DRIMA: DIFFERENTIAL REWARD INTERACTION FOR COOPERATIVE MULTI-AGENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-agent reinforcement learning (MARL) owning to its potent capabilities in complex systems has gained remarkable research attention nowadays, in which collaborative decision-making and control for multi-agent systems is one of the key research focuses. A prevalent learning framework is centralized training with decentralized execution (CTDE), in which the decentralized execution realizes strategy flexibility, and the use of centralized training ensures stationarity and goal consistency while becoming incapable when facing scalability and complicated situations. To address this issue, we follow the concept of distributed training with decentralized execution (DTDE). Decentralization is naturally accompanied by the game during the learning process, which has not been entirely studied in related work, resulting in the constrained strategy combination of MARL. In this paper, we devise a novel approach of differential reward interaction (DRI) with conflict-triggered for the distributed evaluation that enables overall goal consistency through highly efficient local information exchange. With this collaboration of learning, the DRI-based MARL can eliminate the notorious issue of converging to saddle equilibriums of stochastic games. Meanwhile, it possesses provable convergence and is well compatible for general value-based and policy-based algorithms. Experiments in several benchmark scenarios demonstrate that DRIMA realizes collaborative strategy learning with enhanced global goal-achieving.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

1 INTRODUCTION

034 Multi-agent reinforcement learning (MARL) has been extensively studied in recent years owing to its capability of assisting multi-agent systems (MASs) in complex real-world applications, in-035 cluding multi-robot systems (Mou et al., 2021; Marchesini & Farinelli, 2022), autonomous driving (Palanisamy, 2020; Vinitsky et al., 2022), traffic signal control (Ault & Sharon, 2021; Du et al., 037 2023), and energy network regulation (Wang et al., 2021; Ma et al., 2024). MARL has been facing critical scalability issues since the size of the state and action space grows exponentially with the number of agents. The CTDE framework consists of centralized evaluation to collect information, 040 such as observations, actions, or rewards of all agents, and decentralized execution that allows each 041 agent to determine its own behavior according to individual observation. It largely mitigates the 042 curse of dimensionality on the action space and has been adopted by many popular algorithms such 043 as MADDPG (Lowe et al., 2017), VDN (Sunehag et al., 2017), COMA (Foerster et al., 2018), MAP-044 PO (Yu et al., 2022), etc. Such a centralized evaluator, however, still has limitations in the scalability that assembles information from a mass of agents during training, and in complex collaborative tasks where the overall performance evaluator is insufficient to distinguish the global and local situation. 046 A class of DTDE research on the distributed communication-based MARL then emerged, which 047 mainly focuses on achieving global-level objectives with local-level collaborative learning. 048

Distributed MARL mainly differs in the way of information transmission and fusion. Qu and Lin (Qu et al., 2020; Lin et al., 2021; Qu et al., 2022) inventively proposed the exponential decay condition for general networked MARL that quantitatively bound the interaction strength between agents, under which the evaluation for each agent is scaled by the neighboring state-action space size instead of the global one, which realizes scalability. Mean Field Reinforcement Learning (MFRL) (Yang et al., 2018; Qiu et al., 2022; Yu, 2023) takes individual and average neighbor actions as

054 components of the respective state-action value function. The Mean-field theory, that each agent 055 in a gridded world is influenced only by the mean effect of its neighbors, ensures respective value 056 function is a good approximation to the joint one. On the other hand, the consensus-based method 057 (Wai et al., 2018; Zhang et al., 2018; Qu et al., 2019; Chen et al., 2022) has been a popular topic 058 in distributed MARL with parameterized functions. Zhang et al. (2018) pioneeringly introduces parameter consensus to the process of multi-agent learning. Based on the theory of distributed optimization and stochastic approximation, the achievement of global consistency by distributed value 060 function has been proved. The convergence of this consensus approach is guaranteed for linear 061 function approximation while is unreliable for neural networks (i.e. the nonlinear parameterized 062 functions). There is also a class of works (Jiang et al., 2018; Blumenkamp & Prorok, 2021; Nayak 063 et al., 2023) based on Graph Neural Networks (GNNs), designing sophisticated network structures 064 to model multi-agent interactions and achieve collaboration in automatic learning ways. 065

An inevitable issue facing the decentralized decision-making MASs is the game between agents, 066 which has not been fully studied in previous DTDE works and potentially leads the strategy com-067 bination to poor situations. Simultaneous independent decision-making and conflicting payoffs are 068 two essential causes we identified. The former usually ensures strategy diversity, thus a promising 069 direction is to work on the reconstruction of payoffs (or rewards) to solve the dilemma. Chu et al. (2020) proposed a reward reshaping method naturally based on the physical distance between neigh-071 boring agents for the network systems. Hostallero et al. (2020) designed a peer evaluation signal 072 calculated by respective temporal difference (TD) that reflects the value of joint action, and agents 073 receive it from peers to reshape their rewards. This method possesses more generality while less 074 convergence property. Other learn-to-reshape methods were based on deep NNs (Yang et al., 2020; 075 Yi et al., 2022), with special structures designed to automatically learn the reward-reshaping terms. The existing reshaping methods alleviate game dilemmas and promote cooperations to a certain 076 extent, while the absence of identification and effective approaches to the critical causes results in 077 limited generality and capability.

079 In this paper, a novel approach named DRIMA is proposed for distributed MARL that solves the 080 dilemma of games for DTDE and achieves enhanced strategy collaborations. A conflict-triggered 081 differential reward interaction (DRI) method is designed that allow agents to reshape personal rewards according to the opposite signal from neighbors. The conflict sign of DRs essentially reveals 083 strategies contradiction, and the collaboration when it occurs enables the policy combinations to avoid saddle equilibriums in stochastic games. The interaction of DRs, which are scalars, compared 084 with the transmission of NN parameters and the training of extra NN structures, realizes cooperative 085 learning in a highly efficient way. DRIMA is naturally compatible with general multi-agent (deep) reinforcement learning algorithms and possesses provable convergence. The performance of our 087 approach was examined on several benchmark scenarios, including the matrix games, Multi-Agent 088 Particle Environment (MPE) (Lowe et al., 2017), and Star-Craft II (Samvelyan et al., 2019).

089 090 091

2 PRELIMINARIES

092 093 094

2.1 MULTI-AGENT MARKOV GAME

095 Multiple agents making sequential decisions in the same environment can be modeled as a N-agent 096 Markov game (or N-player stochastic game) Γ , defined as a tuple $\langle S, A^1, \dots, A^N, r^1, \dots, r^N, P \rangle$. \mathcal{S} is the state space, \mathcal{A}^i is the action space of agent $i \ (i = 1, \cdots, N), r^i : \mathcal{S} \times \mathcal{A}^1 \times \cdots \times \mathcal{A}^N \to \mathbb{R}$ 098 is the reward function for agent *i*, and $P: \mathcal{S} \times \mathcal{A}^1 \times \cdots \times \mathcal{A}^N \to \Delta(\mathcal{S})$ is the transition probability map where $\Delta(S)$ is the set of probability distributions over state space S. Agent policy is a mapping 100 $\pi^i: \mathcal{S} \times \mathcal{A}^i \to \Delta(\mathcal{A}^i)$, where $\Delta(\mathcal{A}^i)$ denotes the space of probability distributions over the agent 101 *i*'s actions. The joint policy of all agents is denoted as $\pi = (\pi^1, \dots, \pi^N)$, and the probability value 102 of it is $\pi(a|s) = \prod_{i=1}^{N} \pi^i(a^i|s)$. At each time step t, given current state s_t , each agent i choose 103 action a_t^i according to own policy $\pi^i(\cdot|s_t)$. All agents form a joint action $a_t = (a_t^1, \cdots, a_t^N)$ 104 and each agent receives a reward $r^i(s_t, a_t)$. The state then transits to the next s_{t+1} based on the 105 transition probability P of the environment. Multiple agents in a Markov game will rationally try to 106 find strategies that maximize their own expected returns. This paper aims to design a collaborative 107 approach for distributed MARL to maximize the global return of MASs.

108 2.2 NETWORKED MARL FORMULATION

110 Networking is a common framework for studying MASs in their entirety. Agents in a networked 111 system are denoted as a set $\mathcal{N} = \{1, \dots, N\}$. The system can be modeled as an undirected graph 112 $\mathcal{G}(\mathcal{N}, \mathcal{E})$, where each agent *i* serves as vertex *i* and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ is the set of all edges. The edge 113 $e_{ij} = (i, j) \in \mathcal{E}$ indicates the connection relationship of two agents $i, j \in \mathcal{N}$. Two agents associated 114 with an edge are neighbors, agent *i* and its all neighbors in the graph together is denoted as a set 115 \mathcal{N}^i . Note that the connection between agents could be explicit physical connection, communication 116 relationship, implicit mutual influence, etc.

The global objective of the MAS is defined as maximizing the overall average reward per time step:

$$\underset{\boldsymbol{\pi}}{\text{maximize }} J = \lim_{T} \frac{1}{T} \mathop{\mathbb{E}}_{s \sim P, \boldsymbol{a} \sim \boldsymbol{\pi}} \left(\sum_{t=0}^{T-1} \frac{1}{N} \sum_{i \in \mathcal{N}} r^{i}(s, \boldsymbol{a}) \right) = \sum_{s \in \mathcal{S}} d_{\boldsymbol{\pi}}(s) \sum_{\boldsymbol{a} \in \mathcal{A}} \pi(\boldsymbol{a}|s) \cdot \bar{r}(s, \boldsymbol{a}), \quad (1)$$

127 128 129

130 131

137 138

117

118

where $\bar{r}(\cdot) = \frac{1}{N} \sum_{i \in \mathcal{N}} r^i(s, a)$ is the average function of all agents' rewards.

To solve this optimization problem in a decentralized way, agents are set to make decisions independently. Based on the continuing RL problem formulation that utilizes the *average reward* in Sutton & Barto (2018), and considering that the action value of each agent *i* reflects the expected value of its policy π^i affected by other policies, the differential action-value function is formulated as:

$$Q_{\pi}^{i}(s,a^{i}) = \sum_{t=0}^{\infty} \mathbb{E}_{s \sim P, a^{-i} \sim \pi^{-i}} [r^{i}(s_{t}, a_{t}) - \mu_{\pi}^{i} | s_{0} = s, a_{0} = a^{i}],$$

where $\mu_{\pi}^{i} = \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{s \sim P, \mathbf{a} \sim \pi} [r^{i}(s_{t}, \mathbf{a}_{t})] = \sum_{s \in S} d_{\pi}(s) \sum_{\mathbf{a} \in \mathcal{A}} \pi(\mathbf{a}|s) \cdot r^{i}(s, \mathbf{a})$ denotes the average reward per time step. $r^{i}(s_{t}, \mathbf{a}_{t}) - \mu_{\pi}^{i}$ denotes the *differential reward* term of agent *i*. $d_{\pi}(s)$

average reward per time step. $r^{i}(s_{t}, a_{t}) - \mu_{\pi}^{2}$ denotes the *aliferential reward* term of agent *i*. $a_{\pi}(s)$ represents a stationary distribution over *S* induced by joint policy π , which is assumed to exist in general problem formulations to ensure convergence theoretically (Zhang et al., 2018; Qu et al., 2020). -i denotes agents excluding *i*. According to Bellman equation, the differential state-value function of *i* is

$$v^{i}_{\boldsymbol{\pi}}(s) = \sum_{a^{i} \in \mathcal{A}^{i}} \pi^{i}(s, a^{i}) \cdot Q^{i}_{\boldsymbol{\pi}}(s, a^{i})$$

Note that the average reward formulation can be generalized to the discount-reward one, that the differential reward term in above definitions is discounted by a time-decaying factor and becomes $\gamma^t \cdot (r^i(s_{t+1}, a_{t+1}) - \mu_{\pi}^i), \gamma \in [0, 1)$. In such setting, a discounted stationary distribution is assumed to exist (Nachum et al., 2019).

143 This paper focuses on general-sum games, which allows the agents' rewards to be arbitrarily relat-144 ed. Zero-sum games and coordination games are special cases where agents' rewards are always 145 negatively and positively related. The baseline solution concept for general-sum games is Nash Equilibrium (NE), which denotes a joint strategy where each agent's response is the best to the others'. 146 We focus on the special types of NE points: global optima, and saddles, which were well-studied in 147 Hu & Wellman (2003). The detailed definitions refer to Appendix A.1. In the following, a collabo-148 rative evaluation method is proposed to avoid saddle equilibriums in Markov games for multi-agent 149 reinforcement learning. 150

151 152

153

3 DIFFERENTIAL REWARD INTERACTION-BASED MARL

154 We identify that the global performance of distributed MARL is critically restricted by two factors. 155 One is the simultaneous and independent decision-making, which ensures the diversity of strategy 156 and the flexibility of decision, should be practically maintained. The other is the distinct payoffs 157 individuals receive after joint action executions, which lead agents to make selfish evaluations and 158 decisions that might be detrimental to the overall. We propose a method of Differential Reward 159 Interaction (DRI) with conflict-triggered, whose key idea is based on the exchange of the average reward, a scalar reflecting a policy's overall performance, to identify and reshape conflicting pay-160 offs between agents, that efficiently trade off joint goal consistency and individual diversity during 161 learning and achieve the enhanced strategy combination.

162 3.1 CONFLICT-TRIGGERED DIFFERENTIAL REWARD INTERACTION METHOD

For every agent i $(i=1, \dots, N)$ with original reward $r^i(s, a)$, abbreviated as r^i in the following, the conflict-triggered interaction rule for differential rewards is designed as follows

$$\begin{pmatrix} r^{i} - \bar{\mu}^{N^{i}}, \operatorname{sgn}\left(r^{i} - \bar{\mu}^{N^{i}}\right) = \operatorname{sgn}\left(\bar{r}^{N^{-i}} - \bar{\mu}^{N^{i}}\right) \\ \bar{r}^{N^{i}} - \bar{\mu}^{N^{i}}, \operatorname{sgn}\left(r^{i} - \bar{\mu}^{N^{i}}\right) \neq \operatorname{sgn}\left(\bar{r}^{N^{-i}} - \bar{\mu}^{N^{i}}\right),$$

$$(2)$$

where sgn (\cdot) denotes Signum function

$$\operatorname{sgn}(x) = \begin{cases} -1 : x < 0\\ 0 : x = 0\\ 1 : x > 0 \end{cases}$$

 $\bar{r}^{N^{i}}$ and $\bar{r}^{N^{-i}}$ are the averaged rewards of *i*'s neighbors with and without *i* itself, defined as

$$\bar{r}^{N^{i}} = \frac{1}{N^{i}} \sum_{j \in \mathcal{N}^{i}} r^{j} \text{ and } \bar{r}^{N^{-i}} = \frac{1}{N^{i} - 1} \sum_{j \in \mathcal{N}^{-i}} r^{j}.$$

 $\bar{\mu}^{N^i}$ indicates averaging the average reward of *i*'s neighbors, defined as

$$\bar{\mu}^{N^i} = \frac{1}{N^i} \sum_{j \in \mathcal{N}^i} \mu^j_{\pi},$$

where μ_{π}^{j} is the time-average reward of agent j, i.e. , $\mu_{\pi}^{j} = \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{s \sim P, \mathbf{a} \sim \pi} [r^{j}(s_{t}, \mathbf{a}_{t})],$ $j \in \mathcal{N}^{i}$, and it can be incrementally estimated by respective reward samples through $\mu_{t}^{j} \leftarrow \mu_{t-1}^{j} + \alpha_{t} \cdot (r^{j}(s_{t}, \mathbf{a}_{t}) - \mu_{t-1}^{j})$ where α_{t} is an appropriate step-size.

188 In Eqn. (2), the opposite sign of the differential reward terms in the Signum function indicates 189 the "conflict" between agents about their value assessments of the group's behaviors. To get an 190 intuition of this idea, consider first a *single-agent* system. At each iteration step t ($t=1, 2, \cdots$), 191 the agent obtains a differential reward $r(s_t, a_t) - \mu_t$ according to the instant reward sample. Since 192 μ_t estimates the expected average reward and reflects the overall assessment of the current policy, 193 the immediate differential reward in fact indicates a (positive or negative) contribution value to 194 the current policy by taking action a_t at state s_t . Relating to the *multi-agent* system, based on 195 Eqn. (2), each agent *i* is allowed to constantly communicate average rewards with neighbors to 196 obtain $\bar{\mu}_t^{N^i}$, which reflects the performance on average of current neighbor-joint policy $\pi_t^{N^i}$. And the opposite sign between $r^i(s_t, \boldsymbol{a}_t) - \bar{\mu}_t^{N^i}$ and $\bar{r}^{N^{-i}}(s_t, \boldsymbol{a}_t) - \bar{\mu}_t^{N^i}$ indicates the conflict contribution of action decision relating to the individual *i* and its other neighbors N^{-i} to the current neighbor-197 199 joint policy. When conflict signs occur, the agent's personal reward $r^i(s_t, a_t)$ be reshaped as the 200 neighbor-averaged reward $\bar{r}^{N^{i}}(s_{t}, a_{t})$ to keep joint goal consistency. Based on differential reward 201 interaction and reward consensus at critical steps, the conflict-triggered DRI method enables agents 202 to maintain collaboration during learning, balance joint goal consistency and individual behavior 203 diversity, and eliminate joint strategy to saddle equilibriums in a highly efficient way. We present the detailed study in the following. 204

205 206

207

208

181 182 183

184 185

186

187

3.2 Two-Player Matrix Game

Without loss of generality, a Prisoner's dilemma game is provided with payoffs (i.e. rewards) in Fig. 1. It originally has a unique NE which is a saddle point valued at (2, 2), and the global optimum, valued at (10, 10), is not an NE. For the problem formulation of MARL, this game can be seen as a continuing task that is repeatedly played by two agents. It is a non-state transition task, i.e., it can be regarded as having only one state s_0 , and no matter what joint action performed at each step will lead to the same state. The action spaces of Red and Blue agents are \mathcal{A}^{re} , $\mathcal{A}^{bl} = \{a_C, a_D\}$, where a_C means cooperate and a_D is defect. The joint action is denoted as $\mathbf{a} = (a^{re}, a^{bl}), \mathbf{a} \in \mathcal{A}$ where $\mathcal{A} = \mathcal{A}^{re} \times \mathcal{A}^{bl}$. Following the networked formulation, two agents are set to be neighbors and allow

224

233

234 235 236

241 242 243

253 254

Blue Red	Cooperate	Defect
Cooperate	10	0
Defect	18 0	2 2

Figure 1: Prisoner's dilemma payoff matrix.

communication. According to DRI method, their action and state value functions are as

$$\tilde{Q}_{\pi}^{i}\left(s_{0},a^{i}\right) = \sum_{t=0}^{\infty} \mathbb{E}_{a^{-i} \sim \pi^{-i}} \left[\tilde{r}^{i}(s_{0}, \boldsymbol{a}_{t}) - \bar{\mu}_{\pi} | s_{0}, a_{0} = a^{i}\right] \text{ and } \tilde{v}_{\pi}^{i}\left(s_{0}\right) = \sum_{a^{i} \in \mathcal{A}^{i}} \pi^{i}\left(a^{i} | s_{0}\right) \cdot \tilde{Q}_{\pi}^{i}\left(s_{0}, a^{i}\right),$$

where $\tilde{r}_t^i, i \in \{\text{re, bl}\}\$ is the reshaped reward following the interaction rule Eqn. (2), and $\bar{\mu}_{\pi} = (\mu_{\pi}^{\text{re}} + \mu_{\pi}^{\text{bl}})/2$. The objective of this two-agent system is

$$\underset{\boldsymbol{\pi}_{*}}{\text{maximize } J_{*}} = \sum_{\boldsymbol{a} \in \mathcal{A}} \pi_{*}(\boldsymbol{a}) \cdot \frac{1}{2} \left(r^{\text{re}}(\boldsymbol{a}) + r^{\text{bl}}(\boldsymbol{a}) \right) = \frac{1}{2} \left(\mu_{\boldsymbol{\pi}_{*}}^{\text{re}} + \mu_{\boldsymbol{\pi}_{*}}^{\text{bl}} \right) = \bar{\mu}_{\boldsymbol{\pi}_{*}}$$

where $\pi_* = (\pi^{re}_*, \pi^{bl}_*)$ is Nash equilibrium.

239 We take the policy gradient Actor-Critic algorithm here to analyze the learning process to the equilibrium. Set two decision-making policies defined by two parameters p and q, respectively, as

$$\begin{cases} \pi^{\rm re}(a_C) = 1 - p \\ \pi^{\rm re}(a_D) = p \end{cases}, \begin{cases} \pi^{\rm bl}(a_C) = 1 - q \\ \pi^{\rm bl}(a_D) = q \end{cases}, \ p, q \in (0, 1). \end{cases}$$
(3)

 $\pi = (\pi^{re}, \pi^{bl})$ denotes the joint strategy. For each agent *i*, Actor-Critic iteratively advances by a 244 faster Critic step to estimate value function $\tilde{v}^i_{\boldsymbol{\pi}_t}$, alternating with a slower Actor step for policy π^i_t 245 246 improvement based on current value function. At the t-th iteration, regarding the faster speed of the evaluation step, the policy improvement step can be analyzed separately under the assumption that 247 $\tilde{v}^i_{\pi_t}$ and $\bar{\mu}_{\pi_t}$ of the current policy π_t are already the expected values. To the value function defined 248 by the differential return, $\mathbb{E}_{s\in\mathcal{S}}[v_{\pi}(s)] = 0$ holds (Sutton & Barto, 2018). We have $\tilde{v}^{i}_{\pi_{t}}(s_{0}) = 0$ in 249 this case, thus the one-step TD term $\tilde{r}^i(a_t) - \bar{\mu}_{\pi_t} + \tilde{v}^i_{\pi_t}(s_{t+1}) - \tilde{v}^i_{\pi_t}(s_t)$ degenerates as $\tilde{r}^i(a_t) - \bar{\mu}_{\pi_t}$. 250 The policy parameter update for each agent at t is then as follows 251

$$p_{t+1} = p_t + \beta_t \cdot (\tilde{r}^{\text{re}}(\boldsymbol{a}_t) - \bar{\mu}_{\boldsymbol{\pi}_t}) \cdot \nabla_p \log \pi_{p_t}^{\text{re}}(\boldsymbol{a}_t^{\text{re}}),$$

$$q_{t+1} = q_t + \beta_t \cdot (\tilde{r}^{\text{bl}}(\boldsymbol{a}_t) - \bar{\mu}_{\boldsymbol{\pi}_t}) \cdot \nabla_q \log \pi_{q_t}^{\text{bl}}(\boldsymbol{a}_t^{\text{bl}}),$$
(4)

where β_t is an appropriate step-size of the Actor. Based on the policy definition of Eqn. (3), and to facilitate analysis, the gradient terms in Eqn. (4) can always be transferred to be positive no matter what action is sampled. With this view, the differential reward terms are directly affects the gradient ascent direction and size. We draw the geometry of this optimization problem in Fig. 2 to illustrate the conflict-triggered mechanism of DRI.

As shown in Fig. 2(a), the x-axis and y-axis represent the domains of p and q, respectively, and the zaxis represents the value range of μ_{π} . In feasible region, μ_{π}^{re} is displayed as the Hot-colored surface, and μ_{π}^{bl} is as Cold-colored surface. $\bar{\mu}_{\pi}$ denotes the average of μ_{π}^{re} and μ_{π}^{bl} shown as the Rainbowcolored one between them. One can see that independent learning will lead the Red and Blue agent approach to p=1 and q=1 (valued at (2, 2)) along the x-axis and y-axis respectively, since it is the highest point in the direction of their respective gradient ascent. But the global optimum, valued at (10, 10) when p=0, q=0, is the highest point of the joint policy solution on Rainbow surface $\bar{\mu}_{\pi}$.

Conflict-triggered differential reward interaction allows each agent to reconstruct the joint surface at the necessary steps and approach the joint optimum. To see this, we arbitrarily assume that the current policy is at $p_t = 0.4$, $q_t = 0.8$, then $\bar{\mu}_{\pi_t} = 6.88$ is depicted as a green triangle on the Rainbow surface. Agents take actions with the probabilities based on respective current policies. And the



Figure 2: Solution space visualization of the Prisoner's Dilemma.

formed joint action $a_t = (a_t^{\text{re}}, a_t^{\text{bl}})$ produces a related joint reward sample $(r^{\text{re}}(a_t), r^{\text{bl}}(a_t))$, i.e. one of (10, 10), (0, 18), (18, 0), and (2, 2) according to the payoff matrix. The conflict-triggered mechanism works when the two differential reward values calculated by the associated reward samples possess different signs, i.e. $\operatorname{sign}(r_t^{\text{re}} - \bar{\mu}_{\pi_t}) \neq \operatorname{sign}(r_t^{\text{bl}} - \bar{\mu}_{\pi_t})$. Here, samples (0, 18) and (18, 0) trigger interactions, while samples (10, 10) and (2, 2) do not.

288 Taking (0, 18) for an example, it and its average value of 9 are shown in red, blue, and yellow dots in 289 Fig. 2(b). The corresponding differential values for 0, 18, 9 are -6.88, 11.12, 2.12, that indicates the 290 sample's contribution to joint policy, from the view of the Red, Blue, and (hypothetic) global agent. 291 Three colored dashed lines represent the respective sight line of each dot observing the current $\bar{\mu}_{\pi_t}$ 292 (drawn as the green triangle). One can see that, on the Rainbow joint surface, from the global 293 view (i.e., the yellow dot), the gradient ascent direction projection along the x-axis and y-axis (that respectively indicate the policy parameters' feasible direction for agent Red and Blue) are the yellow and blue vectors. However, the ascent direction from the Red dot view on the Hot surface, which 295 is along the red dash line, projects to its feasible direction and then to the joint surface (drawn as 296 the red vector), opposite to the global view (the yellow vector). That is to say, the direction of Red 297 agent's strategy optimization, based on the current personal differential reward sampling, is contrary 298 to the global optimization direction. 299

300 Specifically, conflict arises due to opposite viewing angles when two points above and below the 301 horizontal of the $\bar{\mu}_{\pi}$,-point observe it respectively, causing goal inconsistency. The opposite sign of differential rewards between agents exactly represents such conflict. When it occurs, agents ex-302 change and average respective differential rewards through the method of Eqn. (2) to obtain a global 303 view, thus maintaining a consistent perspective of the global goal during learning. Moreover, agents 304 still use respective original rewards to learn for those samples that do not cause conflicts, such as 305 (2, 2) (shown in Fig. 2(c)), more generally, its value could be (a, b), 0 < a, b < 9 in this game. This is 306 critical for the generalization of tasks with dense rewards, and the achievement of efficient trade-offs 307 between local-level and global-level learning. Independent learning, at non-conflict steps, enables 308 agents to fully explore respective strategies that promise robust strategy combinations with enhanced 309 performances. The contrastive ablation experiments conducted in MPE (Section 4.2) validated this.

Remark 1. It can be inferred that, for general multi-dimensional action spaces that are (linearly or nonlinearly) denoted with multi-dimensional parameterized policies, the conflict-triggered mechanism still works in high-dimensional spaces.

313 314

281 282

3.3 MULTI-PLAYER MARKOV GAME

316 Section 3.2 investigates the problem formulated as a repeated matrix game involving two players, 317 without considering state transitions. For general Markov games with a number of participants, in 318 which state transitions are involved and rewards are associated with state-action pairs at each step, 319 we next show that DRI method remains valid. The mean-field theory (MFT) (Stanley, 1971) is em-320 ployed here to analyze agents in a networked system, it claims that each agent is affected only by 321 the mean effect from others around it that converts many-agent interactions into two-agent interactions. Since multi-agent MDP evolves in a networked environment, the pairwise approximation still 322 preserves global interactions between any pair of agents implicitly (Blume, 1993; Yang et al., 2018). 323 Note that Yang et al. (2018) utilize MFT to construct respective scalable joint action-value functions by assembling the average neighbor actions to achieve collaborative evaluation. This method is still incapable of dealing with saddle equilibriums due to value functions are learned by distinct payoffs. We further exploit the MFT concept to reconstruct payoffs for value functions. Accordingly, the conflict detection and reward reshaping of DRI method work between each agent $i, i \in \mathcal{N}$ and its averaged others in \mathcal{N}^{-i} , where \mathcal{N}^{-i} denoted the set of agents except *i*. Let *j* denotes the averaged \mathcal{N}^{-i} , and consider two agents *i* and *j* with the state value functions as

$$v_{\pi}^{i}(s) = \sum_{t=0}^{\infty} \mathbb{E}_{s \sim P, \boldsymbol{a} \sim \pi} \left[\left(r^{i}(s_{t}, \boldsymbol{a}_{t}) - \bar{\mu}_{\pi} | s_{0} = s \right) \right] \text{ and } v_{\pi}^{j}(s) = \sum_{t=0}^{\infty} \mathbb{E}_{s \sim P, \boldsymbol{a} \sim \pi} \left[\left(r^{j}(s_{t}, \boldsymbol{a}_{t}) - \bar{\mu}_{\pi} | s_{0} = s \right) \right],$$

333 where $\bar{\mu}_{\pi} = \frac{1}{2} (\mu_{\pi}^{i} + \mu_{\pi}^{j}).$

330 331 332

343

In cases, the previous matrix game for example, where the reward scale of both agents is equal, i.e., $\frac{1}{|\mathcal{S}|\cdot|\mathcal{A}|} \sum_{s \in \mathcal{S}, a \in \mathcal{A}} r^i(s, a) = \frac{1}{|\mathcal{S}|\cdot|\mathcal{A}|} \sum_{s \in \mathcal{S}, a \in \mathcal{A}} r^j(s, a), \text{ we have } \mathbb{E}\left[v_{\pi}^i(s)\right] = \mathbb{E}\left[v_{\pi}^j(s)\right] = 0. \text{ In}$ general cases, we have $\mathbb{E}\left[v_{\pi}^i(s)\right] = \frac{\mu_{\pi}^i - \mu_{\pi}^j}{2}, \mathbb{E}\left[v_{\pi}^j(s)\right] = \frac{\mu_{\pi}^j - \mu_{\pi}^i}{2}.$ To our method, the state value \tilde{v} is defined by the reshaped reward \tilde{r} of Eqn. (2), then we can obtain

$$\begin{cases} \bar{\mu}_{\boldsymbol{\pi}} = \frac{\tilde{\mu}_{\boldsymbol{\pi}}^{i} + \tilde{\mu}_{\boldsymbol{\pi}}^{j}}{2}, \quad \mathbb{E}\left[\tilde{v}_{\boldsymbol{\pi}}^{i}(s)\right] = \frac{\tilde{\mu}_{\boldsymbol{\pi}}^{i} - \tilde{\mu}_{\boldsymbol{\pi}}^{j}}{2}, \quad \mathbb{E}\left[\tilde{v}_{\boldsymbol{\pi}}^{j}(s)\right] = \frac{\tilde{\mu}_{\boldsymbol{\pi}}^{j} - \tilde{\mu}_{\boldsymbol{\pi}}^{i}}{2}, \\ \mathbb{E}\left[\tilde{v}_{\boldsymbol{\pi}_{t}}^{i}(s_{t+1}) - \tilde{v}_{\boldsymbol{\pi}_{t}}^{i}(s_{t})\right] = \mathbb{E}\left[\tilde{r}^{i}(\boldsymbol{a}_{t})\right] = \tilde{\mu}_{\boldsymbol{\pi}}^{i}, \quad \mathbb{E}\left[\tilde{v}_{\boldsymbol{\pi}_{t}}^{j}(s_{t+1}) - \tilde{v}_{\boldsymbol{\pi}_{t}}^{j}(s_{t})\right] = \mathbb{E}\left[\tilde{r}^{j}(\boldsymbol{a}_{t})\right] = \tilde{\mu}_{\boldsymbol{\pi}}^{j}, \end{cases}$$

where $\tilde{\mu}_{\pi}^{i}, \tilde{\mu}_{\pi}^{j}$ are the reshaped average reward. Now, the general one-step TD policy parameter update at step t is as follows

$$\theta_{t+1}^{i} = \theta_{t}^{i} + \beta_{t} \cdot \left[\tilde{r}_{t}^{i} - \bar{\mu}_{\pi_{t}} + \tilde{v}_{\pi_{t}}^{i}(s_{t+1}) - \tilde{v}_{\pi_{t}}^{i}(s_{t}) \right] \cdot \nabla_{\theta_{t}^{i}}, \\
\theta_{t+1}^{j} = \theta_{t}^{j} + \beta_{t} \cdot \left[\tilde{r}_{t}^{j} - \bar{\mu}_{\pi_{t}} + \tilde{v}_{\pi_{t}}^{j}(s_{t+1}) - \tilde{v}_{\pi_{t}}^{j}(s_{t}) \right] \cdot \nabla_{\theta_{t}^{j}},$$
(5)

where $\nabla_{\theta_i} = \nabla_{\theta_i} \log \pi(a_t^i | s_t)$ denotes the respective policy gradient term. One can see that in Eqn. (4), the conflict detection mechanism works on the horizontal $\bar{\mu}_{\pi_t}$. While in Eqn. (5), it works on the horizontal $\bar{\mu}_{\pi_t} - \tilde{\mu}_{\pi_t}^i$ for agent *i* and the $\bar{\mu}_{\pi_t} - \tilde{\mu}_{\pi_t}^j$ for agent *j*, it accurately eliminates the mismatch of different reward scales. An exponential decay concept was investigated by Qu et al. (2020) that ensures the scalability and provable convergence of networked MARL. Based on this work, we obtain the convergence Proposition of DRIMA, theoretical details are presented in Appendix A.2.

354 Note that the Nash equilibrium approached by DRIMA is the reshaped one of the original game. 355 DRI method reshapes original payoff (Q_t^1, \dots, Q_t^N) to $(\tilde{Q}_t^1, \dots, \tilde{Q}_t^N)$, $t = 1, 2, \dots$, that reshapes the solution space of each agent to avoid the saddle equilibrium point. Hence, the stationary points 356 of solution space only contains local optimum, global optimum, and inflection ones. Compared 357 with the saddle points that decentralized executive MARL generally can not escape, the sub-optimal 358 points can be avoided. The ability of sub-optimal avoidance mainly relies on the stochastic gradi-359 ent algorithm agents utilized, the nonlinear level of solution space concerning the task complexity, 360 and the way of collaboration. Conflict-triggered DRI enables the joint strategy to approach better 361 optimum with effective collaboration, Appendix B.1 presents an experimental examination. The 362 pseudo-code of DRI-based Actor-Critic is provided in Appendix C. 363

4 EXPERIMENTS

365 366

364

In this section, experiments of three scenarios are presented – the matrix game, the Multi-Agent 367 Particle Environment (MPE) (Lowe et al., 2017) and the StarCraft Multi-Agent Challenge (SMAC), 368 containing state and action spaces as discrete to continuous. DRIMA-based algorithms are labeled 369 as '-Dri', CTDE labeled as '-Ctde', and independent ones as '-Ind'. To draw training curves, all 370 algorithms run 5 times using 5 random seeds. All the value and policy functions in the matrix game 371 and MPE are parameterized by Multi-Layer Perceptrons without the recurrent structures and leave 372 off the parameter sharing to diminish interfering factors. For SMAC where the scenario and strategy 373 complexity increased, the recurrent NN embeds historical observations and parameter sharing en-374 sures training efficiency, we adopt these two techniques to test the efficacy of DRIMA. Experiments 375 are built on the PARL framework¹ and MAPPO benchmark (Yu et al., 2022), and our code will be 376 available on GitHub after this article is published.

¹https://github.com/PaddlePaddle/PARL

4.1 MATRIX GAME

378

379 380

382

384

385

394

397 398 Table 1: Maintain game where B-Y and C-Z are the NE, A-X is the global optimum.



Figure 3: Total rewards in matrix games.

Two 2-player matrix games, the Prisoner's Dilemma and the Maintain (Zhang et al., 2020), are tested 399 here. Their payoffs are shown in Fig. 1 and Table 1 respectively. The Dilemma has a unique NE, 400 in this case a saddle point, valued at (2, 2), while (10, 10) is the global optimum. For the Maintain 401 in which each agent has 3 actions, two NEs are valued at (10, 5) and (5, 10). The global optimum 402 at (20, 15) is a Pareto optimality but not a Nash equilibrium. Fig. 3(a) shows the learning curves of 403 total reward related to the independent and DRIMA algorithms in Dilemma, and Fig. 3(b) denotes 404 the results of Maintain. One can see that, the curves of value-based DQN show large fluctuations 405 due to the lack of strategy exploration and stochasticity. There always is a certain probability of 406 exploration, predefined by a (usually decreasing) factor ϵ , to the points that make the deviation of 407 trajectories. Overall, DRIMA enables DQN-Dri approaches to the global optimum in two tasks, 408 while independent DQNs all converge to poor equilibriums. For the deterministic policy, DDPG-409 Ind converge to the saddle NE in both tasks. With the DRI method, DDPG-Dri converges to the global optimum in Dilemma. In Maintain, it converges to a probability distribution of optimal and 410 suboptimal points due to insufficient exploration, but the result is still superior to saddle NE. The 411 stochastic policy A2C performs best with DRI that achieves the global optimum, without it, A2C-412 Ind falls into saddle points. Note that in Zhang et al. (2020), the optimum is achieved through a 413 centralized Bi-level training approach and requires the joint action vector. A 3-player matrix game 414 was further conducted, see Appendix B.1 for more discussion. 415

415 416 417

4.2 MULTI-AGENT PARTICLE ENVIRONMENT

418 For continuous space scenarios, we examine DRIMA in MPE. In this case, the interactions are al-419 lowed in a certain communication range, and agents are practically observe state from their own 420 perspective. Two types of tasks named Cooperative Hunting and Cooperative Collection are de-421 signed shown in Fig. 4(a) and (d) respectively. For Cooperative Hunting, the predators (blue and 422 green) are slower and try to catch the faster prey (red). The prey's moving strategy is pre-trained on 423 the rules to avoid being hit, it is allowed to observe within a limited view (marked by the red area) to ensure positive activity. Two types of agents consist of the predator group, in which the attacker 424 (blue) is assumed to have the firepower to hit the prey and finish the task, while the interceptors 425 (greens) are poorly equipped which can only hinder the movement of the target and avoid being hit. 426 The predator group needs to learn cooperative strategies to hit prey within a limited episode length 427 (30 steps). The rewards for all predators are set to negative values of their respective distances to 428 the prey, with the attacker receiving an additional positive reward for hitting the prey. The saddle 429 situation it might encounter is discussed in Appendix B.2.1 in detail. 430

431 We first design a hunting task, CH-I, composed of 1 prey, 1 attacker, and 2 interceptors. For multiagent DDPG and A2C, in CH-I, the predator group reward over the number of episodes are shown

449 450

451 452



Figure 4: Experimental results in Multi-Agent Particle Environment.

453 as Fig. 4 (b) and (c). One can see that both A2C-Dri and DDPG-Dri successfully learned the group 454 strategy of hitting the prey that obtains positive rewards. Based on A2C-Ind and DDPG-Ind, how-455 ever, the predator group can hardly achieve the capture goal and the final reward is almost zero. The '-Ctde', despite using a global reward that sums up all agents' rewards for learning, cannot guaran-456 tee learning a globally optimal policy. A2C-Ctde succeeded and DDPG-Ctde failed in this task. The 457 reason might be that higher stochasticity brings A2C moving policy greater flexibility and probabil-458 ity to explore a successful combination. The lack of stochasticity and reward distinction between 459 multiple agents might cause the failure of DDPG-Ctde. We further design a CH-II composed of 2 460 prey, 2 attackers, and 2 interceptors to exam an increased scale. Considering limited space, we put 461 the demonstration of the learned group policy, and the results of CH-II in Appendix B.2.2. 462

For Cooperative Collection (Fig. 4(d)), 6 agents are set to search and occupy 6 randomly located food. Both agents and food have three types, purple, blue, and green, and set the number of each type as 2. Agents need to collaboratively occupy as much food as possible and, meanwhile, occupy the food that matches their type as much as possible during one episode. Each agent can only observe objects within its field of view, it receives a negative reward based on the distance to the nearest food location, a positive reward when it occupies any food, and a maximum reward if the occupied food is of the same type. The saddle situation it might encounter is discussed in Appendix B.2.1 in detail.

Fig. 4(e) and (f) present the learning curves of DDPGs. The '-Dri', compared with '-Ctde' and '-470 Ind', achieves the best overall performance, and it does learn a collaborative strategy of occupying as 471 much food as possible while promoting agent-food type matching (Fig. 4(f)). Interestingly, a class 472 of A2C algorithms, which performs well in Cooperative Hunting, performs poorly in Cooperative 473 Collection. We attribute this to the existence of symmetric games, in which the identities of the 474 players can be changed without changing the payoff to the strategies. It occurs in tasks such as 475 occupying the same target and collision avoidance between homogeneous agents. More specifically, 476 if in a continuous space environment with multiple homogeneous agents, each adopts a policy with 477 high stochasticity, it would be difficult for agents to break the symmetric situation encountered. We provide a detailed discussion in Appendix B.2.3. 478

Meanwhile, the '-Dri-naive' algorithms that constantly average neighbors' rewards without the conflict trigger are performed, in contrast to the '-Dri'. For the previous matrix games with simple reward constructions, the constant average method can also obtain results comparable with DRIMA.
In the continuous space MPE, however, the general inferior performance of '-Dri-naive' shows the critical role of the conflict-triggered mechanism. The greater the scale and complexity of a task, the more significant it is (Fig. 4 and Fig. 7). Integrating information only when necessary can effectively balance individual-level and global-level learning, helping to learn a diversified and robust policy combination, that ensures an enhanced performance of MARL.

486 4.3 STARCRAFT MULTI-AGENT CHALLENGE

In this section, based on the MAPPO benchmark (Yu et al., 2022), we perform DRIMA combined with MAPPO in StarCraft II. In this battle case, agents in the Ally group strive to learn collaborative strategies to eliminate all agents in the Enemy group to win. All the learning agents are set to have access to a group reward that reflects the overall health status of the Enemy group. The distinctions are that for the '-Ctde', each agent receives a reward containing all ally members, while the '-Ind' and '-Dri' agents receive a reward only reflecting individual health status, moreover, the '-Dri' allows differential reward interactions within sight ranges.

495 Several representative tasks are performed, covering the homogeneous and symmetrical, homoge-496 neous and asymmetrical, heterogeneous and symmetrical, and heterogeneous and asymmetrical for-497 mations from the perspective of role composition and number of members; crossing easy, hard, and 498 super hard levels of map difficulty. The median win rate curves are shown in Fig. 5. It can be seen that in general tasks, the DRI algorithms possess a relatively slow rate at the beginning of learning. 499 One can infer that distributed information interactions for goal consensus during the learning pro-500 cess would affect the learning rate within an acceptable range. The learning results, nevertheless, 501 were not degraded. DRIMA based on its effective collaboration achieves at least comparable results 502 with the centralized benchmark MAPPO-Ctde, in a distributed way. It should be noticed that in the hard task 5m vs. 6m, DRIMA achieves an outperforming win rate, and we attribute this to its saddle 504 elimination capability. For complicated multi-agent tasks in which it is difficult to identify the sad-505 dle situations constructed by original reward functions, DRIMA enables MARL to avoid potential 506 saddle equilibriums during learning and approach to better optimum. As to the Ctde, a global reward 507 setting in such tasks may make it difficult to distinguish the individual contribution from the joint 508 behavior and to promote strategy combination.



Figure 5: Results in SMAC.

5 CONCLUSION

517

525

526 527

532 533

In this paper, we introduced DRIMA, a cooperative learning approach based on the conflict-triggered differential reward interaction for distributed MARL. We identify that the game induced by decentralized executions of multiple agents tends to limit overall performance, our approach efficiently eliminates strategies from converging to saddle equilibriums to solve this dilemma. Distributed learning promises great flexibility and robust capability, achieving ideal collaboration performance on this basis is an ongoing challenge. In the future, we will continue to study the remaining issues, such as breaking the symmetry of the game for multi-agent stochastic policies.

540 REFERENCES

547

553

554

555

559

560

561

- James Ault and Guni Sharon. Reinforcement learning benchmarks for traffic signal control. In
 Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks
 Track (Round 1), 2021.
- Lawrence E Blume. The statistical mechanics of strategic interaction. *Games and economic behavior*, 5(3):387–424, 1993.
- Jan Blumenkamp and Amanda Prorok. The emergence of adversarial communication in multi-agent reinforcement learning. In *Conference on Robot Learning*, pp. 1394–1414. PMLR, 2021.
- Ziyi Chen, Yi Zhou, Rong-Rong Chen, and Shaofeng Zou. Sample and communication-efficient decentralized actor-critic algorithms with finite-time analysis. In *International Conference on Machine Learning*, pp. 3794–3834. PMLR, 2022.
 - Tianshu Chu, Sandeep Chinchali, and Sachin Katti. Multi-agent reinforcement learning for networked system control. arXiv preprint arXiv:2004.01339, 2020.
- Wenlu Du, Junyi Ye, Jingyi Gu, Jing Li, Hua Wei, and Guiling Wang. Safelight: A reinforcement
 learning method toward collision-free traffic signal control. In *Proceedings of the AAAI confer- ence on artificial intelligence*, volume 37, pp. 14801–14810, 2023.
 - Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- David Earl Hostallero, Daewoo Kim, Sangwoo Moon, Kyunghwan Son, Wan Ju Kang, and Yung
 Yi. Inducing cooperation through reward reshaping based on peer evaluations in deep multi-agent
 reinforcement learning. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 520–528, 2020.
- Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.
- Jiechuan Jiang, Chen Dun, Tiejun Huang, and Zongqing Lu. Graph convolutional reinforcement
 learning. *arXiv preprint arXiv:1810.09202*, 2018.
- 572 Yiheng Lin, Guannan Qu, Longbo Huang, and Adam Wierman. Multi-agent reinforcement learning in stochastic networked systems. *Advances in neural information processing systems*, 34:7825– 7837, 2021.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multiagent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017.
- Chengdong Ma, Aming Li, Yali Du, Hao Dong, and Yaodong Yang. Efficient and scalable rein forcement learning for large-scale network control. *Nature Machine Intelligence*, pp. 1–15, 2024.
- 581
 582
 583
 584
 584
 585
 586
 587
 588
 588
 588
 589
 589
 589
 580
 581
 580
 581
 581
 581
 582
 583
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
 584
- Zhiyu Mou, Yu Zhang, Feifei Gao, Huangang Wang, Tao Zhang, and Zhu Han. Deep reinforcement
 learning based three-dimensional area coverage with uav swarm. *IEEE Journal on Selected Areas in Communications*, 39(10):3160–3176, 2021.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in neural information processing systems*, 32, 2019.
- Siddharth Nayak, Kenneth Choi, Wenqi Ding, Sydney Dolan, Karthik Gopalakrishnan, and Hamsa
 Balakrishnan. Scalable multi-agent reinforcement learning through intelligent information aggregation. In *International Conference on Machine Learning*, pp. 25817–25833. PMLR, 2023.

612

626

- Praveen Palanisamy. Multi-agent connected autonomous driving using deep reinforcement learning. In 2020 International Joint Conference on Neural Networks (IJCNN), pp. 1–7. IEEE, 2020.
- Dawei Qiu, Jianhong Wang, Zihang Dong, Yi Wang, and Goran Strbac. Mean-field multi-agent rein forcement learning for peer-to-peer multi-energy trading. *IEEE Transactions on Power Systems*, 2022.
- Chao Qu, Shie Mannor, Huan Xu, Yuan Qi, Le Song, and Junwu Xiong. Value propagation for de centralized networked deep multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Guannan Qu, Yiheng Lin, Adam Wierman, and Na Li. Scalable multi-agent reinforcement learning
 for networked systems with average reward. *Advances in Neural Information Processing Systems*, 33:2074–2086, 2020.
- Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning for multiagent networked
 systems. *Operations Research*, 70(6):3601–3628, 2022.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardel li, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The
 starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- H Eugene Stanley. *Phase transitions and critical phenomena*, volume 7. Clarendon Press, Oxford, 1971.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max
 Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition
 networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction.* MIT press, 2018.
- Eugene Vinitsky, Nathan Lichtlé, Xiaomeng Yang, Brandon Amos, and Jakob Foerster. Nocturne:
 a scalable driving benchmark for bringing multi-agent learning one step closer to the real world.
 Advances in Neural Information Processing Systems, 35:3962–3974, 2022.
- Hoi-To Wai, Zhuoran Yang, Zhaoran Wang, and Mingyi Hong. Multi-agent reinforcement learning via double averaging primal-dual optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Jianhong Wang, Wangkun Xu, Yunjie Gu, Wenbin Song, and Tim C Green. Multi-agent reinforce ment learning for active voltage control on power distribution networks. *Advances in Neural Information Processing Systems*, 34:3271–3284, 2021.
- Jiachen Yang, Ang Li, Mehrdad Farajtabar, Peter Sunehag, Edward Hughes, and Hongyuan Zha.
 Learning to incentivize other learning agents. *Advances in Neural Information Processing Systems*, 33:15208–15219, 2020.
- Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multiagent reinforcement learning. In *International conference on machine learning*, pp. 5571–5580. PMLR, 2018.
- Yuxuan Yi, Ge Li, Yaowei Wang, and Zongqing Lu. Learning to share in multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 35:15119–15131, 2022.
- Chao Yu. Hierarchical mean-field deep reinforcement learning for large-scale multiagent systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 11744–11752, 2023.
- Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- Haifeng Zhang, Weizhe Chen, Zeren Huang, Minne Li, Yaodong Yang, Weinan Zhang, and Jun
 Wang. Bi-level actor-critic for multi-agent coordination. In *Proceedings of the AAAI Conference* on Artificial Intelligence, volume 34, pp. 7325–7332, 2020.

648 649 650	Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi- agent reinforcement learning with networked agents. In <i>International Conference on Machine</i>
651	<i>Learning</i> , pp. 5872–5881. PMLR, 2018.
650	
652	
000	
054	
000	
000	
007	
000	
609	
000	
661	
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
080	
600	
002	
003	
004	
680	
080	
687	
688	
689	
690	
691	
692	
604	
094	
695	
090	
097	
600	
700	
700	
7.01	

A THEORETICAL MATERIALS

A.1 DEFINITIONS OF NASH EQUILIBRIUMS 705

The Nash equilibrium of the stochastic game denotes a joint strategy where each agent's response is the best to the others', according to Hu & Wellman (2003), the definition is as follows.

Definition 1. (Nash equilibrium) In stochastic game Γ , a Nash equilibrium point is a tuple of Nstrategies $(\pi_*^1, \dots, \pi_*^N)$ such that for all $s \in S$ and $i = 1, \dots, N$,

710 711

714 715 716

720

729

738 739 740

741

702

703

$$v^{i}_{\pi_{*}^{1},\cdots,\pi_{*}^{N}}(s) \geq v^{i}_{\pi_{*}^{1},\cdots,\pi_{*}^{i-1},\pi^{i},\pi_{*}^{i+1},\cdots,\pi_{*}^{N}}(s)$$

for all $\pi^i \in \Pi^i$, where Π^i is the set of strategies available to agent *i*. Correspondingly, the Nash *Q*-function of agent *i* is

$$Q_*^i(s, a^1, \cdots, a^N) = r^i(s, a^1, \cdots, a^N) + \sum_{s' \in S} P(s'|s, a^1, \cdots, a^N) v^i(s', \pi_*^1, \cdots, \pi_*^N),$$

717 where $r^i(s, a^1, \dots, a^N)$ is agent *i*'s one-period reward in state *s* and under joint action 718 $(a^1, \dots, a^N), v^i(s', \pi_*^1, \dots, \pi_*^N)$ is agent *i*'s expected reward over infinite periods starting from 719 state *s'* given that agents follow the equilibrium strategies.

The strategies are assumed to be stationary, that is, the rule of choosing an action is the same in every stage t, to ensure the existence of the Nash equilibrium point. Meanwhile, a stochastic game can be seen as composed of many stage games (one-period games) with the state transition. We focus on stage games with special types of Nash equilibrium points: global optima, and saddles, the related definitions are given.

Definition 2. (Global optimal point) A joint strategy $(\sigma^1, \dots, \sigma^N)$ of the stage game (M^1, \dots, M^N) is a global optimal point if every agent receives its highest payoff at this point. That is, for all k, $=M^k \geq \hat{a}M^k$ for all $\hat{a} \in \sigma(A)$

$$\sigma M^k \geq \hat{\sigma} M^k$$
 for all $\hat{\sigma} \in \sigma(\mathcal{A})$

730 A global optimal point is always a Nash equilibrium.

Definition 3. (Saddle point) A joint strategy $(\sigma^1, ..., \sigma^N)$ of the stage game $(M^1, ..., M^N)$ is a saddle point if (1) it is a Nash equilibrium, and (2) each agent would receive a higher payoff when at least one of the other agents deviates. That is, for all k,

$$\sigma^{k} \sigma^{-k} M^{k} \ge \hat{\sigma}^{k} \sigma^{-k} M^{k} \text{ for all } \hat{\sigma}^{k} \in \sigma(\mathcal{A}^{k}),$$
$$\sigma^{k} \sigma^{-k} M^{k} \le \sigma^{k} \hat{\sigma}^{-k} M^{k} \text{ for all } \hat{\sigma}^{-k} \in \sigma(\mathcal{A}^{-k}).$$

 M^k is agent k's payoff function over the space of joint actions, $M^k = \{r^k(a^1, ..., a^N) | a^1 \in \mathcal{A}^1, ..., a^N \in \mathcal{A}^N\}$, and r^k is the reward of agent k.

A.2 CONVERGENCE ANALYSIS OF DRIMA

(0,

Qu et al. (2020) investigated the exponential decay property that ensures the scalability of networked
 MARL through the proposed theoretical conditions. Following this work, we obtain the Lemma below.

if $j \notin N^i$,

745 Lemma 1. Define 746

747 748

751

$$C_{i,j} = \begin{cases} \sup_{s,a^{i},a^{j}} \sup_{s,a^{i},a^{\prime j}} \operatorname{TV}\left(P\left(\cdot|s,a^{i},a^{j}\right), P\left(\cdot|s,a^{i},a^{\prime j}\right)\right), & \text{if } j \in N^{i}, j \neq i, \\ \sup_{s,a^{i}} \sup_{s,a^{\prime i}} \operatorname{TV}\left(P\left(\cdot|s,a^{i}\right), P\left(\cdot|s,a^{\prime i}\right)\right), & \text{if } j = i, \end{cases}$$

where $TV(\cdot, \cdot)$ is the total variation distance between two distributions. Further, define the DRI *Q*-function as

755
$$\tilde{Q}_{\theta}^{i}(s,a^{i}) = \sum_{t=0}^{\infty} \mathbb{E}_{a^{-N^{i}} \sim \pi^{-N^{i}}, s \sim P} \Big[\tilde{r}^{i}(s,a) - \bar{\mu}_{\theta}^{N^{i}} | s_{0} = s, a_{0}^{i} = a^{i} \Big],$$

756 where \tilde{r}^i follows Eqn. (2). If for all $i \in \mathcal{N}, \sum_{j=1}^N C_{ij} \leq \rho < 1$ and $|r^i(s, \boldsymbol{a})| \leq \bar{r}, \forall (s, \boldsymbol{a}) \in \mathcal{S} \times \mathcal{A}$, 757 then the $\left(\frac{\bar{r}}{1-\rho},\rho\right)$ exponential decay property holds, i.e., 758

759

760 761 762

763

791

793

794 795

796

797 798

805

 $\left|\tilde{Q}_{\theta}^{i}\left(s,a^{i},\boldsymbol{a}^{N^{-i}}\right) - \tilde{Q}_{\theta}^{i}\left(s,a^{i},\boldsymbol{a}^{N^{-i}}\right)\right| \leq \frac{\bar{r}}{1-a}\rho^{2}$

holds for any $i \in \mathcal{N}$, $s \in \mathcal{S}$, $a^i \in \mathcal{A}^i$, $a^{N^{-i}}$, $a'^{N^{-i}} \in \mathcal{A}^{N^{-i}}$.

764 *Proof.* The DRI method of Eqn. (2) reshapes original reward r^i to \tilde{r}^i and changes distribution of states and actions. Two conditions that ensures the exponential decay property of original Q-765 functions still establish for the reshaped \tilde{Q} , i.e., $|\tilde{r}^i(s, a)| \leq \bar{r}$ and $\sum_{i=1}^N C_{ij} \leq \rho < 1$ hold. From the 766 proof of Theorem 1 in Qu et al. (2020), our Lemma 1 hold. 767

768 In Lemma 1, C_{ij} is the maximum possible of the distribution of next state as a result of a change in 769 node j'th current actions. The "change" is measured in total-variation distance, it can be interpreted 770 as the strength of node j's interaction with node i. The Lemma gives constraints on multi-agent 771 networks that require the interaction strength between each pair of agents to be small enough and 772 each agent to have a small enough neighborhood. It ensures the truncated Q-function $\tilde{Q}^i_{A}(s, a^i)$ is 773 a good approximation of the global one $\tilde{Q}^i_{\theta}(s, a)$. Numerical experiments in Qu et al. (2020) show 774 the broad establishment of the exponential decay property, and the empirical results in our paper 775 demonstrate its practicability. In the following, standard assumptions in reinforcement learning are 776 listed as the basis of theoretical convergence. The first one is that the rewards are bounded.

777 Assumption 1. For all *i*, and $s \in S$, $a \in A^1 \times \cdots \times A^N$, we have $0 \le r^i$ $(s, a) \le \overline{r}$. 778

779 The next is the condition on the step sizes for the two-time scale algorithm. 780

Assumption 2. The positive step size α_t , β_t are deterministic, non-increasing, square summable but 781 not summable, i.e. $\sum_{t=0}^{\infty} \alpha_t = \sum_{t=0}^{\infty} \beta_t = \infty$, $\sum_{t=0}^{\infty} (\alpha_t^2 + \beta_t^2) < \infty$, and satisfy $\sum_{t=0}^{\infty} (\frac{\beta_t}{\alpha_t})^d < \infty$ 782 783 for some d > 0. 784

785 The next one is the uniformly ergodic of the underlying MDP. 786

Assumption 3. Under any policy $\theta = (\theta^1, \dots, \theta^N)$, the induced Markov chain over the state-787 action space $S \times A$ is ergodic with stationary distribution π^{θ} . Further, (a) For all $(s, a) \in S \times A$, 788 $\pi(a|s) \ge \sigma$ for some $\sigma > 0$. (b) $\left\| P^{\theta} - \mathbf{1} \left(\pi^{\theta} \right)^T \right\|_{D^{\theta}} \le \mu_D$ for some $\mu_D \in (0, 1)$, where $D^{\theta} = 0$ 789 790 diag $(\pi^{\theta}) \in \mathbb{R}^{S \times A \times S \times A}$ and $\|\cdot\|_{D^{\theta}}$ is the weighted Euclidean norm $\|x\|_{D^{\theta}} = \sqrt{x^T D^{\theta} x}$ for vectors $x \in \mathbb{R}^{S \times A \times S \times A}$, and the corresponding induced norm for matrices. 792

The following is that the gradient is bounded and is Lipschitz continuous.

Assumption 4. For each i, a_i, s , $\|\nabla_{\theta^i} \log \pi_{\theta^i}(a^i | s)\| \leq L_i$. Let $L := \sqrt{\sum_{i \in \mathcal{N}} L_i^2}$. Further, $\nabla_{\theta^i} \log \pi_{\theta^i}(a^i | s)$ is L'-Lipschitz continuous in θ^i , and $\nabla J(\theta)$ is L'-Lipschitz continuous in θ .

799 Moreover, following the analysis path in Hu & Wellman (2003), our convergence result requires 800 that the (original) stage games encountered during learning have global optima, or alternatively, that 801 they all have saddle points. 802

Assumption 5. Every stage game (Q_t^1, \dots, Q_t^N) , for all t and s, has a global optimal point or a 803 saddle point. 804

It would be unusual for stage games during learning to maintain adherence to Assumption 5. Never-806 theless, the experiments in Hu & Wellman (2003) and in our paper practically show that convergence 807 is not necessarily so sensitive to the properties of the stage games during learning. This condition 808 established for convergence proof may have the potential to relax, which is beyond the scope of this 809 paper. Finally, we obtain the following convergence proposition.

Proposition 1. For an N-player stochastic game, Lemma 1 and Assumption 1-5 hold, the sequence $\tilde{Q}_t = (\tilde{Q}_t^1, \cdots, \tilde{Q}_t^N)$, updated by

$$\begin{split} \mu_{t+1}^{i} &= (1 - \alpha_{t}) \, \mu_{t}^{i} + \alpha_{t} \cdot r^{*}(s_{t}, \boldsymbol{a}_{t}) \,, \\ \tilde{Q}_{t+1}^{i} \left(s_{t}, a_{t}^{i}\right) &= (1 - \alpha_{t}) \, \tilde{Q}_{t}^{i} \left(s_{t}, a_{t}^{i}\right) + \alpha_{t} \cdot \left(\tilde{r}^{i}(s_{t}, \boldsymbol{a}_{t}) - \bar{\mu}_{t}^{N^{i}} + \tilde{Q}_{t}^{i} \left(s_{t+1}, a_{t+1}^{i}\right)\right) \\ \theta_{t+1}^{i} &= \theta_{t}^{i} + \beta_{t} \cdot \tilde{Q}_{t}^{i} \left(s_{t}, a_{t}^{i}\right) \cdot \nabla_{\theta^{i}} \log \pi_{\theta_{t}^{i}} \left(a_{t}^{i}|s_{t}\right) \,, \end{split}$$

818 where \tilde{r}^i follows Eqn. (2), converges to an approximated Nash equilibrium point $(\pi^1_*, \cdots, \pi^N_*)$ 819 with an asymptotic bound as $\liminf_{t\to\infty} \|\nabla J(\theta_t)\| \leq L \frac{\bar{r}\rho^2}{(1-\mu_D)(1-\rho)}$, and the $\tilde{Q}_* = (\tilde{Q}^1_*, \cdots, \tilde{Q}^N_*)$ 821 constitutes a Nash Q-value.

From Proposition 1, the DRI-based Actor-Critic algorithm can find an approximated stationary point with gradient size $O(\rho^2)$. Utilizing the game theory in Hu & Wellman (2003) and the stochastic approximation in Qu et al. (2020), Proposition 1 can be proved through the analogous procedure, that we omit here.

Remark 2. Each average reward tracker μ_t^i is updated by the respective original rewards r_{t+1}^i rather than the reshaped one \tilde{r}_{t+1}^i . This is critical for agents to restructure the original global objective of Eqn. (1) by $\bar{\mu}^{N^i} = \frac{1}{N^i} \sum_{i \in N^i} \mu_{\pi}^j$.

B Additional Experimental Results

B.1 A THREE-PLAYER MATRIX GAME



A 3-player matrix game is provided in Fig. 6(a) from Matthew Rousd's Game Theory course. ² It is a question about meeting during COVID. Two students, Gus and Henry, want to hang out together on a Friday night. They are deciding where to go to drink some juices. They also know Scooter and don't mind being around him, but they are a bit leery as Scooter is suspected to be a member of the Secret COVID Policy Force. They think Scooter might bust students for COVID violations if multiple people try to gather in an unapproved indoor space together. Gus and Henry could choose to

²www.youtube.com/watch?v=DdLqwsMKjMc.

go in Henry's room, in an empty classroom, or in the cornfield. Scooter doesn't know the classroom
 is an option, so will either visit Henry's room or the cornfield.

To solve this question, the communication network of the '-Dri's is set to be fully connected, and different '-Ind' and '-Dri' algorithms are performed. The total reward curves in Fig. 6 show that DRIMA enables three agents to learn the destination strategies with the optimal global reward, i.e., to meet in the cornfield valued at (7,7,6). The '-Ind's may converge to the global optimum or the inferior ones, both are Nash equilibriums. In other words, in general games with multiple equilibria, the strategy convergence result of decentralized MARL is sensitive to the initial condition and lacks collaboration to approach the global optimum. The DRI method enables agents to collaboratively achieve better solutions, in this case, the global optimal one.

874 875

876

877 878 879

880

882 883

885

887 888

889

B.2 MULTI-AGENT PARTICLE ENVIRONMENT

B.2.1 SADDLE SITUATIONS IN MPE TASKS

Table 2: Payoff matrix for Cooperative Hunting, where (5,0) is the saddle point and (15,-5) is the global optimum.

		Intercepter	
		Approach	Avoidance
Attacker	Approach	15, -5	5,0
ALIACKEI	Avoidance	0, -5	0, 0

Table 3: Payoff matrix for Cooperative Collection, where (10,5) is the saddle point and (20,0) is the global optimum.

		Purple	
		Approach	Leave
Groon	Approach	10, 5	20, 0
Green	Leave	0, 10	0,0

For the continuous scenario defined with the continuous reward function, we can simplify and dis-896 cretize agents' rewards as the matrix game to facilitate analysis while not affecting the validity. 897 Table 2 denotes a saddle situation agents in the hunting task would encounter. When the attacker 898 and interceptor learn independently, the former (with firepower) would actively approach and attack 899 prey, while the latter (without firepower) would approach but avoid getting too close to the target and 900 receiving collision damage. This could lead to two of policies to the saddle point (valued at (5, 0)). 901 To get the optimal point(valued at (15, -5)), agents are required to learn to choose the "Approach" 902 cooperatively. Table 3 is set to denote a saddle issue in the collection task, where two types of a-903 gents (Green and Purple) surround a Green food. In this case, the lack of cooperation agents would 904 separately attempt to occupy the food (which one would succeed is stochastic) leading to the saddle equilibrium (valued at(10,5)). To achieve the global optimum (valued at(20,0)), the Purple agent 905 needs to learn to leave this food occupied by the Green agent to obtain the highest type-matching 906 reward. The actual situation of these two tasks is more complex since more agents are involved 907 and rewards are dense, highly effective collaborative learning is required to achieve the ideal overall 908 performance. 909

910 911 B.2.2 COOPERATIVE HUNTING

The group strategies of success and failure in Cooperative Hunting I are demonstrated in Fig. 7(a) and (b) respectively. As can be seen in Fig. 7(a), the two interceptors (greens) of the predators learned the cooperative strategy that surrounding the prey from outside the field of view and blocking its escape route to help the attacker(blue) hit the prey. Fig. 7(b), however, shows that three predators chased the prey from the same direction without forming an encirclement. Since the prey is faster, it can easily escape from blank directions. The DRI method enables multiple agents to learn collaborative strategies successfully.



Figure 7: Results in Cooperative Hunting.

Fig. 7(c) and (d) are the results of Cooperative Hunting II, in which 2 attackers and 2 interceptors hunt 2 prey. It can be seen that the performance of the centralized algorithm '-Ctde' decreases significantly as the number of agents increases, and the global reward can no longer accurately discriminate the situation of strategy combination. The '-Ind's are more likely to catch prey and obtain positive rewards due to the number of agents in the field increasing. The '-Dri's receive the highest reward and learn the best overall strategies through effective local information interactions, which are not affected by the growth number of agents.

B.2.3 RESULTS OF A2C IN COOPERATIVE COLLECTION



Figure 8: Results of A2Cs in Cooperative Collection.

In the Cooperative Collection, where 6 agents randomly search and occupy 6 foods, the stochastic policy algorithm fails to learn a good strategy combination. The situation A2C faced is pictured in Fig. 8(a), the total reward and match count are shown in Fig. 8(b) and (c). From Fig. 8(a), one food happens to be within the view of the four agents, and its position is considered to be the closest to each agent. These agents simultaneously approach and attempt to occupy that food. Since the space is continuous and each agent adopts a highly random exploration strategy, no agent has the advantage in occupying the target, i.e., no one of them successfully settles in the target without being pushed out by others. Agents are at an impasse. In contrast with A2C, multiple agents using DDPG are easy to break this deadlock due to their more deterministic and less stochastic policies. Theoretically, this situation constitutes a symmetric game in which the identities of the players can be changed without changing the payoff to the strategies. In the cases that coexist symmetric games, continuous space, and stochastic policy, the independent and simultaneous policy MARL will face challenges. We leave this problem as future work.

C PSEUDO-CODE OF DRIMA

Algorithm 1 represents the pseudocode of the Differential Reward Interaction-based Actor-Critic, where the individual V-function and policy are parameterized by ω^i and θ^i respectively. Other algorithm variants of RL can be obtained by analogy.

972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982		
983		
984		
985		
986		
987		
988		
989	Alge	orithm 1 Differential reward interaction-based Actor-Critic
990	1:	Input: Initial value of parameters $\mu_0^i, \omega_0^i, \theta_0^i, \forall i \in \mathcal{N}$, the initial state s_0 , and step-sizes $\{\alpha_t\}_{t>0}$
991		and $\{\beta_t\}_{t>0}$.
992		Initialize the iteration counter $t \leftarrow 0$.
993	2:	repeat
994	3:	for all $i \in \mathcal{N}$ do
995	4:	Each agent <i>i</i> executes action $a_t^i \sim \pi_{\theta_t^i}^i (\cdot s_t)$.
990	5:	Observe state s_{t+1} , and reward $r_t^i = r^i (s_t, a_t^i)$.
998	6:	Update $\mu_{t+1}^i \leftarrow \mu_t^i + \alpha_t \cdot (r_t^i - \mu_t^i)$.
999	7:	end for
1000	8:	for all $i \in \mathcal{N}$ do
1001	9:	Receive μ_t^j, r_t^j from neighbors $j \in \mathcal{N}^i$.
1002	10:	Calculate $\tilde{r}_t^i - \bar{\mu}_t^{N^*}$ according to Eqn. 2.
1003	11:	Update $\delta_t^i \leftarrow \tilde{r}_t^i - \bar{\mu}_t^{N^i} + V_{\omega_t^i}^i(s_{t+1}) - V_{\omega_t^i}^i(s_t).$
1004	12:	Critic step: $\omega_{t+1}^i \leftarrow \omega_t^i + \alpha_t \cdot \delta_t^i \cdot \nabla_{\omega_t^i} V_{i,i}^i (s_t).$
1005	13.	Actor step: $\theta^i \rightarrow \theta^i + \beta_i \cdot \delta^i \cdot \nabla_{\alpha_i} \log \pi^i \cdot (a^i s_i)$
1006	14.	and for
1007	14. 15·	Undate the iteration counter $t \leftarrow t + 1$
1008	16:	until Convergence
1009		
1010		
1011		
1012		
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020		
1021		
1022		
1023		
1024		
1025		