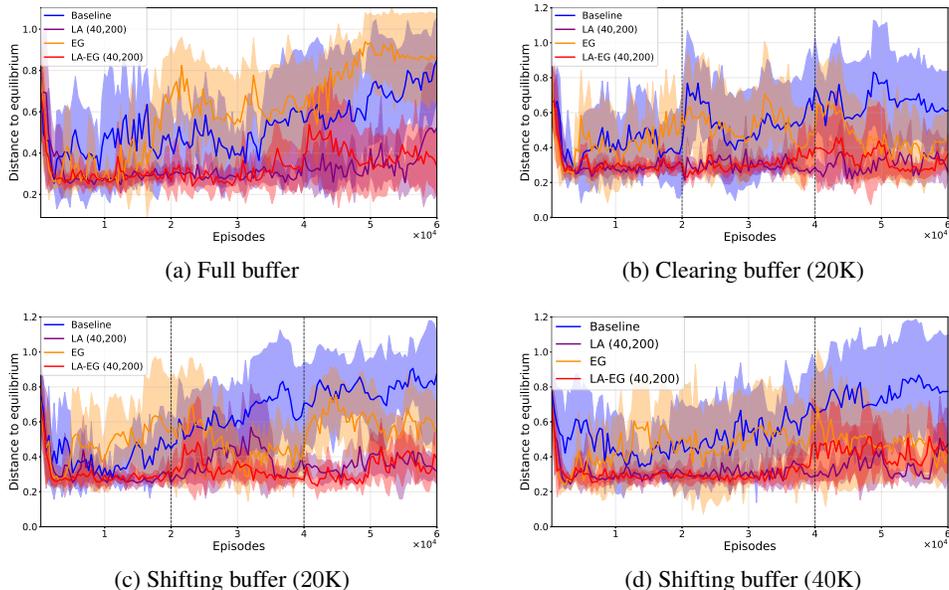
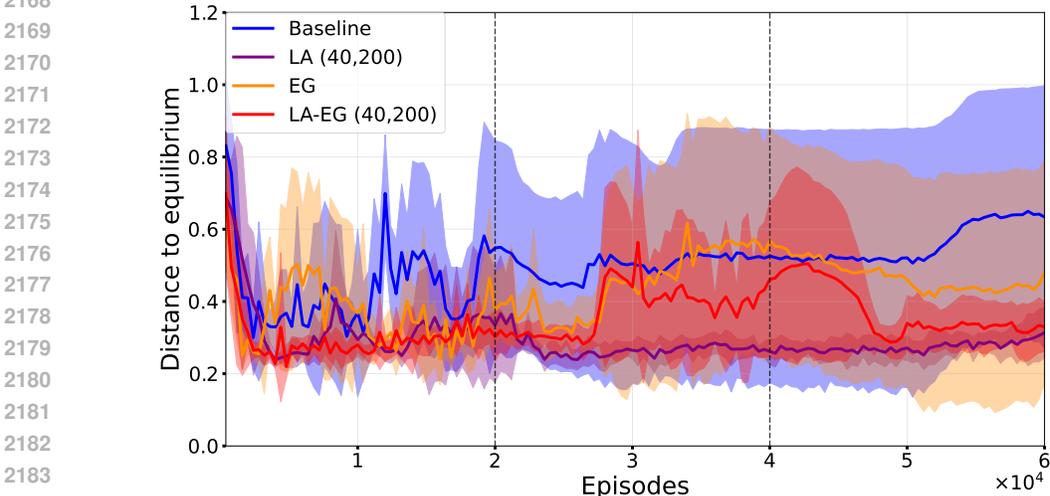


Figure 7: **Comparison of different buffer configurations (see Appendix F.2) and methods on Rock-paper-scissors game.** x -axis: training episodes. y -axis: 5-seed average norm between the two players’ policies and equilibrium policy $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^2$. The dotted line indicates the point at which the buffer begins to change, either through shifting or clearing.

2160 F.3 ROCK-PAPER-SCISSORS: SCHEDULED LEARNING RATE

2161 We experimented with gradually decreasing the learning rate (LR) during training to see if it would
 2162 aid convergence to the optimal policy in RPS. While this approach reduced noise in the results, it
 2163 also led to increased variance across all methods except for LA-MADDPG.

2164 Figure 8 depicts the average distance to the equilibrium policy over 5 different seeds for each methods,
 2165 using periodically decreased step sizes.



2166 Figure 8: Compares MADDPG with different LA-MADDPG configurations to the baseline MADDPG with
 2167 (Adam) in Rock-paper-scissors with a scheduled learning rate. x -axis: training episodes. y -axis: 5-seed
 2168 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
 2169 times when the learning rate was decreased by a factor of 10.

2190 F.4 MPE: PREDATOR-PREY FULL RESULTS

2191 We also evaluated the trained models of all methods on an instance of the environment that runs for 50
 2192 steps to compare learned policies. We present snapshots from it in Figure 10. Here, you can clearly
 2193 anticipate the difference between the policies from baseline and our optimization methods. As in the
 2194 baseline, only one agent will chase at the beginning of episode. Moreover, for the baseline (topmost
 2195 row), the agents move further away from the landmarks and the good agent, which is suboptimal.
 2196 This can be noticed from the decreasing agents’ size in the figures. While in ours, both adversary
 2197 agents engage in chasing the good agent until the end.

2199 F.5 MPE: PREDATOR-PREY AND PHYSICAL DECEPTION TRAINING FIGURES

2200 In figures 11a and 11b we include the rewards achieved during the training of GD-MADDPG and
 2201 LA-MADDPG resp. for MPE: Predator-prey. The figures show individual rewards for the agent
 2202 (prey) and one adversary (predator). Blue and green show the individual rewards received at each
 2203 episode while the orange and red lines are the respective running averages with window size of 100
 2204 of those rewards.

2205 Figures 12a and 12b demonstrate same results but for MPE: Physical deception. In this game, We
 2206 have two good agents, 'Agent 0 and 1' but since they are both receive same rewards, we only show
 2207 agent 0.

2210 F.6 ADDITIONAL METRICS

2211 For completeness, we include the NashConv and Exploitability results for Rock-paper-scissors (RPS).
 2212 NashConv and Exploitability are commonly used in game-theoretic papers to measure the incentive
 2213 of players to deviate quantized by the gain they get compared to best-response utility. In the optimal

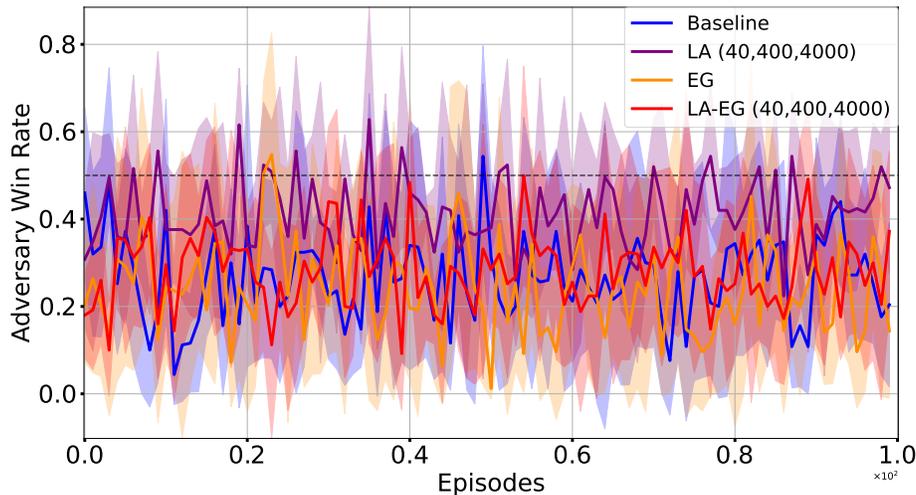


Figure 9: Comparison on the MPE–Predator-prey game between the *GD-MADDPG*, *LA-MADDPG*, *EG-MADDPG* and *LA-EG-MADDPG* optimization methods, denoted as *Baseline*, *LA*, *EG*, *LA-EG*, resp. x -axis: evaluation episodes. y -axis: mean adversaries win rate, averaged over 5 runs with different seeds.

case, both measures are equal to zero and closer to zero means less ability to deviate and less chance of being exploited.

In figure 13 we can see the results for using GD vs. LA in RPS, we can clearly see that LA has much lower values for the metrics compared to the baseline. Which confirms that it finds better policies.

F.7 VARIANCE OF OUTERMOST LOOKAHEAD LEVEL

We also ran an additional experiment with longer episodes, where we evaluated and recorded the distance to equilibrium only after the outermost lookahead step, depicted in Figure 14. In these runs, the variance across seeds is noticeably smaller and the trajectories are more stable (compared to Figure 3a, further supporting the claim that Lookahead stabilizes training. This also reinforces the design choice in our pseudocode and in (Chavdarova et al., 2021) to always use the iterate obtained immediately after the outermost lookahead level (Algorithm 1 & 3).

F.8 ABLATION: NUMBER OF SEEDS

In Figure 15, we compare the mean and variance of the distance to equilibrium in Rock–paper–scissors for *LA-MADDPG* trained with 5 versus 10 runs. The curves closely overlap, indicating that the results are consistent and robust to the number of seeds.

F.9 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 16 (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

2268
 2269
 2270
 2271
 2272
 2273
 2274
 2275
 2276
 2277
 2278
 2279
 2280
 2281
 2282
 2283
 2284
 2285
 2286
 2287
 2288
 2289
 2290
 2291
 2292
 2293
 2294
 2295
 2296
 2297
 2298
 2299
 2300
 2301
 2302
 2303
 2304
 2305
 2306
 2307
 2308
 2309
 2310
 2311
 2312
 2313
 2314
 2315
 2316
 2317
 2318
 2319
 2320
 2321

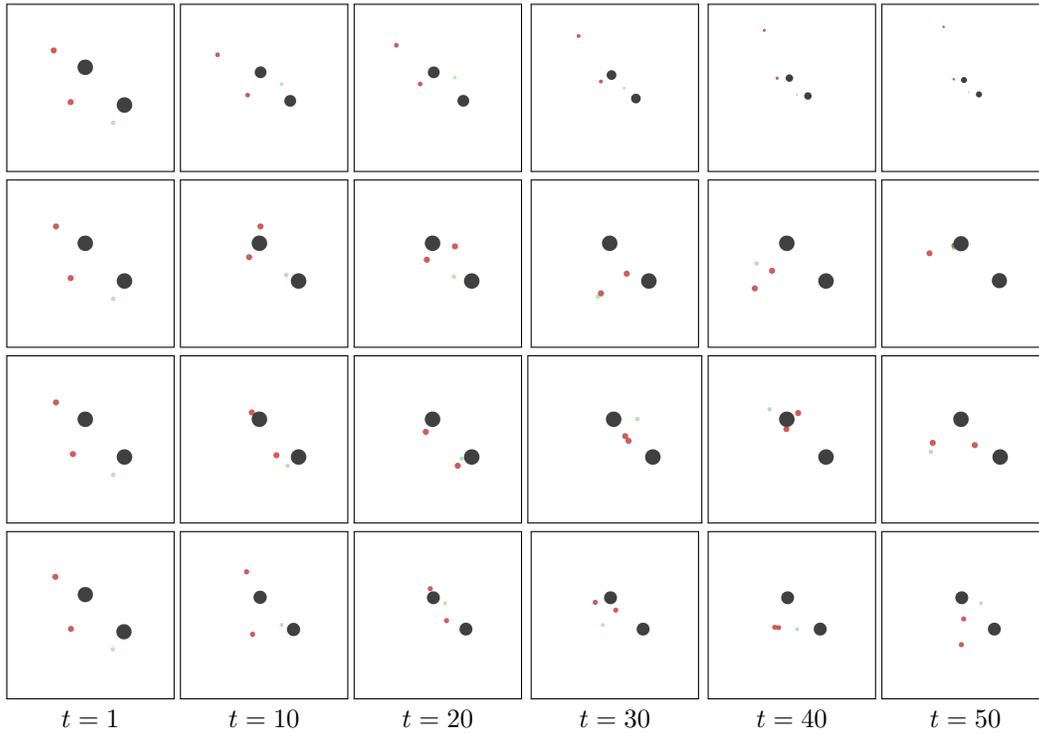


Figure 10: **Agents’ trajectories of fully trained models with all considered optimization methods on the same environment seed of MPE: Predator-prey.** Snapshots show the progress of agents as time progresses in a 50 steps long environment. Each row contains snapshots of one method, from top to bottom: *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.

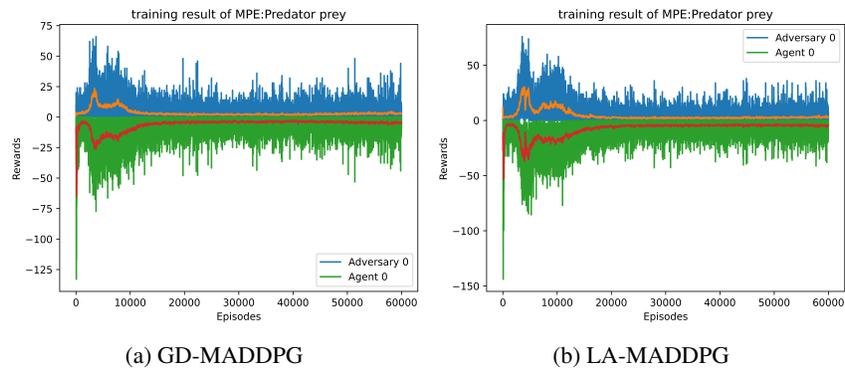
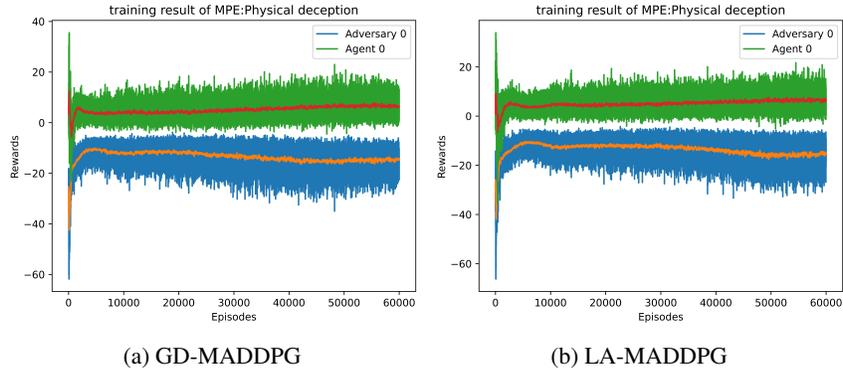


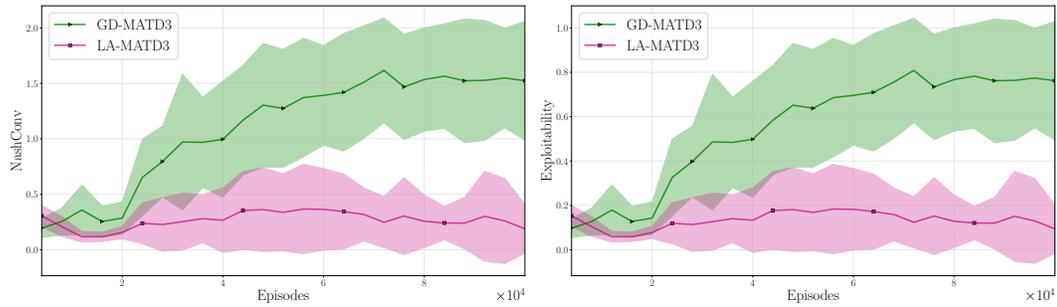
Figure 11: **The figure shows the learning curves during training of GD-MADDPG and LA-MADDPG for MPE: Predator-Prey.** *x*-axis: training episodes. *y*-axis: agents’ rewards and their moving average with a window size of 100, calculated over 5-seeds over 5 seeds.

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334



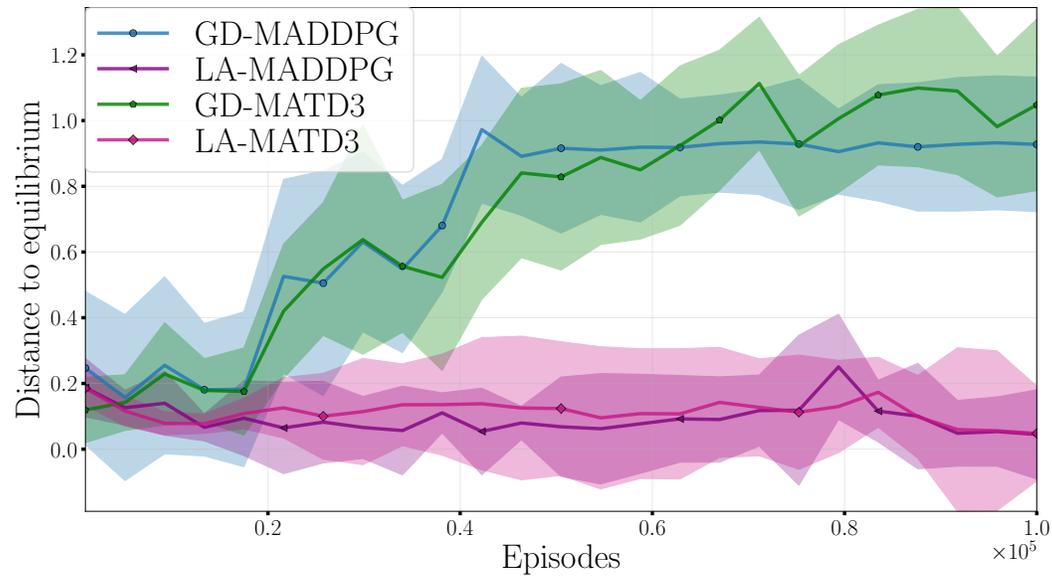
2335 **Figure 12: The figure shows the learning curves during training of GD-MADDPG and LA-MADDPG for MPE: Physical deception.** *x*-axis: training episodes. *y*-axis: agents’ rewards and their moving average with a window size of 100, calculated over 5-seeds over 5 seeds.

2339
2340
2341
2342
2343
2344
2345
2346
2347
2348



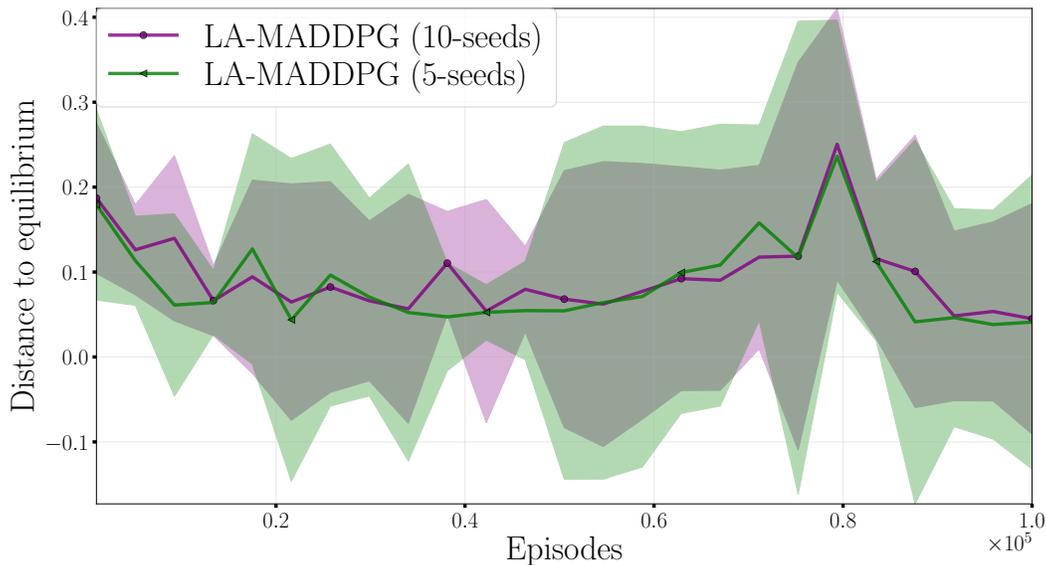
2349 **Figure 13: NashConv (left) and exploitability (right) during training episodes of MATD3 on Rock-Paper-Scissors for GD-MATD3 and LA-MATD3.** Curves show the mean over 10 random seeds, and shaded regions indicate one standard deviation across seeds.

2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372



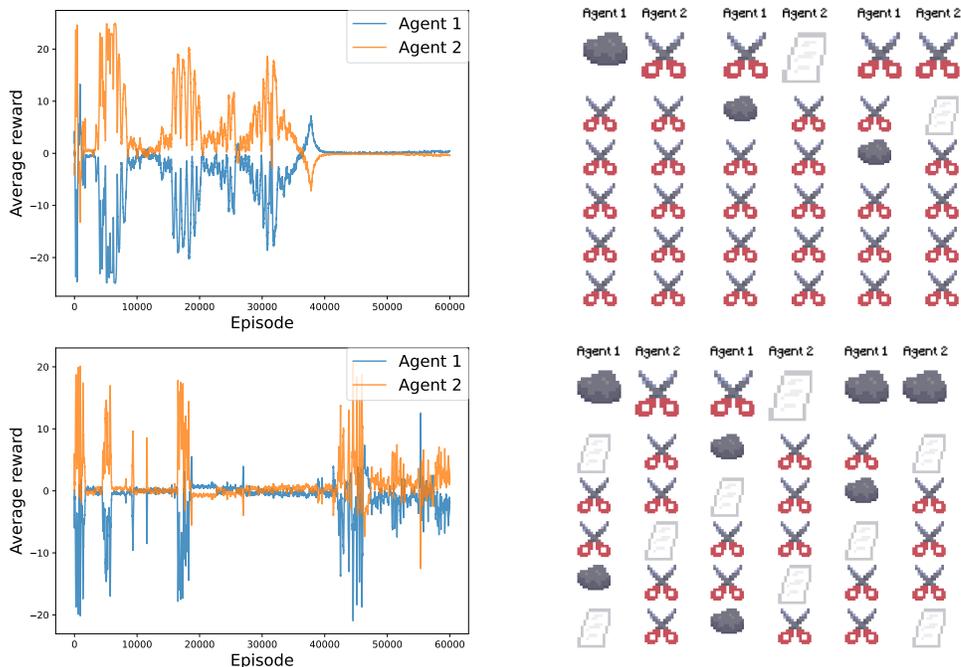
2373 **Figure 14: Comparison of GD-(MADDPG/MATD3) and LA-(MADDPG/MATD3) on Rock-paper-scissors.** The *x*-axis shows training episodes, and the *y*-axis shows the distance between the agents’ policies and the equilibrium policy. Results are measured after the outermost lookahead level. Curves are averaged over 10 random seeds, and shaded regions indicate ± 1 standard deviation across these 10 seeds.

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395



2396 **Figure 15: Effect of the number of random seeds on LA-MADDPG in Rock–Paper–Scissors.** The x -axis
2397 shows training episodes, and the y -axis shows the distance between the learned joint policy and the Nash
2398 equilibrium. The two curves correspond to LA-MADDPG trained with 10 and 5 random seeds, respectively,
2399 illustrating the consistency of the method under different levels of averaging across seeds.

2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423



2424 **Figure 16: Detailed Figure 5.** Saturating rewards (left) versus actions of the learned policies at the end (right)
2425 in the Rock–paper–scissors game. **Top row: GD-MADDPG; bottom row: LA-MADDPG.** In the left column,
2426 blue and orange show the running average of rewards through a window of 100 episodes. In the right column,
2427 we depict actions from the respective learned policies evaluation after training is completed, where each row
2428 represents what actions players have chosen in one step of the episode. Saturating rewards do not imply good
2429 performance, as evidenced by the top row; refer to Section 5.2 for discussion.

2430
 2431
 2432
 2433
 2434
 2435
 2436
 2437
 2438
 2439
 2440
 2441
 2442
 2443
 2444
 2445
 2446
 2447
 2448
 2449
 2450
 2451
 2452
 2453
 2454
 2455
 2456
 2457
 2458
 2459
 2460
 2461
 2462
 2463
 2464
 2465
 2466
 2467
 2468
 2469
 2470
 2471
 2472
 2473
 2474
 2475
 2476
 2477
 2478
 2479
 2480
 2481
 2482
 2483

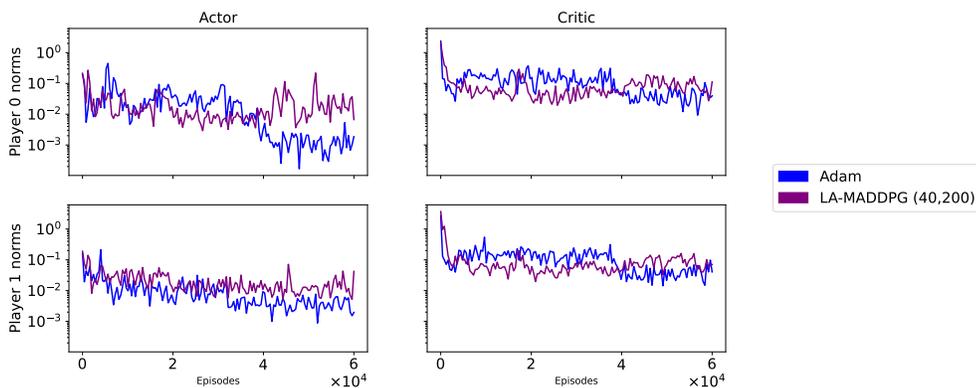


Figure 17: Gradient norms across training in the *Rock-paper-scissors* game.

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 3a.

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 17, 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.