

---

# Scalable Constrained Policy Optimization for Safe Multi-agent Reinforcement Learning

---

Lijun Zhang<sup>1</sup>, Lin Li<sup>1</sup>, Wei Wei<sup>1\*</sup>, Huizhong Song<sup>1</sup>, Yaodong Yang<sup>2</sup>, Jiye Liang<sup>1</sup>

1. Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, School of Computer and Information Technology, Shanxi University, Taiyuan, Shanxi, China.
2. Institute for AI, Peking University, Beijing, China.

## Abstract

A challenging problem in seeking to bring multi-agent reinforcement learning (MARL) techniques into real-world applications, such as autonomous driving and drone swarms, is how to control multiple agents safely and cooperatively to accomplish tasks. Most existing safe MARL methods learn the centralized value function by introducing a global state to guide safety cooperation. However, the global coupling arising from safety constraints and the exponential growth of the state-action space size limit their applicability in instant communication or computing resource-constrained systems and larger multi-agent systems. In this paper, we develop a novel scalable and theoretically-justified multi-agent constrained policy optimization method. This method integrates the rigorous bounds of the trust region method and the bounds of the truncated advantage function to provide a new local policy optimization objective for each agent. Also, we prove that the safety constraints and the joint policy improvement can be met when each agent adopts a sequential update scheme to optimize a  $\kappa$ -hop policy. Furthermore, we propose a practical algorithm called Scalable MAPPO-Lagrangian (Scal-MAPPO-L). The proposed method’s effectiveness is verified on a collection of benchmark tasks, and the results support our theory that decentralized training with local interactions can still improve reward performance and satisfy safe constraints.

## 1 Introduction

With the advanced and rapid developments of reinforcement learning technology, many researchers have gradually shifted their focus from virtual simulation to real-world cyber-physical applications [1, 2]. In this process, safety challenges are inevitable, especially in multi-agent safety-critical scenarios, e.g., autonomous vehicle navigation [3], power grids [4], and drone swarms [5], in which agents perform complex cooperative tasks while adhering to a variety of local and system-wide limitations or constraints. These constraints can be derived from domain-specific knowledge and are intended to prevent damage to people or other environmental elements, such as equipment and infrastructure, or to prevent the inability to accomplish specific tasks or objectives. Take multi-robot control as an example. Each running robot must not take certain actions or not visit certain states, which may imply unsafe for itself, its collaborators, or the infrastructure of its environment [6]. These widespread potential dangers exacerbate the difficulty of safety decision-making when applying MARL. Consequently, it is necessary to research the safe decision-making problem in MARL to ensure that agents can work together safely and cooperatively to accomplish tasks.

---

\*Correspondence to <weiwei@sxu.edu.cn>.

There are two main approaches concerning safe MARL techniques in the existing literature. The first type is shielded-based reactive methods [7, 8], which combine environmental dynamics and safety specification constraints to predict whether the actions chosen by agents will violate cost constraints. Nevertheless, due to the reliance on precise modeling knowledge, these methods may lead to poor performance when the accurate state transition model is unavailable. The second type formulates the safe MARL problem as a constrained Markov game, which requires agents to solve a constrained optimization problem, i.e., maximize total reward while avoiding violating cost constraints. To mention a few, several safe MARL variants, such as CMIX [9] and MAPPO-L [10], have been proposed, which learn the centralized value function to overcome policy conflicts caused by the partially observable and non-stationarity nature of the environment faced by each agent. Unfortunately, the global coupling arising from agents’ safety constraints and the exponential growth of the state-action space size make the usability of these algorithms in instant communication or computing resource-constrained systems and the scalability in larger multi-agent systems become a bottleneck, limiting their applicability.

A promising approach for avoiding these shortcomings, which has received attention in recent years, is to exploit networked application-specific structures. For example, Safe Dec-PG [11] employs a primal-dual framework to find the saddle point between maximizing decoupled rewards and minimizing costs under a consensus network. However, it is worth noting that this approach still assumes each agent can access the global state and requires that the actions of all neighboring agents on the network be available. Recent research [12] proposes a scalable safe MARL approach based on the spatial decay assumption of the environment dynamics, which updates the policies of agents by the truncated gradient estimators depending on the local states and actions of the  $\kappa$ -hop neighboring agents. However, due to the dependence on the actions and states of its neighbors, this method necessarily involves joint training in a local area, which is still plagued by non-stationary issues. Motivated by the urgent desire for scalable learning in practical applications and the fact that meeting both safety constraints and joint policy improvement is challenging for most methods, we investigate a novel scalable safe MARL with theoretical analysis, practical algorithm, and simulation verification.

Specifically, we focus on decentralized learning settings without global observability, where each agent can only access the local state information of itself and its neighbors. Our main contributions are summarized as follows.

- We develop a novel scalable multi-agent constrained policy optimization method that eliminates dependence on the global state and other agent actions during each agent’s training. Furthermore, we parameterize each agent’s policy and propose a practical algorithm called Scalable MAPPO-Lagrangian (Scal-MAPPO-L).
- We quantify the maximum information loss regarding the advantage truncation based on two assumptions about the transition dynamics and policies. Then, each agent’s new local policy optimization objective is provided by integrating the rigorous bounds of the trust region method and the bounds of the truncated advantage function. In addition, we prove that the safety constraints and the joint policy improvement can be guaranteed when updating the local policy with a sequential update scheme.
- Experimentally, we provide the results on several safe MARL tasks to evaluate the effectiveness of our proposed method and the sensitivity for the parameter  $\kappa$ . The results support our theory that decentralized training with local interactions can still improve reward performance and satisfy safe constraints.

## 2 Preliminaries

### 2.1 Constrained Markov game

Consider a safe MARL problem subject to multiple constraints, where each agent are associated with an underlying undirected graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ . Here,  $\mathcal{N} = \{1, \dots, n\}$  is the set of  $n$  agents and  $\mathcal{E} \subset \mathcal{N} \times \mathcal{N}$  is the set of edges. The problem can be formulated as a constrained Markov game  $\langle \mathcal{N}, \mathcal{S}, \mathcal{A}, P, \rho_0, \gamma, \mathbf{R}, \mathbf{C}, \mathbf{c} \rangle$ .  $\mathcal{S} = \times_{i \in \mathcal{N}} \mathcal{S}^i$  and  $\mathcal{A} = \times_{i \in \mathcal{N}} \mathcal{A}^i$  are the state and action spaces, which are the product of local spaces; global state  $\mathbf{s} = (s^1, \dots, s^n)$  and joint action  $\mathbf{a} = (a^1, \dots, a^n)$  for any  $\mathbf{s} \in \mathcal{S}$  and  $\mathbf{a} \in \mathcal{A}$ .  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the probabilistic transition dynamics function, which

satisfies the Dobrushin condition [13] as follows:

$$W^{ij} = \sup_{z^j, z'^j, z^{-j}} \left\| P^i(\cdot | z^j, \mathbf{z}^{-j}) - P^i(\cdot | z'^j, \mathbf{z}^{-j}) \right\|_1, \quad (1)$$

where  $z^j = (s^j, a^j)$  and  $z'^j = (s'^j, a'^j)$  represent two different state-action pairs of the agent  $j$  respectively, and  $\mathbf{z}^{-j}$  represents the state-action pair of the agent other than  $j$ . The value of  $W^{ij}$  reflects the extent to which the local transition probability of agent  $i$  is affected by the state and action of agent  $j$ .  $\rho_0$  is the initial state distribution,  $\gamma \in [0, 1)$  is the discount factor.  $\mathbf{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the joint reward function,  $\mathbf{C} = \{C_j^i\}_{1 \leq j \leq m^i}^{i \in \mathcal{N}}$  is the sets of cost functions (every agent  $i$  has  $m^i$  cost functions) of the form  $C_j^i : \mathcal{S}^i \times \mathcal{A}^i \rightarrow \mathbb{R}$ , and finally the set of corresponding cost values is given by  $\mathbf{c} = \{c_j^i\}_{1 \leq j \leq m^i}^{i \in \mathcal{N}}$ .

At each timestep  $t$ , every agent  $i$  is in a state  $s_t^i$ , and takes an action  $a_t^i$  according to its policy  $\pi^i = (a^i | s_t^i)$ . Together with other agents actions, it gives a joint action  $\mathbf{a}_t = (a_t^1, \dots, a_t^n)$  and the joint policy  $\pi = \prod_{i=1}^n \pi^i(a^i | s_t^i)$ . The agents receive the reward  $\mathbf{R}(s_t, \mathbf{a}_t)$ , meanwhile each agent  $i$  pays the costs  $C_j^i(s_t^i, a_t^i), \forall j = 1, \dots, m^i$ , and all agents have a joint goal, i.e., maximizing the expected total reward of

$$J(\pi) \triangleq \mathbb{E}_{\mathbf{s}_0 \sim \rho_0, \mathbf{a}_0, \infty \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{R}(s_t, \mathbf{a}_t) \right], \quad (2)$$

meanwhile satisfying every agent  $i$ 's safety constraints, written as

$$J_j^i(\pi) \triangleq \mathbb{E}_{\mathbf{s}_0 \sim \rho_0, \mathbf{a}_0, \infty \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t C_j^i(s_t, \mathbf{a}_t) \right] \leq c_j^i, \forall j = 1, \dots, m^i. \quad (3)$$

## 2.2 Spatial correlation decay

Exponential decay property [13, 14], also known as spatial correlation decay, is a powerful property associated with local interactions, which says that the impact of agents on each other decays exponentially in their graph distance. More information about spatial correlation decay is presented in Appendix B.1. Here, inspired by [15], we make the following two assumptions for the spatial correlation of the transition dynamics and policies. We use the notation  $\pi^i(\cdot | s_{\mathcal{N}_\kappa^i})$  for  $\kappa$ -hop policies, where  $s_{\mathcal{N}_\kappa^i}$  represents the state of agent  $i$ 's  $\kappa$ -hop neighbors. It may be replaced with  $\pi_\kappa^i$  for simplicity when it is clear from context.

**Assumption 2.1.** (Spatial Decay of Correlation for the Dynamics) Assume that there exist  $\beta > 0$  in (1), for any agents  $i, j \in \mathcal{N}$ , such that

$$\max_{i \in \mathcal{N}} \sum_{j \in \mathcal{N}} e^{\beta d(i, j)} W^{ij} \leq \zeta, \quad (4)$$

where  $d(i, j)$  represents the distance between agent  $i$  and agent  $j$ , and  $\zeta \in [0, 2/\gamma)$  is a constant.

**Assumption 2.2.** (Spatial Decay of Correlation for the Policies) Assume that there exist  $\xi, \beta \geq 0$  such that for any agent  $i \in \mathcal{N}$ ,  $s_{\mathcal{N}_\kappa^i} \in \mathcal{S}_{\mathcal{N}_\kappa^i}$ ,  $s_{\mathcal{N}_\kappa^{-i}} \in \mathcal{S}_{\mathcal{N}_\kappa^{-i}}$ ,  $s'_{\mathcal{N}_\kappa^{-i}} \in \mathcal{S}_{\mathcal{N}_\kappa^{-i}}$ , one have

$$\sup_{s_{\mathcal{N}_\kappa^i}, s_{\mathcal{N}_\kappa^{-i}}, s'_{\mathcal{N}_\kappa^{-i}}} \left| \pi^i \left( \cdot | s_{\mathcal{N}_\kappa^i}, s_{\mathcal{N}_\kappa^{-i}} \right) - \pi^i \left( \cdot | s_{\mathcal{N}_\kappa^i}, s'_{\mathcal{N}_\kappa^{-i}} \right) \right| \leq \xi e^{-\beta \kappa}. \quad (5)$$

Assumption 2.2 reveals how much information is lost compared with access to the global state and allows us to consider a policy class with the necessary properties for the optimal policy under Assumption 2.1. More information is stated in Appendix B.2.

## 3 Scalable constrained policy optimization

This section develops a novel scalable and theoretically-justified multi-agent constrained policy optimization method and proposes a practical algorithm, i.e., Scal-MAPPO-L, by parameterizing each agent's policy. Specifically, we first quantify the maximum information loss regarding the

advantage truncation based on the spatial correlation decay property of the transition dynamics and policies. Then, the rigorous bounds of the trust region method and the bounds of the truncated advantage function are integrated to provide a new local policy optimization objective for each agent. Further, we prove that the safety constraints and the joint policy improvement can be guaranteed when updating the local policy with a sequential update scheme, in which the policy update only depends on its action and the state of its  $\kappa$ -hop neighbors for each agent.

### 3.1 Truncated advantage function estimator

For a standard safe MARL, the state-action value function (the definition can be seen in Appendix C.1) and advantage function of agent  $i$  yield that

$$Q_{\pi}^i(\mathbf{s}, \mathbf{a}^i) = \mathbb{E}_{\mathbf{a}^{-i} \sim \pi^{-i}} Q_{\pi}^i(\mathbf{s}, \mathbf{a}^{-i}, \mathbf{a}^i), \quad (6)$$

$$A_{\pi}^i(\mathbf{s}, \mathbf{a}^j, \mathbf{a}^i) = Q_{\pi}^{j,i}(\mathbf{s}, \mathbf{a}^j, \mathbf{a}^i) - Q_{\pi}^j(\mathbf{s}, \mathbf{a}^j). \quad (7)$$

where  $\mathbf{s}$  represents the global state,  $\mathbf{a}^{-i}$  represents the actions of all other agents, and  $Q_{\pi}^{j,i}(\mathbf{s}, \mathbf{a}^j, \mathbf{a}^i)$  represents the state-action value function of agent  $i$  and agent  $j$ . Then, updating agents' policies with a sequential update scheme [16], the multi-agent joint advantage function  $A_{\pi}(\mathbf{s}, \mathbf{a})$  can be written as a sum of sequentially unfolding multi-agent advantages of individual agents, as stated by the following lemma.

**Lemma 3.1.** (*Multi-agent advantage decomposition*). *For any action  $\mathbf{a}^i$ ,  $i \in \mathcal{N}$ , and the state  $\mathbf{s} \in \mathcal{S}$ , the following identity holds*

$$A_{\pi}(\mathbf{s}, \mathbf{a}) = \sum_{i=1}^n A_{\pi}^i(\mathbf{s}, \mathbf{a}^{-i}, \mathbf{a}^i). \quad (8)$$

Similar result to Lemma 3.1 can be seen in [10], and the proof is reported in Appendix C.2. Specifically, based on the multi-agent advantage decomposition in Lemma 3.1, the "surrogate" return is given as follows.

**Definition 3.2.** Let  $\pi$  be a joint policy,  $\bar{\pi}^{1:i-1}$  be some other joint policy of agents  $1 : i - 1$ , and  $\hat{\pi}^i$  be a policy of agent  $i$ . Then, the surrogate return can be defined as

$$L_{\pi}^{1:i}(\bar{\pi}^{1:i-1}, \hat{\pi}^i) \triangleq \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^{1:i-1} \sim \bar{\pi}^{1:i-1}, \mathbf{a}^i \sim \hat{\pi}^i} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)]. \quad (9)$$

Building on Lemma 3.1 and Definition 3.2, one can obtain

$$L_{\pi}^{1:i}(\bar{\pi}^{1:i-1}, \bar{\pi}^i) = \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} \left[ \sum_{h=1}^i A_{\pi}^h(\mathbf{s}, \mathbf{a}^{1:h-1}, \mathbf{a}^h) \right]. \quad (10)$$

Further, recalling Assumption 2.1 and Assumption 2.2, we can quantify the maximum information loss regarding the advantage function as stated by the following proposition.

**Proposition 3.3.** *For any agent  $i \in \mathcal{N}$ , let the parameters  $(\eta, \phi) = \left( \frac{\xi\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ . If Assumption 2.1 and Assumption 2.2 hold, for any  $\mathbf{z}_{\mathcal{N}_{\kappa}^i} = (\mathbf{s}_{\mathcal{N}_{\kappa}^i}, \mathbf{a}_{\mathcal{N}_{\kappa}^i}) \in \mathcal{S}_{\mathcal{N}_{\kappa}^i} \times \mathcal{A}_{\mathcal{N}_{\kappa}^i}$ , the exponential decay property of the advantage function holds, i.e., we have*

$$\sup_{\mathbf{z}_{\mathcal{N}_{\kappa}^i}, \mathbf{z}_{\mathcal{N}_{\kappa}^{-i}}, \mathbf{z}'_{\mathcal{N}_{\kappa}^{-i}}} \left| A^i(\mathbf{z}_{\mathcal{N}_{\kappa}^i}, \mathbf{z}_{\mathcal{N}_{\kappa}^{-i}}) - A^i(\mathbf{z}_{\mathcal{N}_{\kappa}^i}, \mathbf{z}'_{\mathcal{N}_{\kappa}^{-i}}) \right| \leq \eta\phi^{\kappa}. \quad (11)$$

Proposition 3.3 shows that when the transition dynamics and policies correlation satisfy the exponential correlation decay property, the advantage functions also have exponential decay dependence on the states and actions of the more distant agents. The proof of Proposition 3.3 is reported in Appendix C.3. In addition, based on this proposition, we can obtain the following corollary.

**Corollary 3.4.** *For any agent  $i \in \mathcal{N}$ , let the parameters  $(\eta', \phi) = \left( \frac{M^i\xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ ,  $M^i$  is a constant. If Proposition 3.3 holds, the exponential decay property of the surrogate return holds, i.e., we have*

$$\left| L_{\pi}^{1:i}(\bar{\pi}^{1:i-1}, \bar{\pi}^i) - L_{\pi_{\kappa}^i}^{1:i}(\bar{\pi}_{\kappa}^i) \right| \leq \eta'\phi^{\kappa}. \quad (12)$$

The proofs of Corollary 3.4 is reported in Appendix C.4.

Corollary 3.4 shows that the approximation error of  $L_{\pi_\kappa^i}^i(\bar{\pi}_\kappa^i)$  decreases exponentially with  $\kappa$  when the truncated advantage functions are bounded. The main advantage of using the estimator  $L_{\pi_\kappa^i}^i(\bar{\pi}_\kappa^i)$  lies in that every agent  $i$  only needs to know the action and state of its  $\kappa$ -hop neighbors, which can significantly reduce the communication burden and expand its application scenarios.

### 3.2 Scalable constrained policy optimization

With the Definition 3.2, we see that Lemma 3.1 allows for decomposing the joint surrogate return  $L_\pi(\bar{\pi}) \triangleq \mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a} \sim \bar{\pi}} [A_\pi(\mathbf{s}, \mathbf{a})]$  into a sum over surrogates of  $L_\pi^{1:i}(\bar{\pi}^{1:i-1}, \hat{\pi}^i)$ . Then, combining the rigorous bounds of the trust region method [17] and the bounds of the truncated advantage function, we can obtain the following proposition.

**Proposition 3.5.** *Let  $\pi$  and  $\bar{\pi}$  be joint policies. Let each agent  $i \in \mathcal{N}$  sequentially solves the following optimization problem:*

$$\bar{\pi}_\kappa^i = \arg \max_{\hat{\pi}_\kappa^i} \left( L_{\pi_\kappa^i}^i(\hat{\pi}_\kappa^i) - \eta' \phi^\kappa - \nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) \right), \quad (13)$$

where  $(\eta', \phi) = \left( \frac{M^i \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ ,  $\nu_\kappa^i = \frac{2\gamma \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i} |A_{\pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i)|}{(1-\gamma)^2}$ , and  $D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) = \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}} D_{\text{KL}}(\pi^i(\cdot | \mathbf{s}_{\mathcal{N}_\kappa^i}), \hat{\pi}^i(\cdot | \mathbf{s}_{\mathcal{N}_\kappa^i}))$ , then the resulting joint policy  $\bar{\pi}$  will improve the expected return, i.e.,

$$J(\bar{\pi}) - J(\pi) \geq \sum_{i=1}^N \left( L_{\pi_\kappa^i}^i(\hat{\pi}_\kappa^i) - \eta' \phi^\kappa - \nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) \right). \quad (14)$$

The proof of Proposition 3.5 is reported in Appendix C.5. Similarly, by generalizing the result about the surrogate return in Equation (12), we can derive how the expected costs change when the agents update their policies. Specifically, we provide the following corollary.

**Corollary 3.6.** *Let  $\pi$  and  $\bar{\pi}$  be joint policies. For any agent  $i \in \mathcal{N}$  and its cost index  $j \in \{1, \dots, m^i\}$ , the following inequality holds*

$$J_j^i(\bar{\pi}) \leq J_j^i(\pi) + L_{j, \pi_\kappa^i}^i(\bar{\pi}_\kappa^i) + \eta'' \phi^\kappa + \nu_{j, \kappa}^i \sum_{h=1}^{i-1} D_{\text{KL}}^{\max}(\pi_\kappa^h, \bar{\pi}_\kappa^h), \quad (15)$$

where  $L_{j, \pi_\kappa^i}^i(\bar{\pi}_\kappa^i) = \mathbb{E}_{\mathbf{s}_{\mathcal{N}_\kappa^i} \sim \rho_{\pi_\kappa^i}, \mathbf{a}^i \sim \bar{\pi}_\kappa^i} [A_{j, \pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i)]$ ,  $\nu_{j, \kappa}^i = \frac{2\gamma \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i} |A_{j, \pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i)|}{(1-\gamma)^2}$ ,  $(\eta'', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ , and  $M_j$  is a constant.

The proofs of Corollary 3.6 is reported in Appendix C.6.

From (14), we can derive that the lower bound for the difference between the new joint policy  $\bar{\pi}$  and the old joint policy  $\pi$  in terms of expected return can be decomposed into a cumulative sum of local surrogate TRPO policy objectives. From (15), we can derive the upper bound for the new joint policy  $\bar{\pi}$ , which can be used to restrict agents only to choose safe actions. Therefore, we use the objective, i.e., maximize the lower bound for the reward performance and minimize the upper bound for the safety constraints with a proper update size, as a surrogate for each agent. Then, we can obtain the following theorem.

**Theorem 3.7.** *The joint policy  $\pi$  has the monotonic improvement property,  $J(\bar{\pi}) \geq J(\pi)$ , as well as it satisfies the safety constraints,  $J_j^i(\bar{\pi}) \leq c_j^i$ , for any agent  $i \in \mathcal{N}$  and its cost index  $j \in \{1, \dots, m^i\}$ , when the policy is updated by following a sequential update scheme, that is, each agent sequentially solves the following optimization problem:*

$$\begin{aligned} \bar{\pi}_\kappa^i &= \arg \max_{\hat{\pi}_\kappa^i \in \bar{\Pi}_\kappa^i} \left( L_{\pi_\kappa^i}^i(\hat{\pi}_\kappa^i) - \eta' \phi^\kappa - \nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) \right), \\ \text{s.t. } &\{ \hat{\pi}_\kappa^i \in \bar{\Pi}_\kappa^i \mid D_{\text{KL}}^{\max}(\pi_\kappa^i, \hat{\pi}_\kappa^i) \leq \delta_\kappa^i, \text{ and} \\ &J_j^i(\pi_\kappa^i) + L_{j, \pi_\kappa^i}^i(\hat{\pi}_\kappa^i) + \eta'' \phi^\kappa + \nu_{j, \kappa}^i D_{\text{KL}}^{\max}(\pi_\kappa^i, \hat{\pi}_\kappa^i) \leq c_j^i - \nu_{j, \kappa}^i \sum_{h=1}^{i-1} D_{\text{KL}}^{\max}(\pi_\kappa^h, \hat{\pi}_\kappa^h) \}, \end{aligned} \quad (16)$$

$$\begin{aligned}
\text{where } \delta_\kappa^i &= \min \left\{ \min_{h \leq i-1} \min_{1 \leq j \leq m^h} \frac{\Xi_j^h - L_{j, \pi_\kappa^h}^h(\bar{\pi}_\kappa^h) - \eta'' \phi^\kappa}{\nu_{j, \kappa}^i}, \min_{h \geq i+1} \min_{1 \leq j \leq m^h} \frac{\Xi_j^h}{\nu_{j, \kappa}^i} \right\}, \\
\nu_\kappa^i &= \frac{2\gamma \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i} |A_{\pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i)|}{(1-\gamma)^2}, \nu_{j, \kappa}^i = \frac{2\gamma \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i} |A_{j, \pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i)|}{(1-\gamma)^2}, (\eta', \phi) = \\
&\left( \frac{M^i \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right), (\eta'', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right), \Xi_j^h = c_j^h - J_j^h(\pi_\kappa^h) - \\
&\nu_{j, \kappa}^i \sum_{l=1}^{i-1} D_{\text{KL}}^{\max}(\pi_\kappa^l, \hat{\pi}_\kappa^l).
\end{aligned}$$

The proof of Theorem 3.7 is reported in Appendix C.7. It assures that if one follows (16) to update policies, agents will not only explore safe policies independently; meanwhile, every new policy will be guaranteed to result in performance improvement. It is worth mentioning that these two properties hold only under the condition that the only policy update restriction, i.e.,  $\bar{\pi}_\kappa^i \in \bar{\Pi}_\kappa^i$ , is satisfied; this is due to the KL-penalty term in every agent's objective, i.e.,  $\nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i, \bar{\pi}_\kappa^i)$ , as well as the constraints on cost surrogates.

### 3.3 Algorithm

In this section, we focus on how to practically implement policy updates in Theorem 3.7 for each agent. Specifically, we parameterize each local policy  $\pi_{\theta_\kappa^i}^i$  by a neural network with parameter  $\theta_\kappa^i$ . At each policy update, every agent  $i$  maximizes its surrogate return subject to surrogate cost constraints and a form of expected KL-divergence constraint  $\tilde{D}_{\text{KL}}(\pi_\kappa^i, \bar{\pi}_\kappa^i) \leq \delta_\kappa^i$ , which avoids computing KL-divergence at every state. Then, we introduce a scalar variable  $\lambda^i$  for any agent  $i \in \mathcal{N}$  and convert the constrained optimization problem from (16) into a min-max optimization problem with Lagrangian multipliers by subsuming the cost constraints. As such, the new optimization problem for any agent  $i \in \mathcal{N}$  is as follows:

$$\begin{aligned}
&\max_{\theta_\kappa^i} \min_{\lambda_{1:m^i}^i \geq 0} \left[ \mathbb{E}_{\mathbf{s}_{\mathcal{N}_\kappa^i} \sim \rho_{\pi_{\theta_\kappa^i}^i}, \mathbf{a}^i \sim \pi_{\theta_\kappa^i}^i} \left[ A_{\pi_{\theta_\kappa^i}^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) \right] \right. \\
&\quad \left. - \sum_{u=1}^{m^i} \lambda_u^i \left( \mathbb{E}_{\mathbf{s}_{\mathcal{N}_\kappa^i} \sim \rho_{\pi_{\theta_\kappa^i}^i}, \mathbf{a}^i \sim \pi_{\theta_\kappa^i}^i} \left[ A_{u, \pi_{\theta_\kappa^i}^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) \right] + d_u^i \right) \right], \quad (17) \\
&\text{s.t. } \tilde{D}_{\text{KL}}(\pi_{\theta_\kappa^i}^i, \bar{\pi}_{\theta_\kappa^i}^i) \leq \delta_\kappa^i.
\end{aligned}$$

where  $\lambda_{1:m^i}^i$  is a scalar variable,  $\theta_\kappa^i$  is a parameter of neural network, and  $d_u^i$  is the cost-constraining value for agent  $i$ .

Further, denoting

$$A_{\pi_{\theta_\kappa^i}^i}^{i, (\lambda)}(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) = A_{\pi_{\theta_\kappa^i}^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) - \sum_{u=1}^{m^i} \lambda_u^i \left( A_{u, \pi_{\theta_\kappa^i}^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) + d_u^i \right), \quad (18)$$

then the optimization problem in (17) can be rewritten as

$$\max_{\theta_\kappa^i} \min_{\lambda_{1:m^i}^i \geq 0} \left[ \mathbb{E}_{\mathbf{s}_{\mathcal{N}_\kappa^i} \sim \rho_{\pi_{\theta_\kappa^i}^i}, \mathbf{a}^i \sim \pi_{\theta_\kappa^i}^i} \left[ A_{\pi_{\theta_\kappa^i}^i}^{i, (\lambda)}(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) \right] \right], \text{ s.t. } \tilde{D}_{\text{KL}}(\pi_{\theta_\kappa^i}^i, \bar{\pi}_{\theta_\kappa^i}^i) \leq \delta_\kappa^i. \quad (19)$$

To alleviate the complications caused by computing the KL-divergence constraint, we simplify it by adopting the PPO-clip objective [18], i.e., replacing the KL-divergence constraint with the clip operator and updating the policy parameter with first-order methods. The final optimization problem takes the form

$$\max_{\theta_\kappa^i} \min_{\lambda_{1:m^i}^i \geq 0} \mathbb{E}_{\mathbf{s}_{\mathcal{N}_\kappa^i} \sim \rho_{\pi_{\theta_\kappa^i}^i}, \mathbf{a}^i \sim \pi_{\theta_\kappa^i}^i} \left[ \min \left( \frac{\bar{\pi}_{\theta_\kappa^i}^i}{\pi_{\theta_\kappa^i}^i} A_{\pi_{\theta_\kappa^i}^i}^{i, (\lambda)}(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i), \left( \frac{\bar{\pi}_{\theta_\kappa^i}^i}{\pi_{\theta_\kappa^i}^i}, 1 \pm \epsilon \right) A_{\pi_{\theta_\kappa^i}^i}^{i, (\lambda)}(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i) \right) \right], \quad (20)$$

where the clip operator replaces the policy ratio with  $1 + \epsilon$ , or  $1 - \epsilon$ , depending on whether its value is below or above the threshold interval. As such, agent  $i$  can learn within its trust region by updating

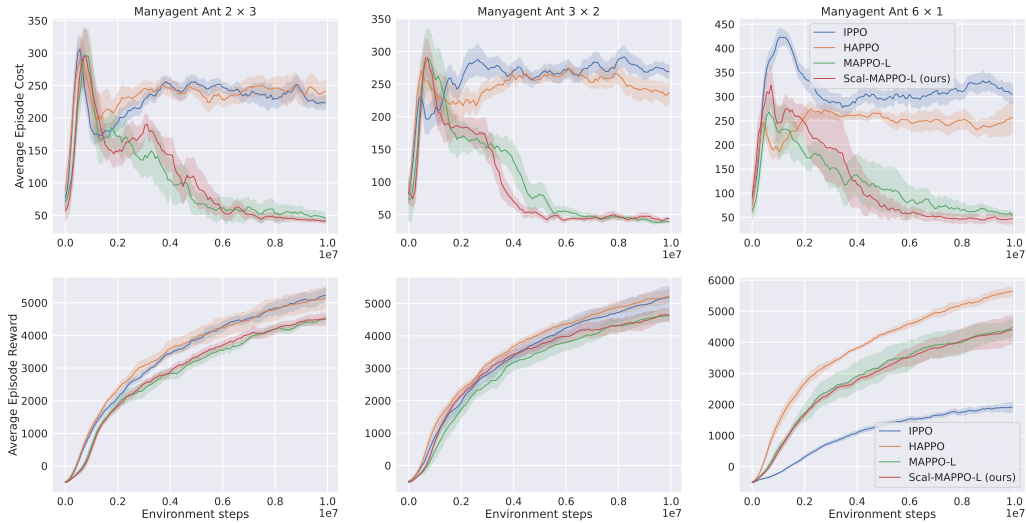


Figure 1: Performance comparisons in terms of cost and reward on three Safe ManyAgent Ant tasks. Each column subfigure represents a different task, and we plot the cost curves (the lower the better) in the upper row and the reward curves (the higher the better) in the bottom row for each task.

$\theta_{\kappa}^i$  to maximize Equation (20), which only depends on its action and the state of its  $\kappa$ -hop neighbors and can be computed analytically.

To summarize, we give a procedure for each agent  $i$ , name Scalable MAPPO-Lagrangian (Scal-MAPPO-L), and provide its pseudocode (Algorithm 1) in Appendix C.8. The algorithm has a simple idea that each agent independently optimizes the surrogate objective (20), which only depends on its action and the state of its  $\kappa$ -hop neighbors for each agent. In the actual execution, some approximations of the surrogate objective are employed, the same as the MAPPO-L [10]. Most of these approximations are traditional practices in RL, yet they may make it impossible for the practical algorithm to rigorously maintain the theoretical guarantees in Theorem 3.7.

## 4 Experiments

In this section, we evaluate our method via several numerical experiments. Our experiments aim to answer the following questions: First, how does the cost and reward performance of Scal-MAPPO-L compare with existing methods on challenging multi-agent safe tasks? Second, how does the different  $\kappa$  affect the performance of Scal-MAPPO-L, and could the advantage truncation effectively alleviate computational load?

### 4.1 Experimental setup

Safe MAMuJoCo [10] is an extension of MAMuJoCo [19], which preserves the agents, physics simulator, background environment, and reward function and comes with obstacles, like walls or pitfalls. To answer the first question, we compare our method against the other PPO family algorithms, i.e., IPPO [20], HAPPO [16], and MAPPO-L [10] and choose three games from Safe MAMuJoCo: Safe ManyAgent Ant task with 2 agents ( $2 \times 3$ ), 3 agents ( $3 \times 2$ ) and 6 agents ( $6 \times 1$ ) to evaluate their performance. Concerning the second question, we choose three games with different tasks and agent numbers from Safe MAMuJoCo: Safe ManyAgent Ant task with 6 agents ( $6 \times 1$ ), Safe Ant task with 8 agents ( $8 \times 1$ ), and Safe Coupled HalfCheetah task with 12 agents ( $12 \times 1$ ). We train Scal-MAPPO-L with the same network architecture and hyperparameters as the original MAPPO-L implementation. All reported results are averaged over three or more random seeds, and the curves are smooth over time.

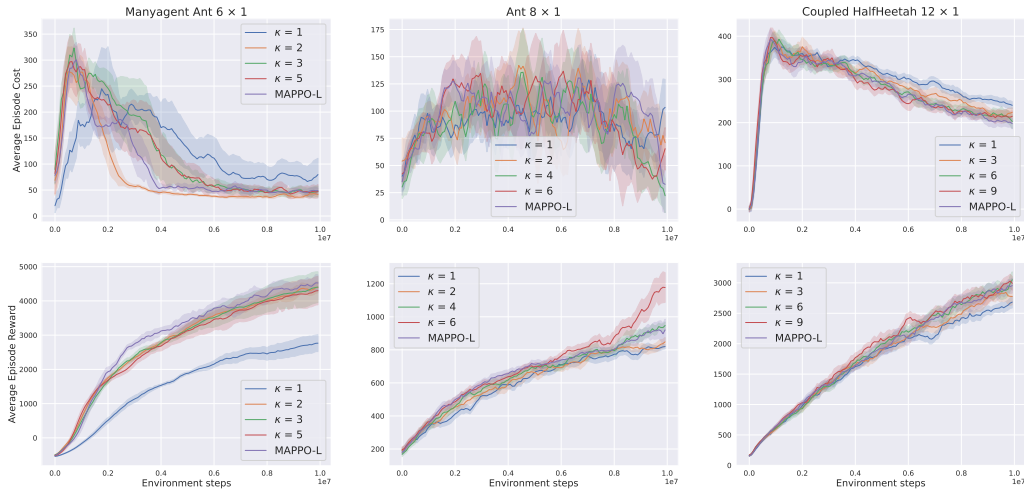


Figure 2: Performance comparisons in terms of cost and reward on Safe ManyAgent Ant task, Safe Ant task, and Safe Coupled HalfCheetah task. In each task, the performance of Scal-MAPPO-L with different  $\kappa$  and MAPPO-L are demonstrated.

## 4.2 Results

**Comparisons with baselines:** Figure 1 shows the cost and reward performance of Scal-MAPPO-L and other PPO family algorithms on three Safe ManyAgent Ant tasks, where each agent in Scal-MAPPO-L is set to access the state of about half of the agents by adjusting the value of  $\kappa$ . Specifically,  $\kappa = 1$  in Safe ManyAgent Ant ( $2 \times 3$ ),  $\kappa = 2$  in Safe ManyAgent Ant ( $3 \times 2$ ), and  $\kappa = 3$  in Safe ManyAgent Ant ( $6 \times 1$ ). From Figure 1, we can see that compared to IPPO and HAPPO, on all three tasks, both Scal-MAPPO-L and MAPPO-L have fewer constraint violations and good performance (in terms of reward), i.e., they keep their explorations within the feasible policy space and quickly learn to satisfy safety constraints, which show that the safe learning algorithm is effective. Moreover, it should be further pointed out that Scal-MAPPO-L only accesses half of the state information on all tasks; it exhibits almost identical performance and constraint violations with MAPPO-L (which accesses the global state). This means that the sensitivity of each agent to the states and actions perturbations of distant agents is minimal, and Scal-MAPPO-L is effective. More experimental results are in Appendix D.

**Performance with different  $\kappa$ :** Figure 2 shows the performance of Scal-MAPPO-L in different environments with varying values of  $\kappa$ , where MAPPO-L accesses the global state. We have noticed that the algorithm’s performance is consistently the lowest, and the cost is nearly the highest when  $\kappa = 1$ . However, when the truncation with  $\kappa \geq 3$ , i.e., each agent has access to the states of at least two neighbors, we can observe that the performance of Scal-MAPPO-L improves considerably and can approach or even outperform MAPPO-L in some environments, such as  $\kappa = 6$  in the Safe Ant task ( $8 \times 1$ ). This may be due to the fact that the impact of far-away agents’ states and actions on the agent’s decision is almost negligible in many cases. However, for algorithms with global communication, such as MAPPO-L, the difficulty of extracting useful information from many messages may lead to lower performance. Overall, these results underscore the efficiency of Scal-MAPPO-L since it employs a smaller communication radius that can significantly reduce the computation.

## 5 Related work

### 5.1 Safe RL

Safety is one of the bottlenecks preventing RL use in real-life applications, such as physical robotics [21], medical applications [22] and autonomous driving [23]. It has become a research hotspot in recent years and a growing number of safe RL approaches, such as primal-dual methods [24], formal methods [25], Lyapunov methods [26], Gaussian processes methods [27], and safety-augmented



methods [28], have been developed. However, when it comes to multi-agent systems, a great challenge is exacerbated by policy conflicts caused by multiple agents interacting within a shared environment and learning simultaneously. In other words, each agent has to not only satisfy its safety constraints but also consider the conflicts between its safety constraints and maximization reward as well as the safety constraints of others so that their joint behaviors have a safety guarantee. In order to address the above issue, CMIX [9] and MAPPO-L [10] have been proposed with the in-depth study of MARL. These algorithms follow the centralized training and decentralized execution (CTDE) framework [29, 30, 16], which learns the centralized value function by introducing the global state. Unfortunately, the global coupling arising from agents’ safety constraints and the exponential growth of the state-action space size make the usability in communication or computing resource-constrained systems and the scalability of these algorithms in larger multi-agent systems become a bottleneck, limiting their applicability. Recent works [11, 12] have provided some theoretical results to avoid these shortcomings. However, most of these methods fail to ensure both safety guarantee and joint policy improvement under a decentralized learning framework under a decentralized learning framework, which motivates us to investigate a new scalable and theoretically-justified safe MARL method.

## 5.2 Centralized training

In cooperative MARL settings, the training of agents can be broadly divided into two paradigms, namely centralized and decentralized [31]. The centralized training paradigm describes agent policies updated based on mutual information, which can be further differentiated into the centralized and decentralized execution framework. Centralized training and centralized execution (CTCE) utilize the centralized evaluator and executor to learn the joint policy of all agents [32, 18]. The obvious flaw is that its applicability is limited because its implementation requires the premise that instantaneous and unconstrained information exchange between agents. Recently, centralized training and decentralized execution (CTDE) has become the most popular framework [30, 20, 16, 10], since the fact that it addresses the non-stationarity issue with the centralized value function, and removes the dependency on global state and actions during execution. Many experiment results demonstrate state-of-the-art performance on challenging tasks, such as unit micromanagement in StarCraft II [33]. However, although this framework does not require agents to access the global state during execution, the reliance on the global state only during training still poses a significant barrier to real-world applications, especially in scenarios where communication and computational resources are constrained [34, 35].

## 5.3 Decentralized training

In a decentralized learning paradigm, each agent learns independently and accesses local observations rather than the global state; the idea is direct, comprehensible, and easy to realize in practice [36, 34]. There are two mainline research approaches concerning decentralized learning in the existing literature. One line of research pursues fully decentralized learning, such as independent Q-learning (IQL) [37, 38] and independent actor-critic (IAC) [39, 20], which make agents directly execute the single-agent Q-learning or actor-critic algorithm individually. Another line of research allows agents to establish rational local communication networks, such as setting certain distance or neighbor graphs [40, 41], which is also known as networked MARL. Communication networks expand agents’ perceptual capabilities and mitigate, to some extent, the decision conflicts or errors caused by partial observability. However, it is worth noting that each agent’s decision violates the stationary condition of the Markov Decision Process (MDP) in both lines of research, even though they achieve good experimental results on a collection of benchmark tasks. It poses a significant challenge to the convergence analysis of algorithms in the short term. Recently, motivated by good experiment performance, some studies have tried to provide theoretical support for these phenomena. To mention a few, Qu et al. [42] introduced the spatial correlation decay property into the field of MARL and carried out a series of fundamental results [15, 43, 12], which broadened the research avenues of scalable MARL. However, all of these studies mainly focus on (natural) policy gradient methods with average rewards or general utilities and have not yet been combined with trust region methods, which rigorously enable RL agents to learn monotonically improving policies. Furthermore, only recent research [12] considers both safety and scalability for MARL. Our results build upon the scalable MARL family of works [42, 15, 43, 12] and PPO-based (TRPO-based) MARL family of works [16, 10].

## 6 Conclusion

Safety is a tremendous challenge for MARL when applied to real-world scenarios. In this paper, we quantize the approximation errors arising from policy implementation and advantage truncation and then derive a novel lower bound for joint policy improvement and an upper bound for the safety constraints for every agent. Furthermore, we propose a novel scalable and theoretically justified multi-agent constrained policy optimization method that follows a sequential update scheme to optimize  $\kappa$ -hop policies. Finally, we introduce a practical constrained policy optimization algorithm called Scal-MAPPO-L and experimentally validate the effectiveness of the proposed algorithm on a collection of benchmark tasks.

## Acknowledgements

This work is supported by the National Key Research and Development Program of China (No.2020AAA0106100), the National Natural Science Project of China (Nos.62276160, 62376013), and the Basic Research Program of Shanxi Province (No.202203021211294).

## References

- [1] Lukas Brunke, Melissa Greeff, Adam W Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and Angela P Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement learning. *Annual Review of Control, Robotics, & Autonomous Systems*, 5:411–444, 2022.
- [2] Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*, 2022.
- [3] Wei Zhou, Dong Chen, Jun Yan, Zhaojian Li, Huilin Yin, and Wanchen Ge. Multi-agent reinforcement learning for cooperative lane changing of connected and autonomous vehicles in mixed traffic. *Autonomous Intelligent Systems*, 2(1):5, 2022.
- [4] Wenqi Cui, Jiayi Li, and Baosen Zhang. Decentralized safe reinforcement learning for inverter-based voltage control. *Electric Power Systems Research*, 211:108609, 2022.
- [5] Yu-Jia Chen, Deng-Kai Chang, and Cheng Zhang. Autonomous tracking using a swarm of uavs: A constrained multi-agent reinforcement learning approach. *IEEE Transactions on Vehicular Technology*, 69(11):13702–13717, 2020.
- [6] Kai-Chieh Hsu, Allen Z Ren, Duy P Nguyen, Anirudha Majumdar, and Jaime F Fisac. Sim-to-lab-to-real: Safe reinforcement learning with shielding and generalization guarantees. *Artificial Intelligence*, 314:103811, 2023.
- [7] Wenbo Zhang, Osbert Bastani, and Vijay Kumar. Mamps: Safe multi-agent reinforcement learning via model predictive shielding. *arXiv preprint arXiv:1910.12639*, 2019.
- [8] Daniel Melcer, Christopher Amato, and Stavros Tripakis. Shield decentralization for safe multi-agent reinforcement learning. In *NeurIPS*, 2022.
- [9] Chenyi Liu, Nan Geng, Vaneet Aggarwal, Tian Lan, Yuan Yang, and Mingwei Xu. Cmix: Deep multi-agent reinforcement learning with peak and average constraints. In *ECML-PKDD*, 2021.
- [10] Shangding Gu, Jakub Grudzien Kuba, Yuanpei Chen, Yali Du, Long Yang, Alois Knoll, and Yaodong Yang. Safe multi-agent reinforcement learning for multi-robot control. *Artificial Intelligence*, 319:103905, 2023.
- [11] Songtao Lu, Kaiqing Zhang, Tianyi Chen, Tamer Başar, and Lior Horesh. Decentralized policy gradient descent ascent for safe multi-agent reinforcement learning. In *AAAI*, 2021.
- [12] Donghao Ying, Yunkai Zhang, Yuhao Ding, Alec Koppel, and Javad Lavaei. Scalable primal-dual actor-critic method for safe multi-agent rl with general utilities. In *NuerIPS*, 2023.

- [13] Amir Dembo and Andrea Montanari. Gibbs measures and phase transitions on sparse random graphs. *arXiv preprint arXiv:0910.5460*, 2009.
- [14] David Gamarnik. Correlation decay method for decision, optimization, and inference in large-scale networks. In *Theory Driven by Influential Applications*, pages 108–121. 2013.
- [15] Guannan Qu, Yiheng Lin, Adam Wierman, and Na Li. Scalable multi-agent reinforcement learning for networked systems with average reward. In *NeurIPS*, 2020.
- [16] Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. Trust region policy optimisation in multi-agent reinforcement learning. In *ICLR*, 2022.
- [17] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *ICML*, 2015.
- [18] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [19] Bei Peng, Tabish Rashid, Christian Schroeder de Witt, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Boehmer, and Shimon Whiteson. Facmac: Factored multi-agent centralised policy gradients. In *NeurIPS*, 2021.
- [20] Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. In *NeurIPS*, 2022.
- [21] Javier García and Diogo Shafie. Teaching a humanoid robot to walk faster through safe reinforcement learning. *Engineering Applications of Artificial Intelligence*, 88:103360, 2020.
- [22] Shounak Datta, Yanjun Li, Matthew M Ruppert, Yuanfang Ren, Benjamin Shickel, Tezcan Ozrazgat-Baslanti, Parisa Rashidi, and Azra Bihorac. Reinforcement learning in surgery. *Surgery*, 170(1):329–332, 2021.
- [23] Shangding Gu, Guang Chen, Lijun Zhang, Jing Hou, Yingbai Hu, and Alois Knoll. Constrained reinforcement learning for vehicle motion planning with topological reachability analysis. *Robotics*, 11(4):81, 2022.
- [24] Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo R Jovanovic. Natural policy gradient primal-dual method for constrained markov decision processes. In *NeurIPS*, 2020.
- [25] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. In *NeurIPS*, 2018.
- [26] Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A lyapunov-based approach to safe reinforcement learning. In *NeurIPS*, 2018.
- [27] Yanan Sui, Alkis Gotovos, Joel Burdick, and Andreas Krause. Safe exploration for optimization with gaussian processes. In *ICML*, 2015.
- [28] Aivar Sootla, Alexander I Cowen-Rivers, Taher Jafferjee, Ziyang Wang, David H Mguni, Jun Wang, and Haitham Ammar. Sauté rl: Almost surely safe reinforcement learning using state augmentation. In *ICML*, 2022.
- [29] 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.
- [30] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *ICML*, 2018.
- [31] Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, pages 1–49, 2022.

- [32] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016.
- [33] S Whiteson, M Samvelyan, T Rashid, CS De Witt, G Farquhar, N Nardelli, TGJ Rudner, CM Hung, PHS Torr, and J Foerster. The starcraft multi-agent challenge. In *AAMAS*, 2019.
- [34] Wei Du and Shifei Ding. A survey on multi-agent deep reinforcement learning: from the perspective of challenges and applications. *Artificial Intelligence Review*, 54:3215–3238, 2021.
- [35] Afshin Oroojlooy and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *Applied Intelligence*, 53(11):13677–13722, 2023.
- [36] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Studies in Systems, Decision and Control*, pages 321–384, 2021.
- [37] Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *ICML*, 1993.
- [38] Ardi Tampuu, Tabet Matiisen, Dorian Kodelja, Ilya Kuzovkin, Kristjan Korjus, Juhan Aru, Jaan Aru, and Raul Vicente. Multiagent cooperation and competition with deep reinforcement learning. *arXiv preprint arXiv:1511.08779*, 2015.
- [39] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- [40] Jiechuan Jiang, Chen Dun, Tiejun Huang, and Zongqing Lu. Graph convolutional reinforcement learning. In *ICLR*, 2019.
- [41] Tianshu Chu, Sandeep Chinchali, and Sachin Katti. Multi-agent reinforcement learning for networked system control. In *ICLR*, 2019.
- [42] Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning for multi-agent networked systems. *arXiv preprint arXiv:1912.02906*, 2019.
- [43] Yiheng Lin, Guannan Qu, Longbo Huang, and Adam Wierman. Multi-agent reinforcement learning in stochastic networked systems. In *NeurIPS*, 2021.
- [44] Hans-Otto Georgii. *Gibbs measures and phase transitions*. Walter de Gruyter GmbH & Co. KG, Berlin, 2011.
- [45] Alison L Gibbs and Francis Edward Su. On choosing and bounding probability metrics. *International statistical review*, 70(3):419–435, 2002.
- [46] Wenjun Mei, Shadi Mohagheghi, Sandro Zampieri, and Francesco Bullo. On the dynamics of deterministic epidemic propagation over networks. *Annual Reviews in Control*, 44:116–128, 2017.
- [47] Alessandro Zocca. Temporal starvation in multi-channel csma networks: an analytical framework. *ACM SIGMETRICS Performance Evaluation Review*, 46(3):52–53, 2019.
- [48] Guannan Qu and Na Li. Exploiting fast decaying and locality in multi-agent mdp with tree dependence structure. In *CDC*, 2019.
- [49] Haotian Gu, Xin Guo, Xiaoli Wei, and Renyuan Xu. Mean-field controls with q-learning for cooperative marl: convergence and complexity analysis. *SIAM Journal on Mathematics of Data Science*, 3(4):1168–1196, 2021.
- [50] Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multi-agent reinforcement learning. In *NeurIPS*, 2018.
- [51] Changxi Zhu, Mehdi Dastani, and Shihan Wang. A survey of multi-agent reinforcement learning with communication. *arXiv preprint arXiv:2203.08975*, 2022.

- [52] Junjie Sheng, Xiangfeng Wang, Bo Jin, Junchi Yan, Wenhao Li, Tsung-Hui Chang, Jun Wang, and Hongyuan Zha. Learning structured communication for multi-agent reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 36(2):50, 2022.

## A Preliminary lemmas

Before proving propositions, corollaries, and theorems, we need a series of intermediate results as a foundation. Results similar to Lemmas A.1 and A.2 can be found in Chapter 8 of [44], Lemma A.3 is an extension of results from [15], Lemma A.4 is an extension of Lemma A.2 found in [12]. We state these lemmas and provide the corresponding proofs for completeness as follows.

**Lemma A.1.** *Let  $f : \mathcal{S} \rightarrow [m, M]$ , where  $\mathcal{S} = \times_{i \in \mathcal{N}} \mathcal{S}^i$  and  $m, M \in \mathbb{R}$ . For every  $i \in \mathcal{N}$ , let  $\mu^i$  and  $\nu^i$  be two distributions on  $\mathcal{S}^i$ . Let  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  be the respective product distributions. Let  $\delta^i(f(\mathbf{s})) = \sup_{\mathbf{s}^i, \mathbf{s}^{-i}, \mathbf{s}'^i} |f(\mathbf{s}^i, \mathbf{s}^{-i}) - f(\mathbf{s}'^i, \mathbf{s}^{-i})|$ . Then, one have*

$$|\mathbb{E}_{\mathbf{s} \sim \boldsymbol{\mu}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim \boldsymbol{\nu}} f(\mathbf{s})| \leq \sum_{i \in \mathcal{N}} D_{\text{TV}}(\mu^i, \nu^i) \delta^i(f). \quad (21)$$

*Proof.* We prove Lemma A.1 by induction. Note that

$$D_{\text{TV}}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \frac{1}{2} \max_{|h| \leq 1} |\mathbb{E}_{\boldsymbol{\mu}}(h) - \mathbb{E}_{\boldsymbol{\nu}}(h)|$$

is an equivalent formulation of the total variation distance [45].

For  $|\mathcal{N}| = 1$ , one have

$$\begin{aligned} & |\mathbb{E}_{\mu^1}(f) - \mathbb{E}_{\nu^1}(f)| \\ &= \left| \mathbb{E}_{\mu^1} \left( f - \frac{M+m}{2} \right) - \mathbb{E}_{\nu^1} \left( f - \frac{M+m}{2} \right) \right| \\ &= \frac{M-m}{2} \left| \mathbb{E}_{\mu^1} \left( \frac{2f}{M-m} - \frac{M+m}{M-m} \right) - \mathbb{E}_{\nu^1} \left( \frac{2f}{M-m} - \frac{M+m}{M-m} \right) \right| \\ &\leq \frac{M-m}{2} \max_{|h| \leq 1} |\mathbb{E}_{\mu^1}(h) - \mathbb{E}_{\nu^1}(h)| \\ &= D_{\text{TV}}(\mu^1, \nu^1) \delta^1(f). \end{aligned}$$

As induction assumption, assume that Lemma A.1 holds for  $|\mathcal{N}| > 1$ . Then, one have

$$\begin{aligned} & |\mathbb{E}_{\mathbf{s} \sim \boldsymbol{\mu}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim \boldsymbol{\nu}} f(\mathbf{s})| \\ &= |\mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\mu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^1 \sim \nu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| \\ &= |\mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\mu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s}) \\ &\quad + \mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^1 \sim \nu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| \\ &\leq |\mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\mu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| \\ &\quad + |\mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^1 \sim \nu^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| \\ &\leq \mathbb{E}_{\mathbf{s}^1 \sim \mu^1} |\mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\mu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| + |\mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \tilde{f}(\mathbf{s}^1) - \mathbb{E}_{\mathbf{s}^1 \sim \nu^1} \tilde{f}(\mathbf{s}^1)|, \end{aligned}$$

where  $\tilde{f}(\mathbf{s}^1) = \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})$ .

By induction assumption, one have

$$\begin{aligned} |\mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\mu}^{2:n}} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s})| &\leq \sum_{i \neq 1 \in \mathcal{N}} D_{\text{TV}}(\mu^i, \nu^i) \delta^i(f(\mathbf{s}^1, \cdot)) \\ &\leq \sum_{i \neq 1 \in \mathcal{N}} D_{\text{TV}}(\mu^i, \nu^i) \delta^i(f). \end{aligned}$$

Since

$$\begin{aligned} \delta^1(\tilde{f}) &= \sup_{\mathbf{s}^1, \mathbf{s}'^1} |\mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s}^1, \mathbf{s}^{2:n}) - \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} f(\mathbf{s}'^1, \mathbf{s}^{2:n})| \\ &\leq \sup_{\mathbf{s}^1, \mathbf{s}'^1} \mathbb{E}_{\mathbf{s}^{2:n} \sim \boldsymbol{\nu}^{2:n}} |f(\mathbf{s}^1, \mathbf{s}^{2:n}) - f(\mathbf{s}'^1, \mathbf{s}^{2:n})| \\ &\leq \sup_{\mathbf{s}^1, \mathbf{s}'^1, \mathbf{s}^{2:n}} |f(\mathbf{s}^1, \mathbf{s}^{2:n}) - f(\mathbf{s}'^1, \mathbf{s}^{2:n})| \\ &= \delta^1(f), \end{aligned}$$

one have

$$\begin{aligned}
& |\mathbb{E}_{\mathbf{s} \sim \mu} f(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim \nu} f(\mathbf{s})| \\
& \leq \mathbb{E}_{\mathbf{s}^1 \sim \mu^1} \sum_{i \neq 1 \in \mathcal{N}} D_{\text{TV}}(\mu^i, \nu^i) \delta^i(f) + D_{\text{TV}}(\mu^1, \nu^1) \delta^1(f) \\
& \leq \sum_{i \in \mathcal{N}} D_{\text{TV}}(\mu^i, \nu^i) \delta^i(f),
\end{aligned}$$

which concludes the induction.  $\square$

**Lemma A.2.** Consider a Markov Chain with state  $\mathbf{s} \in \mathcal{S}$ , where  $\mathcal{S} = \times_{i \in \mathcal{N}} \mathcal{S}^i$ , and  $\mathcal{N} = \{1, \dots, n\}$  is the set of agents. Suppose its transition probability factorizes as

$$P(\mathbf{s}_{t+1} | \mathbf{s}_t) = \prod_{i \in \mathcal{N}} P^i(\mathbf{s}_{t+1}^i | \mathbf{s}_t).$$

Let  $W \in \mathbb{R}^{n \times n}$  be a matrix whose elements respect the condition

$$W^{ij} \geq \sup_{\mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j} D_{\text{TV}}(P^i(\cdot | \mathbf{s}^j, \mathbf{s}^{-j}), P^i(\cdot | \mathbf{s}'^j, \mathbf{s}^{-j})).$$

If  $\sum_{j \in \mathcal{J}} e^{\beta d(j,i)} W^{ij} \leq \zeta$ ,  $\mathcal{J} \subseteq \mathcal{N}$ , then one have

$$\sup_{\mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j} D_{\text{TV}}(P^i(\cdot | \mathbf{s}^J, \mathbf{s}^{-J}), P^i(\cdot | \mathbf{s}'^J, \mathbf{s}^{-J})) \leq \sum_{j \in \mathcal{J}} W^{ij}, \quad (22)$$

and

$$\sup_{\mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j} D_{\text{TV}}(P^i(\cdot | \mathbf{s}^J, \mathbf{s}^{-J}), P^i(\cdot | \mathbf{s}'^J, \mathbf{s}^{-J})) \leq \zeta e^{-\beta d(\mathcal{J}, i)}, \quad (23)$$

where  $d(\mathcal{J}, i) = \min_{j \in \mathcal{J}} d(j, i)$ .

*Proof.* We prove the first claim of Lemma A.2. The first claim clearly holds if  $|\mathcal{J}| = 1$ . As induction assumption, assume that the first claim holds for a set  $\mathcal{J}$ . Then, it holds for  $\mathcal{J}' = \mathcal{J} + \{k\}$

$$\begin{aligned}
& \sup_{\mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j} D_{\text{TV}}(P^i(\cdot | \mathbf{s}^{J'}, \mathbf{s}^{-J'}), P^i(\cdot | \mathbf{s}'^{J'}, \mathbf{s}^{-J'})) \\
& = \sup_{\substack{A \subseteq \mathcal{S}^i \\ \mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j}} \left| P^i(A | \mathbf{s}^{J'}, \mathbf{s}^{-J'}), P^i(A | \mathbf{s}'^{J'}, \mathbf{s}^{-J'}) \right| \\
& \leq \sup_{\substack{A \subseteq \mathcal{S}^i \\ \mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j}} \left| P^i(A | \mathbf{s}^{J'}, \mathbf{s}^{-J'}), P^i(A | \mathbf{s}'^J, \mathbf{s}^{-J}) \right| \\
& \quad + \sup_{\substack{A \subseteq \mathcal{S}^i \\ \mathbf{s}^j, \mathbf{s}^{-j}, \mathbf{s}'^j}} \left| P^i(A | \mathbf{s}'^J, \mathbf{s}^{-J}), P^i(A | \mathbf{s}'^{J'}, \mathbf{s}^{-J'}) \right| \\
& \leq \sum_{j \in \mathcal{J}} W^{ij} + W^{ik} \\
& = \sum_{j \in \mathcal{J}'} W^{ij}.
\end{aligned}$$

The second claim follows immediately, since

$$e^{\beta d(\mathcal{J}, i)} \sum_{j \in \mathcal{J}} W^{ij} \leq \sum_{j \in \mathcal{J}} e^{\beta d(j,i)} W^{ij} \leq \sum_{j \in \mathcal{N}} e^{\beta d(j,i)} W^{ij} \leq \zeta,$$

and

$$\sum_{j \in \mathcal{J}} W^{ij} \leq \zeta e^{-\beta d(\mathcal{J}, i)}.$$

$\square$

**Lemma A.3.** Consider the setting of Lemma A.2. For a generic value of  $\kappa$ , denote by  $\rho_t$  and  $\tilde{\rho}_t$  the distribution of  $\mathbf{s}_t$  with starting state, respectively,  $\mathbf{s} = (s_{\mathcal{N}_\kappa^i}, s_{\mathcal{N}_\kappa^{-i}})$  and  $\tilde{\mathbf{s}} = (s_{\mathcal{N}_\kappa^i}, \tilde{s}_{\mathcal{N}_\kappa^{-i}})$ . Then, if  $\sum_{j \in \mathcal{N}} e^{\beta d(j,i)} W^{ij} \leq \zeta$ , we have that  $D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i) \leq \zeta^t e^{-\beta \kappa}, \forall i \in \mathcal{N}$ .

*Proof.* We prove Lemma A.3 by induction. The case where  $t = 1$  follows from Lemma A.2. As induction assumption, assume that Lemma A.3 holds for  $t$ . Then, one have

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{s} \sim \rho_{t+1}} \mathbf{1}_A(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim \tilde{\rho}_{t+1}} \mathbf{1}_A(\mathbf{s}) \right| \\ &= \left| \mathbb{E}_{\mathbf{s} \sim \rho_t} E_{\mathbf{s} \sim P^i(\cdot|\mathbf{s})} \mathbf{1}_A(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim \tilde{\rho}_t} E_{\mathbf{s} \sim P^i(\cdot|\mathbf{s})} \mathbf{1}_A(\mathbf{s}) \right| \\ &\leq \sum_{j \in \mathcal{N}} D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i) \delta^j (E_{\mathbf{s} \sim P^i(\cdot|\cdot)} \mathbf{1}_A(\mathbf{s})) \\ &\leq \sum_{j \in \mathcal{N}} D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i) W^{ij} \\ &= \zeta^t e^{-\beta \kappa} \sum_{j \in \mathcal{N}} e^{\beta d(j,i)} W^{ij} \\ &\leq \zeta^{t+1} e^{-\beta \kappa}, \end{aligned}$$

where we used Lemma A.1 in the first inequality.  $\square$

**Lemma A.4.** Consider the setting of Lemma A.2. Let  $P^t(\mathbf{s}' | \mathbf{s}) = P(\mathbf{s}_t = \mathbf{s}' | \mathbf{s}_0 = \mathbf{s})$  and

$$\delta^j P^{i,t} = \sup_{s^j, s^{-j}, s'^j} D_{\text{TV}}(P^{i,t}(\cdot | s^j, \mathbf{s}^{-j}), P^{i,t}(\cdot | s'^j, \mathbf{s}^{-j})).$$

If  $\sum_{j \in \mathcal{N}} e^{\beta d(i,j)} W^{ij} \leq \zeta$ , we have

$$\sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \delta^j P^{i,t} \leq \zeta^t, \forall i \in \mathcal{N}. \quad (24)$$

*Proof.* We prove Lemma A.4 by induction. The claim holds for  $t = 1$ ,

$$\sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \delta^j P^{i,t} = \sum_{j \in \mathcal{N}} e^{\beta d(i,j)} W^{ij} \leq \zeta.$$

As induction assumption, we assume that the claim holds for  $t$ . Then, using Lemma A.1,

$$\begin{aligned} \delta^j P^{i,t+1} &= \sup_{\substack{A \subseteq \mathcal{S}^i \\ s^j, s^{-j}, s'^j}} \left| \mathbb{E}_{\mathbf{s} \sim P^{i,t+1}(\cdot|s^j, \mathbf{s}^{-j})} \mathbf{1}_A(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim P^{i,t+1}(\cdot|s'^j, \mathbf{s}^{-j})} \mathbf{1}_A(\mathbf{s}) \right| \\ &= \sup_{\substack{A \subseteq \mathcal{S}^i \\ s^j, s^{-j}, s'^j}} \left| \mathbb{E}_{\mathbf{s} \sim P^t(\cdot|s^j, \mathbf{s}^{-j})} E_{\mathbf{s} \sim P^i(\cdot|\mathbf{s})} \mathbf{1}_A(\mathbf{s}) - \mathbb{E}_{\mathbf{s} \sim P^t(\cdot|s'^j, \mathbf{s}^{-j})} E_{\mathbf{s} \sim P^i(\cdot|\mathbf{s})} \mathbf{1}_A(\mathbf{s}) \right| \\ &\leq \sup_{s^j, s^{-j}, s'^j} \sum_{k \in \mathcal{N}} D_{\text{TV}}(P^{k,t}(\cdot | s^j, \mathbf{s}^{-j}), P^{k,t}(\cdot | s'^j, \mathbf{s}^{-j})) \delta^j (E_{\mathbf{s} \sim P^i(\cdot|\cdot)} \mathbf{1}_A(\mathbf{s})) \\ &\leq \sum_{k \in \mathcal{N}} \delta^j P^{k,t} W^{ik}, \end{aligned}$$

and using the inverse triangle inequality,

$$\begin{aligned} \sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \delta^j P^{i,t+1} &\leq \sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \sum_{k \in \mathcal{N}} \delta^j P^{k,t} W^{ik} \\ &\leq \sum_{k \in \mathcal{N}} e^{\beta d(i,k)} W^{ik} \sum_{j \in \mathcal{N}} e^{\beta(d(i,j) - d(i,k))} \delta^j P^{k,t} \\ &\leq \sum_{k \in \mathcal{N}} e^{\beta d(i,k)} W^{ik} \sum_{j \in \mathcal{N}} e^{\beta d(k,j)} \delta^j P^{k,t} \\ &\leq \zeta^{t+1}, \end{aligned}$$

which concludes the induction.  $\square$



## B Supplementary materials for Section 2

### B.1 Spatial correlation decay

Exponential decay property [13, 14], also known as spatial correlation decay, is a powerful property associated with local interactions, which says that the impact of agents on each other decays exponentially in their graph distance. Over the past decades, many researchers have utilized spatial correlation property to design scalable, distributed algorithms for optimization and control problems in scenarios such as epidemics [46] and wireless communication [47]. Inspired by the studies mentioned above, a recent line of work [48] has formally considered spatial decay of correlation assumptions and proposes a method that finds nearly optimal local policies. An application [49] with the same principles adopts the setting of mean-field MARL [50], which proposes an actor-critic algorithm with global convergence. However, unlike the mean-field setting, which requires an agent's transition scheme to be only affected by the mean effect from its neighbors and effective only when agents are homogeneous, we allow each agent to have different transition probabilities and local policies.

### B.2 Regarding Assumptions 2.1 - 2.2

Assumption 2.1 portrays a common phenomenon: the transition dynamic of each agent is exponentially less sensitive to perturbations of the states and actions of more distant agents. This is commonly seen in scenarios involving wireless communication, epidemics, traffic, and so on [46, 47]. Assumption 2.2 imposes a design constraint for the policy class that encodes a weaker correlation decay property than the assumptions on the nature of Assumption 2.1. Moreover, Assumption 2.2 reveals how much information is lost compared with access to the global state and allows us to consider a policy class with the necessary properties for the optimal policy under Assumption 2.1. Below, we use a mathematical example to illustrate the relationship between the two assumptions.

**Mathematical example:** Firstly, we start from Assumption 2.1, letting  $\tilde{\kappa} = \max_{i,j \in \mathcal{N}} d(i, j)$  be the maximum distance between agent  $i$  and agent  $j$ . Define a set of differentiable functions  $\{f_\kappa : \mathcal{S}_{\mathcal{N}_\kappa^i} \times \mathcal{A}^i \rightarrow \mathcal{K} \mid 0 \leq \kappa \leq \tilde{\kappa}\}$ , where  $\mathcal{K} \subset [-K, K]$ ,  $K > 0$ , and a set of parameters  $\{\alpha_\kappa \geq 0 \mid 0 \leq \kappa \leq \tilde{\kappa}\}$ . Then, for each agent  $i$ , one have

$$f^i(\mathbf{s}, \mathbf{a}^i) = \sum_{\kappa=0}^{\tilde{\kappa}} \alpha_\kappa f_\kappa^i(\mathbf{s}_{\mathcal{N}_\kappa^i}, \mathbf{a}^i),$$

$$\pi^i(\mathbf{a} \mid \mathbf{s}) = \frac{\exp(f^i(\mathbf{s}, \mathbf{a}))}{\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\mathbf{s}, \mathbf{a}'))}.$$

By tuning the parameters  $\alpha_\kappa$ , we can make any policy belonging to this policy class respect Assumptions 2.2, as we show in the following. Let  $\kappa \in \{0, \dots, \tilde{\kappa}\}$ ,  $\mathbf{s}, \tilde{\mathbf{s}} \in \mathcal{S}$  be such that  $\mathbf{s}_{\mathcal{N}_\kappa^i} = \tilde{\mathbf{s}}_{\mathcal{N}_\kappa^i}$ , then one have

$$\begin{aligned} & \|\pi^i(\cdot \mid \mathbf{s}) - \pi^i(\cdot \mid \tilde{\mathbf{s}})\|_1 \\ &= \sum_{\mathbf{a} \in \mathcal{A}^i} |\pi^i(\mathbf{a} \mid \mathbf{s}) - \pi^i(\mathbf{a} \mid \tilde{\mathbf{s}})| \\ &= \sum_{\mathbf{a} \in \mathcal{A}^i} \left| \frac{\exp(f^i(\mathbf{s}, \mathbf{a}))}{\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\mathbf{s}, \mathbf{a}'))} - \frac{\exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}))}{\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}'))} \right| \\ &= \frac{\sum_{\mathbf{a} \in \mathcal{A}^i} |\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\mathbf{s}, \mathbf{a})) \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}')) - \sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a})) \exp(f^i(\mathbf{s}, \mathbf{a}'))|}{\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\mathbf{s}, \mathbf{a}')) \sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}'))} \\ &\leq \frac{\sum_{\mathbf{a} \in \mathcal{A}^i} \sum_{\mathbf{a}' \in \mathcal{A}^i} |\exp(f^i(\mathbf{s}, \mathbf{a})) \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}')) - \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a})) \exp(f^i(\mathbf{s}, \mathbf{a}'))|}{\sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\mathbf{s}, \mathbf{a}')) \sum_{\mathbf{a}' \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}'))} \\ &\leq \frac{\sum_{\mathbf{a} \in \mathcal{A}^i} |\exp(f^i(\tilde{\mathbf{s}}, \mathbf{a})) - \exp(f^i(\mathbf{s}, \mathbf{a}))|}{\sum_{\mathbf{a} \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}))} \\ &\leq \frac{\sum_{\mathbf{a} \in \mathcal{A}^i} |f^i(\tilde{\mathbf{s}}, \mathbf{a}) - f^i(\mathbf{s}, \mathbf{a})| \exp(\sup_{\mathbf{s}' \in \{\mathbf{s}, \tilde{\mathbf{s}}\}} f^i(\mathbf{s}', \mathbf{a}))}{\sum_{\mathbf{a} \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}))} \end{aligned}$$

$$\begin{aligned}
&\leq e^{2K(\tilde{\kappa}-\kappa)} \frac{\sum_{\mathbf{a} \in \mathcal{A}^i} |f^i(\tilde{\mathbf{s}}, \mathbf{a}) - f^i(\mathbf{s}, \mathbf{a})| \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}))}{\sum_{\mathbf{a} \in \mathcal{A}^i} \exp(f^i(\tilde{\mathbf{s}}, \mathbf{a}))} \\
&\leq e^{2K(\tilde{\kappa}-\kappa)} \mathbb{E}_{\pi^i} \left| \sum_{\kappa'=\kappa+1}^{\tilde{\kappa}} \alpha_{\kappa'} \left( f_{\kappa'}^i(\tilde{\mathcal{N}}_{\kappa'}^i, \mathbf{a}) - f_{\kappa'}^i(\mathcal{N}_{\kappa'}^i, \mathbf{a}) \right) \right| \\
&\leq e^{2K(\tilde{\kappa}-\kappa)} \sum_{\kappa'=\kappa+1}^{\tilde{\kappa}} \alpha_{\kappa'} \mathbb{E}_{\pi_{\kappa'}^i} \left| \left( f_{\kappa'}^i(\tilde{\mathcal{N}}_{\kappa'}^i, \mathbf{a}) - f_{\kappa'}^i(\mathcal{N}_{\kappa'}^i, \mathbf{a}) \right) \right| \\
&\leq 2K e^{2K(\tilde{\kappa}-\kappa)} \sum_{\kappa'=\kappa+1}^{\tilde{\kappa}} \alpha_{\kappa'}.
\end{aligned}$$

Denote that  $(\xi, \beta) = \left( 2K e^{2K\tilde{\kappa}} \sum_{\kappa'=\kappa+1}^{\tilde{\kappa}} \alpha_{\kappa'}, 2K \right)$ , and setting the parameters  $\{\alpha_{\kappa'}\}_{\kappa' \in \{\kappa+1, \dots, \tilde{\kappa}\}}$  small enough ensures that the policy respects Assumption 2.2.

*Remark B.1.* The mathematical example illustrates the relationship between the Assumptions 2.1 and 2.2. It is evident from this mathematical example that Assumption 2.2 necessarily holds when Assumption 2.1 holds and the parameters  $\xi$  and  $\beta$  satisfy certain conditions. However, for the sake of more concise presentation, we treat it as a separate assumption.

*Remark B.2.* When Assumption 2.1 holds, the numerical example can provide a reference basis for selecting the values of the parameters in Assumption 2.2. However, accurately determining the spatial decay of correlation for the dynamics remains a challenging engineering task. In this paper, we empirically adopt conservative values.

*Remark B.3.* Assumption 2.2 implies that multi-agent environments must satisfy the requirement that the impact from far-away agents' states and actions is almost negligible for the agent's decision; in other words, an action of an agent has an instantaneous effect on the system only locally. We believe that this formulation realistically describes most multi-agent interactions in the real-world. Take multi-vehicle transportation as an example. For a vehicle traveling on the road, a far-away vehicle taking different actions or being in a different state will affect itself shortly thereafter, but the impact on the current policy is minimal. More examples are seen in wireless communication, epidemics, traffic, and other scenarios [46, 47].

*Remark B.4.* It is worth noting that the assumption of spatial correlation decay is not in direct conflict with well-known phenomena, e.g., Butterfly Effect, since two seemingly unrelated things can also have a significant impact on each other, generally occurring in different time domains.

## C Supplementary materials for Section 3

### C.1 Basic definitions

Regarding the state value function and the state-action value function, we give the following definitions.

*Definition C.1.* We define the state value function and the state-action value function in terms of reward as

$$V_{\pi}(\mathbf{s}) \triangleq \mathbb{E}_{\mathbf{a} \sim \pi} [Q_{\pi}(\mathbf{s}, \mathbf{a})], \quad (25)$$

$$Q_{\pi}(\mathbf{s}, \mathbf{a}) \triangleq \mathbb{E}_{\mathbf{s}_{1:\infty} \sim p, \mathbf{a}_{1:\infty} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{R}(\mathbf{s}_t, \mathbf{a}_t) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right]. \quad (26)$$

Based on C.1, one can expand to derive

$$V_{\pi}^i(\mathbf{s}) = \mathbb{E}_{\mathbf{a}^{-i} \sim \pi^{-i}} [Q_{\pi}^i(\mathbf{s}, \mathbf{a}^i)], \quad (27)$$

$$Q_{\pi}^i(\mathbf{s}, \mathbf{a}^i) = \mathbb{E}_{\mathbf{a}^{-i} \sim \pi^{-i}, \mathbf{s}_{1:\infty} \sim p, \mathbf{a}_{1:\infty} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{R}(\mathbf{s}_t, \mathbf{a}_t^i) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right]. \quad (28)$$

*Definition C.2.* We define the  $j$ th state cost value function and state-action cost value function for agent  $i$  as follows

$$V_{j,\pi}^i(\mathbf{s}) \triangleq \mathbb{E}_{\mathbf{s}_{1:\infty} \sim \mathbf{p}, \mathbf{a}_{1:\infty} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t C_j^i(\mathbf{s}_t, \mathbf{a}_t^i) \mid \mathbf{s}_0 = \mathbf{s} \right], \quad (29)$$

$$Q_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i) \triangleq \mathbb{E}_{\mathbf{a}^{-i} \sim \pi^{-i}, \mathbf{s}_{1:\infty} \sim \mathbf{p}, \mathbf{a}_{1:\infty} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t C_j^i(\mathbf{s}_t, \mathbf{a}_t^i) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right]. \quad (30)$$

## C.2 The proof of Lemma 3.1

*Proof.* We write the multi-agent advantage function as in its definition, and then expand it in a telescoping sum.

$$\begin{aligned} A_{\pi}(\mathbf{s}, \mathbf{a}) &= Q_{\pi}(\mathbf{s}, \mathbf{a}) - V_{\pi}(\mathbf{s}) \\ &= \sum_{i=1}^n [Q_{\pi}^{1:i}(\mathbf{s}, \mathbf{a}^{1:i}) - Q_{\pi}^{1:i-1}(\mathbf{s}, \mathbf{a}^{1:i-1})] \\ &= \sum_{i=1}^n A_{\pi}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i). \end{aligned}$$

□

## C.3 The proof of Proposition 3.3

*Proof.* Let  $\mathbf{s}, \tilde{\mathbf{s}} \in \mathcal{S}$ ,  $\mathbf{a}, \tilde{\mathbf{a}} \in \mathcal{A}$ , such that for any agent  $i \in \mathcal{N}$ ,  $\mathbf{s}_{\mathcal{N}_{\kappa}^i} = \tilde{\mathbf{s}}_{\mathcal{N}_{\kappa}^i}$  and  $\mathbf{a}_{\mathcal{N}_{\kappa}^i} = \tilde{\mathbf{a}}_{\mathcal{N}_{\kappa}^i}$ . According to Equation (7), when only the state and action of the far-away agent are different, one have

$$\begin{aligned} &|A_{\pi}^i(\mathbf{s}, \mathbf{a}) - A_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| \\ &= |(Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - V_{\pi}^i(\mathbf{s})) - (Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}}) - V_{\pi}^i(\tilde{\mathbf{s}}))| \\ &= |(Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})) + (V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}}))| \\ &\leq |Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| + |V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}})|. \end{aligned} \quad (31)$$

Next, we analyze  $|Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})|$  and  $|V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}})|$  separately.

Firstly, for  $|Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})|$ , we have

$$\begin{aligned} &|Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| \\ &= \left| \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathbf{R}(\mathbf{s}_t, \mathbf{a}_t) \mid \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}] - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[\mathbf{R}(\mathbf{s}_t, \mathbf{a}_t) \mid \pi, \mathbf{s}_0 = \tilde{\mathbf{s}}, \mathbf{a}_0 = \tilde{\mathbf{a}}] \right| \\ &\leq \sum_{t=0}^{\infty} \gamma^t |\mathbb{E}[\mathbf{R}(\mathbf{s}_t, \mathbf{a}_t) \mid \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}] - \mathbb{E}[\mathbf{R}(\mathbf{s}_t, \mathbf{a}_t) \mid \pi, \mathbf{s}_0 = \tilde{\mathbf{s}}, \mathbf{a}_0 = \tilde{\mathbf{a}}]| \\ &\leq \sum_{t=1}^{\infty} \gamma^t D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i), \end{aligned}$$

where  $\rho_t^i$  and  $\tilde{\rho}_t^i$  are the distributions at time  $t$  with starting point  $(\mathbf{s}, \mathbf{a})$  and  $(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})$ , respectively. We use the result in Lemma A.3 to bound  $D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i)$ . The structure of our MDP implies that:

$$P(\mathbf{s}_{t+1}, \mathbf{a}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t) = \prod_{i \in \mathcal{N}} \pi^i(\mathbf{a}_{t+1}^i \mid \mathbf{s}_{\mathcal{N}_{\kappa}^i, t+1}) P^i(\mathbf{s}_{t+1}^i \mid \mathbf{s}_{\mathcal{N}_{\kappa}^i, t}, \mathbf{a}_t^i).$$

Then, if Assumption 2.1 holds, the requirements of Lemma A.3 are satisfied, one have  $D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i) \leq \zeta^t e^{-\beta\kappa}$  and

$$|Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| \leq \sum_{t=1}^{\infty} \gamma^t D_{\text{TV}}(\rho_t^i, \tilde{\rho}_t^i) \leq e^{-\beta\kappa} \sum_{t=1}^{\infty} \gamma^t \zeta^t = \frac{\gamma\zeta}{1-\gamma\zeta} e^{-\beta\kappa}, \quad (32)$$

where  $\zeta$  is defined in Assumption 2.1.

Let

$$\delta^j Q_{\pi}^i(\mathbf{s}, \mathbf{a}) = \sup_{z^j, \mathbf{z}^{-j}, z'^j} |Q_{\pi}^i(z^j, \mathbf{z}^{-j}) - Q_{\pi}^i(z'^j, \mathbf{z}^{-j})|,$$

and the MDP satisfies the condition of Lemma A.4, one can obtain

$$\sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \delta^j (Q_{\pi}^i(\mathbf{s}, \cdot)) \leq \sum_{t=1}^{\infty} \gamma^t \sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \delta^j P^i \leq \sum_{t=1}^{\infty} \gamma^t \zeta^t = \frac{\gamma \zeta}{1 - \gamma \zeta}.$$

Secondly, building on Assumption 2.2 and Lemma A.4, we analyze  $|V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}})|$  can further obtain

$$\begin{aligned} & |V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}})| \\ &= |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\tilde{\mathbf{s}}, \mathbf{a})| \\ &= |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\tilde{\mathbf{s}}, \mathbf{a})| \\ &\leq |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a})| + |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\tilde{\mathbf{s}}, \mathbf{a})| \\ &\leq \sum_{j \in \mathcal{N}} D_{\text{TV}}(\pi^j(\cdot|\mathbf{s}), \pi^j(\cdot|\tilde{\mathbf{s}})) \delta^j Q_{\pi}^i(\mathbf{s}, \mathbf{a}) + \frac{\gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa} \\ &\leq \xi e^{-\beta \kappa} \sum_{j \in \mathcal{N}} e^{-\beta d(j,i)} \delta^j Q_{\pi}^i(\mathbf{s}, \mathbf{a}) + \frac{\gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa} \\ &\leq \frac{\gamma \zeta}{1 - \gamma \zeta} \xi e^{-\beta \kappa} + \frac{\gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa} \\ &\leq \frac{(1 + \xi) \gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa}. \end{aligned} \tag{33}$$

Then, bringing (32) and (33) into (31), we have

$$\begin{aligned} & |A_{\pi}^i(\mathbf{s}, \mathbf{a}) - A_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| \\ &\leq |Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| + |V_{\pi}^i(\mathbf{s}) - V_{\pi}^i(\tilde{\mathbf{s}})| \\ &\leq |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a})| + 2 |\mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\tilde{\mathbf{s}})} Q_{\pi}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})| \\ &\leq \frac{\gamma \zeta}{1 - \gamma \zeta} \xi e^{-\beta \kappa} + \frac{2\gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa} \\ &\leq \frac{(2 + \xi) \gamma \zeta}{1 - \gamma \zeta} e^{-\beta \kappa}. \end{aligned} \tag{34}$$

Finally, denoting  $(\eta, \phi) = \left( \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ , we can obtain the Proposition 3.3.  $\square$

#### C.4 The proof of Corollary 3.4

*Proof.* Firstly, according to Lemma 3.1 and Definition 3.2, the following result holds when each agent adopts a sequential update scheme to optimize policy, i.e., we have

$$\begin{aligned} & \left| L_{\pi}^{1:i}(\bar{\pi}^{1:i-1}, \bar{\pi}^i) - L_{\pi_{\kappa}^i}^i(\bar{\pi}_{\kappa}^i) \right| \\ &= \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \mathbb{E}_{\mathbf{s}_{N_{\kappa}^i} \sim \rho_{\pi_{\kappa}^i}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi_{\kappa}^i}^i(\mathbf{s}_{N_{\kappa}^i}, \mathbf{a}^i)] \right| \\ &= \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^i)] \right. \\ &\quad \left. + \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^i)] - \mathbb{E}_{\tilde{\mathbf{s}} \sim \tilde{\rho}_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\tilde{\mathbf{s}}, \mathbf{a}^i)] \right| \\ &\leq \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^i)] \right| \\ &\quad + \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\mathbf{s}, \mathbf{a}^i)] - \mathbb{E}_{\tilde{\mathbf{s}} \sim \tilde{\rho}_{\pi}, \mathbf{a}^i \sim \bar{\pi}^i} [A_{\pi}^i(\tilde{\mathbf{s}}, \mathbf{a}^i)] \right|. \end{aligned} \tag{35}$$

Then, based on Assumptions 2.2, we have

$$\begin{aligned}
& \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^{1:i} \sim \bar{\boldsymbol{\pi}}^{1:i}} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)] \right| \\
&= \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} \left[ \sum_{h=1}^{i-1} (\bar{\pi}^h - \pi^h) A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i) \right] \right| \\
&\leq \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} \left[ \sum_{h=1}^{i-1} |\bar{\pi}^h - \pi^h| |A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)| \right] \\
&\leq \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} \left[ M^i \sum_{h=1}^{i-1} |\bar{\pi}^h - \pi^h| \right] \quad M^i = \max_{\boldsymbol{\pi}^i} |A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)| \\
&\leq \frac{M^i}{1-\gamma} \sum_{h=1}^{i-1} \max_{\mathbf{s}} D_{\text{TV}}(\bar{\pi}^h, \pi^h) \\
&\leq \frac{M^i \xi}{1-\gamma} e^{-\beta \kappa}.
\end{aligned} \tag{36}$$

According to (34), we have

$$\begin{aligned}
& \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)] - \mathbb{E}_{\tilde{\mathbf{s}} \sim \tilde{\rho}_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\tilde{\mathbf{s}}, \mathbf{a}^i)] \right| \\
&\leq \mathbb{E} \left| [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}) - A_{\boldsymbol{\pi}}^i(\tilde{\mathbf{s}}, \tilde{\mathbf{a}})] \right| \\
&\leq \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta} e^{-\beta \kappa}.
\end{aligned} \tag{37}$$

Then, bringing (36) and (37) into (35), we have

$$\begin{aligned}
& \left| L_{\boldsymbol{\pi}}^{1:i}(\bar{\boldsymbol{\pi}}^{1:i-1}, \bar{\boldsymbol{\pi}}^i) - L_{\boldsymbol{\pi}}^i(\bar{\boldsymbol{\pi}}_{\kappa}^i) \right| \\
&\leq \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^{1:i} \sim \bar{\boldsymbol{\pi}}^{1:i}} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)] \right| \\
&\quad + \left| \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^i)] - \mathbb{E}_{\tilde{\mathbf{s}} \sim \tilde{\rho}_{\boldsymbol{\pi}}, \mathbf{a}^i \sim \bar{\boldsymbol{\pi}}^i} [A_{\boldsymbol{\pi}}^i(\tilde{\mathbf{s}}, \mathbf{a}^i)] \right| \\
&\leq \frac{M^i \xi}{1-\gamma} e^{-\beta \kappa} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta} e^{-\beta \kappa} \\
&\leq \left( \frac{M^i \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta} \right) e^{-\beta \kappa}.
\end{aligned}$$

Finally, denoting  $(\eta', \phi) = \left( \frac{M^i \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ , we can obtain the Corollary 3.4.  $\square$

### C.5 The proof of Proposition 3.5

*Proof.* From (12), we can obtain  $-\eta' \phi^{\kappa} \leq L_{\boldsymbol{\pi}}^{1:i}(\bar{\boldsymbol{\pi}}^{1:i-1}, \bar{\boldsymbol{\pi}}^i) - L_{\boldsymbol{\pi}_{\kappa}^i}^i(\bar{\boldsymbol{\pi}}_{\kappa}^i) \leq \eta' \phi^{\kappa}$ . By the trust region theorem in Theorem 1 from [17], we have

$$\begin{aligned}
J(\bar{\boldsymbol{\pi}}) - J(\boldsymbol{\pi}) &\geq \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a} \sim \bar{\boldsymbol{\pi}}} [A_{\boldsymbol{\pi}}(\mathbf{s}, \mathbf{a})] - \nu D_{\text{KL}}^{\max}(\boldsymbol{\pi}, \bar{\boldsymbol{\pi}}) \\
&\geq \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a} \sim \bar{\boldsymbol{\pi}}} [A_{\boldsymbol{\pi}}(\mathbf{s}, \mathbf{a})] - \sum_{i=1}^n \nu D_{\text{KL}}^{\max}(\boldsymbol{\pi}^i, \bar{\boldsymbol{\pi}}^i) \\
&= \sum_{i=1}^n \mathbb{E}_{\mathbf{s} \sim \rho_{\boldsymbol{\pi}}, \mathbf{a}^{1:i} \sim \bar{\boldsymbol{\pi}}^{1:i}} [A_{\boldsymbol{\pi}}^i(\mathbf{s}, \mathbf{a}^{1:i-1}, \mathbf{a}^i)] - \sum_{i=1}^n \nu D_{\text{KL}}^{\max}(\boldsymbol{\pi}^i, \bar{\boldsymbol{\pi}}^i) \\
&= \sum_{i=1}^n (L_{\boldsymbol{\pi}}^{1:i}(\bar{\boldsymbol{\pi}}^{1:i-1}, \bar{\boldsymbol{\pi}}^i) - \nu D_{\text{KL}}^{\max}(\boldsymbol{\pi}^i, \bar{\boldsymbol{\pi}}^i)) \\
&\geq \sum_{i=1}^n \left( L_{\boldsymbol{\pi}_{\kappa}^i}^i(\hat{\boldsymbol{\pi}}_{\kappa}^i) - \eta' \phi^{\kappa} - \nu D_{\text{KL}}^{\max}(\boldsymbol{\pi}^i | \hat{\boldsymbol{\pi}}^i) \right) \\
&\geq \sum_{i=1}^n \left( L_{\boldsymbol{\pi}_{\kappa}^i}^i(\hat{\boldsymbol{\pi}}_{\kappa}^i) - \eta' \phi^{\kappa} - \nu_{\kappa}^i D_{\text{KL}}^{\max}(\boldsymbol{\pi}_{\kappa}^i | \hat{\boldsymbol{\pi}}_{\kappa}^i) \right).
\end{aligned}$$

Then, when each agent sequentially solves the following optimization problem:

$$\bar{\pi}_\kappa^i = \arg \max_{\hat{\pi}_\kappa^i} \left( L_{\pi_\kappa^i}^i(\hat{\pi}_\kappa^i) - \eta' \phi^\kappa - \nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) \right),$$

where  $(\eta', \phi) = \left( \frac{M^i \xi}{1-\gamma} + \frac{(2+\xi)\gamma \zeta}{1-\gamma \zeta}, e^{-\beta} \right)$ ,  $\nu_\kappa^i = \frac{2\gamma \max_{\mathbf{s}_{\mathcal{N}_\kappa^i, \mathbf{a}^i}} |A_{\pi_\kappa^i}^i(\mathbf{s}_{\mathcal{N}_\kappa^i, \mathbf{a}^i})|}{(1-\gamma)^2}$ , and  $D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) = \max_{\mathbf{s}_{\mathcal{N}_\kappa^i}} D_{\text{KL}}(\pi^i(\cdot | \mathbf{s}_{\mathcal{N}_\kappa^i}), \hat{\pi}^i(\cdot | \mathbf{s}_{\mathcal{N}_\kappa^i}))$ , we have  $J(\bar{\pi}) - J(\pi) \geq \sum_{i=1}^n \left( L_{\pi_\kappa^i}^i(\hat{\pi}_\kappa^i) - \eta' \phi^\kappa - \nu_\kappa^i D_{\text{KL}}^{\max}(\pi_\kappa^i | \hat{\pi}_\kappa^i) \right)$ .  $\square$

### C.6 The proof of Corollary 3.6

*Proof.* Firstly, by generalizing the result about the return in (14), one can derive how the expected costs change when the agents update their policies. Inspired by [10], we provide the following lemma.

*Lemma C.3.* Let  $\pi$  and  $\bar{\pi}$  be joint policies. Let  $i \in \mathcal{N}$  be an agent, and  $j \in \{1, \dots, m^i\}$  be an index of one of its costs. The following inequality holds

$$J_j^i(\bar{\pi}) \leq J_j^i(\pi) + L_{j,\pi}^i(\bar{\pi}^i) + \nu_j^i \sum_{h=1}^i D_{\text{KL}}^{\max}(\pi^h, \bar{\pi}^h).$$

where  $L_{j,\pi}^i(\bar{\pi}^i) = \mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a}^i \sim \bar{\pi}^i} [A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)]$ ,  $\nu_j^i = \frac{2\gamma \max_{\mathbf{s}, \mathbf{a}^i} |A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)|}{(1-\gamma)^2}$ .

*Proof.* From the upper bound version of Theorem 1 of [17] applied to joint policies  $\bar{\pi}$  and  $\pi$ , we conclude that

$$J_j^i(\bar{\pi}) \leq J_j^i(\pi) + \mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)] + \frac{4\alpha^2 \gamma \max_{\mathbf{s}, \mathbf{a}^i} |A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)|}{(1-\gamma)^2},$$

where  $\alpha = D_{\text{TV}}^{\max}(\pi^{1:i}, \bar{\pi}^{1:i}) = \max_{\mathbf{s}} D_{\text{TV}}(\pi^{1:i}(\cdot | \mathbf{s}), \bar{\pi}^{1:i}(\cdot | \mathbf{s}))$ .

Then, using Pinsker's inequality  $D_{\text{TV}}(p, q)^2 \leq D_{\text{KL}}(p, q)/2$ , we obtain

$$J_j^i(\bar{\pi}) \leq J_j^i(\pi) + \mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)] + \frac{2\gamma \max_{\mathbf{s}, \mathbf{a}^i} |A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)|}{(1-\gamma)^2} D_{\text{TV}}^{\max}(\pi^{1:i}, \bar{\pi}^{1:i}),$$

where  $D_{\text{KL}}^{\max}(\pi^{1:i}, \bar{\pi}^{1:i}) = \max_{\mathbf{s}} D_{\text{KL}}(\pi^{1:i}(\cdot | \mathbf{s}), \bar{\pi}^{1:i}(\cdot | \mathbf{s}))$ .

Notice that we have  $\mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a}^{1:i} \sim \bar{\pi}^{1:i}} [A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)] = \mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a}^i \sim \bar{\pi}^i} [A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)]$  as the action of agents other than  $i$  do not change the value of the variable inside of the expectation. Furthermore,

$$\begin{aligned} D_{\text{KL}}^{\max}(\pi^{1:i}, \bar{\pi}^{1:i}) &= \max_{\mathbf{s}} D_{\text{KL}}(\pi^{1:i}(\cdot | \mathbf{s}), \bar{\pi}^{1:i}(\cdot | \mathbf{s})) \\ &\leq \max_{\mathbf{s}} \left( \sum_{h=1}^i D_{\text{KL}}(\pi^h(\cdot | \mathbf{s}), \bar{\pi}^h(\cdot | \mathbf{s})) \right) \\ &\leq \sum_{h=1}^i \max_{\mathbf{s}} (D_{\text{KL}}(\pi^h(\cdot | \mathbf{s}), \bar{\pi}^h(\cdot | \mathbf{s}))) \\ &= \sum_{h=1}^i D_{\text{KL}}^{\max}(\pi^h, \bar{\pi}^h). \end{aligned}$$

Setting  $\nu_j^i = \frac{2\gamma \max_{\mathbf{s}, \mathbf{a}^i} |A_{j,\pi}^i(\mathbf{s}, \mathbf{a}^i)|}{(1-\gamma)^2}$ , we finally obtain

$$J_j^i(\bar{\pi}) \leq J_j^i(\pi) + L_{j,\pi}^i(\bar{\pi}^i) + \nu_j^i \sum_{h=1}^{i-1} D_{\text{KL}}^{\max}(\pi^h, \bar{\pi}^h).$$

$\square$

Secondly, from (12), we can obtain  $-\eta' \phi^\kappa \leq L_{\pi}^{1:i}(\bar{\pi}^{1:i-1}, \bar{\pi}^i) - L_{\pi_\kappa^i}^i(\bar{\pi}_\kappa^i) \leq \eta' \phi^\kappa$ . By generalizing the result, we can obtain  $-\eta'' \phi^\kappa \leq L_{j,\pi}^{1:i}(\bar{\pi}^{1:i-1}, \bar{\pi}^i) - L_{j,\pi_\kappa^i}^i(\bar{\pi}_\kappa^i) \leq \eta'' \phi^\kappa$ . Further, we can derive the upper bounds about surrogate cost

$$\begin{aligned} J_j^i(\bar{\pi}) &\leq J_j^i(\pi) + L_{j,\pi}^i(\bar{\pi}^i) + \nu_j^i \sum_{h=1}^n D_{KL}^{max}(\pi^h, \bar{\pi}^h) \\ &\leq J_j^i(\pi) + L_{j,\pi_\kappa^i}^i(\bar{\pi}_\kappa^i) + \eta'' \phi^\kappa + \nu_{j,\kappa}^i \sum_{h=1}^i D_{KL}^{max}(\pi_\kappa^h, \bar{\pi}_\kappa^h). \end{aligned}$$

where  $L_{j,\pi_\kappa^i}^i(\bar{\pi}_\kappa^i) = \mathbb{E}_{s_{\pi_\kappa^i} \sim \rho_{\pi_\kappa^i}, a^i \sim \bar{\pi}^i} [A_{j,\pi_\kappa^i}^i(s_{\mathcal{N}_\kappa^i}, a^i)]$ ,  $(\eta'', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ ,  $\nu_{j,\kappa}^i = \frac{2\gamma \max_{s_{\mathcal{N}_\kappa^i}, a^i} |A_{j,\pi_\kappa^i}^i(s_{\mathcal{N}_\kappa^i}, a^i)|}{(1-\gamma)^2}$ , and  $M_j$  is a constant.  $\square$

### C.7 The proof of Theorem 3.7

*Proof.* Based on the conclusions in Proposition 3.5 and Corollary 3.6, we can derive that in order to realize reward performance improvement and satisfy safety constraints, agents have to sequentially maximize their surrogate returns and ensure that their surrogate costs stay below the corresponding safety thresholds. Meanwhile, they have to constrain the policy search to small local neighborhoods (w.r.t, max-KL distance). Therefore, the size of KL constraint in Equation (16) should be set as

$$\delta_\kappa^i = \min \left\{ \min_{h \leq i-1} \min_{1 \leq j \leq m^h} \frac{c_j^h - J_j^h(\pi) - L_{j,\pi_\kappa^h}^i(\bar{\pi}_\kappa^h) - \eta'' \phi^\kappa - \nu_{j,\kappa}^h \sum_{l=1}^{i-1} D_{KL}^{max}(\pi_\kappa^l, \bar{\pi}_\kappa^l)}{\nu_{j,\kappa}^h}, \min_{h \geq i+1} \min_{1 \leq j \leq m^h} \frac{c_j^h - J_j^h(\pi) - \nu_{j,\kappa}^h \sum_{l=1}^{i-1} D_{KL}^{max}(\pi_\kappa^l, \bar{\pi}_\kappa^l)}{\nu_{j,\kappa}^h} \right\}, \quad (38)$$

where  $h \in \mathcal{N}_\kappa^i$  is the  $\kappa$ -hop neighbors of agent  $i$ , and  $j \in \{1, \dots, m^h\}$  is its cost index.

Note that  $\delta_\kappa^i$  is guaranteed to be non-negative if  $\pi$  satisfies safety constraints; that is because then  $c_j^h \geq J_j^h(\pi)$  for all  $h \in \mathcal{N}$ , and  $j \in \{1, \dots, m^i\}$ , and the set  $\{h \mid h < i\}$  is empty.

This formula for  $\delta_\kappa^i$ , combined with Lemma 3.1, assures that the policies  $\pi_\kappa^i$  within  $\delta_\kappa^i$  max-KL distance from  $\pi_\kappa^i$  will not violate other agents' safety constraints, as long as the base joint policy  $\pi$  did not violate them (which assures  $\delta_\kappa^1 \geq 0$ ). To see this, for every  $h = 1, \dots, i-1$ , and  $j = 1, \dots, m^h$ , we have

$$D_{KL}^{max}(\pi_\kappa^i, \bar{\pi}_\kappa^i) \leq \delta_\kappa^i \leq \frac{c_j^h - J_j^h(\pi) - L_{j,\pi_\kappa^h}^i(\bar{\pi}_\kappa^h) - \eta'' \phi^\kappa - \nu_{j,\kappa}^h \sum_{l=1}^{i-1} D_{KL}^{max}(\pi_\kappa^l, \bar{\pi}_\kappa^l)}{\nu_{j,\kappa}^h},$$

which implies

$$J_j^h(\pi) + L_{j,\pi_\kappa^h}^i(\bar{\pi}_\kappa^h) + \eta'' \phi^\kappa + \nu_{j,\kappa}^h \sum_{l=1}^{i-1} D_{KL}^{max}(\pi_\kappa^l, \bar{\pi}_\kappa^l) + \nu_{j,\kappa}^h D_{KL}^{max}(\pi_\kappa^i, \bar{\pi}_\kappa^i) \leq c_j^h. \quad (39)$$

By Corollary 3.6, the left-hand side of the inequality (39) is an upper bound of  $J_j^h(\bar{\pi}^{1:i-1}, \pi^i)$ , which implies that the update of agent  $i$  does not violate the constraint of  $J_j^h$ . The fact that the constraints of  $J_j^h$  for  $h \geq i+1$  are not violated, i.e.,

$$J_j^h(\pi) + \nu_{j,\kappa}^h \sum_{l=1}^{i-1} D_{KL}^{max}(\pi_\kappa^l, \bar{\pi}_\kappa^l) + \nu_{j,\kappa}^h D_{KL}^{max}(\pi_\kappa^i, \bar{\pi}_\kappa^i) \leq c_j^h.$$

Therefore, let  $(\eta', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ ,  $(\eta'', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\gamma\zeta}{1-\gamma\zeta}, e^{-\beta} \right)$ ,  $\nu_\kappa^i = \frac{2\gamma \max_{s_{\mathcal{N}_\kappa^i}, a^i} |A_{\pi_\kappa^i}^i(s_{\mathcal{N}_\kappa^i}, a^i)|}{(1-\gamma)^2}$ ,  $\nu_{j,\kappa}^i = \frac{2\gamma \max_{s_{\mathcal{N}_\kappa^i}, a^i} |A_{j,\pi_\kappa^i}^i(s_{\mathcal{N}_\kappa^i}, a^i)|}{(1-\gamma)^2}$ ,  $\delta_\kappa^i =$

$$\min \left\{ \min_{h \leq i-1} \min_{1 \leq j \leq m^h} \frac{\Xi_j^h - L_{j, \pi_{\kappa}^h}^h(\bar{\pi}_{\kappa}^h) - \eta'' \phi^{\kappa}}{\nu_{j, \kappa}^h}, \min_{h \geq i+1} \min_{1 \leq j \leq m^h} \frac{\Xi_j^h}{\nu_{j, \kappa}^h} \right\}, \quad \Xi_j^h = c_j^h - J_j^h(\pi_{\kappa}^h) - \nu_{j, \kappa}^h \sum_{l=1}^{i-1} D_{\text{KL}}^{\max}(\pi_{\kappa}^l, \hat{\pi}_{\kappa}^l),$$
 when the policy is updated by following a sequential update scheme, that is, each agent sequentially solves the following optimization problem:

$$\begin{aligned}
 \bar{\pi}_{\kappa}^i &= \arg \max_{\hat{\pi}_{\kappa}^i \in \bar{\Pi}_{\kappa}^i} \left( L_{\pi_{\kappa}^i}^i(\hat{\pi}_{\kappa}^i) - \eta' \phi^{\kappa} - \nu_{\kappa}^i D_{\text{KL}}^{\max}(\pi_{\kappa}^i | \hat{\pi}_{\kappa}^i) \right), \\
 \text{s.t. } &\left\{ \hat{\pi}_{\kappa}^i \in \bar{\Pi}_{\kappa}^i \mid D_{\text{KL}}^{\max}(\pi_{\kappa}^i, \hat{\pi}_{\kappa}^i) \leq \delta_{\kappa}^i, \text{ and} \right. \\
 &\left. J_j^i(\pi_{\kappa}^i) + L_{j, \pi_{\kappa}^i}^i(\hat{\pi}_{\kappa}^i) + \eta'' \phi^{\kappa} + \nu_{j, \kappa}^i D_{\text{KL}}^{\max}(\pi_{\kappa}^i, \hat{\pi}_{\kappa}^i) \leq c_j^i - \nu_{j, \kappa}^i \sum_{h=1}^{i-1} D_{\text{KL}}^{\max}(\pi_{\kappa}^h, \hat{\pi}_{\kappa}^h) \right\},
 \end{aligned}$$

the joint policy  $\pi$  has the monotonic improvement property,  $J(\bar{\pi}) \geq J(\pi)$ , as well as it satisfies the safety constraints,  $J_j^i(\bar{\pi}) \leq c_j^i$ , for any agent  $i \in \mathcal{N}$  and its cost index  $j \in \{1, \dots, m^i\}$ .  $\square$

## C.8 Algorithm

In this subsection, we provide the main pseudocode for Scalable MAPPO-Lagrangian (Scal-MAPPO-L), as outlined in Algorithm 1.

---

### Algorithm 1 Scalable MAPPO-Lagrangian

---

**Input:** Stepsizes  $\alpha_{\theta}, \alpha_{\lambda}$ , batch size  $B$ , number of agents  $n$ , episodes  $Z$ , steps per episode  $T$ , discount factor  $\gamma$ , parameter  $\kappa$ .

**Initialize:** Actor networks  $\theta_{\kappa, 0}^1, \dots, \theta_{\kappa, 0}^n$ , V-value network  $\chi_{\kappa, 0}^1, \dots, \chi_{\kappa, 0}^n$ , V-cost networks  $\{\phi_{j, 0}^i\}_{1 \leq j \leq m^i, \forall i \in \mathcal{N}}$ , Replay buffer  $\mathcal{B}$ .

- 1: **for**  $z = 0, 1, \dots, Z - 1$  **do**
  - 2:   Collect a set of trajectories by running the policies  $\pi_{\theta_{\kappa}^1}, \dots, \pi_{\theta_{\kappa}^n}$ .
  - 3:   Push transitions  $\{(o_t^i, a_t^i, o_{t+1}^i, r_t^i), \forall i \in \mathcal{N}, t \in T\}$  into  $\mathcal{B}$ .
  - 4:   Sample a random minibatch of  $B$  transitions from  $\mathcal{B}$ .
  - 5:   **for**  $i = 1 : n$  **do**
  - 6:     Initialize a policy parameter  $\theta_{\kappa, 0}^i$  and Lagrangian multipliers  $\lambda_j^i, \forall i \in \mathcal{N}, j \in 1, \dots, m^i$ .
  - 7:     Compute advantage function  $\hat{A}^i(\mathbf{s}, \mathbf{a}^i)$  and cost advantage functions  $\hat{A}_j^i(\mathbf{s}, \mathbf{a}^i)$ .
  - 8:     Compute the parameters  $\eta', \eta'', \phi, \nu_{\kappa}^i$  and  $\nu_{j, \kappa}^i, \forall j \in \{1, \dots, m^i\}$ .
  - 9:     Compute the radius of the KL-constraint  $\delta_{\kappa}^i$ .
  - 10:    Compute the advantage function in (18).
  - 11:    Update policy according to (20).
  - 12:    Update V-value network and V-cost networks.
  - 13:    **end for**
  - 14: **end for**
- 

The algorithm has a simple idea that each agent independently optimizes the surrogate objective, which only depends on its action and the state of its  $\kappa$ -hop neighbors for each agent. In the actual execution, we adopt the surrogate objective (20) instead of (16). It actually uses some approximations for the decentralized surrogate objective, the same as the MAPPO-L [10]. Most of these approximations are traditional practices in RL, yet they may make it impossible for the practical algorithm to rigorously maintain the theoretical guarantees in Theorem 3.7. However, we need to argue that we should go one step further and provide a decentralized surrogate for decentralized learning with a convergence guarantee. We believe and expect that a better practical method can be found based on this objective in future work.



## D Supplementary materials for Section 4

### D.1 Additional experimental results

In this paper, we compare the algorithm of our proposed (i.e., Scal-MAPPO-L in Algorithm 1) against other PPO family algorithms on several safe MARL tasks to evaluate their performance. Here, we provide some additional experimental results, which are illustrated in Figures 3-4.

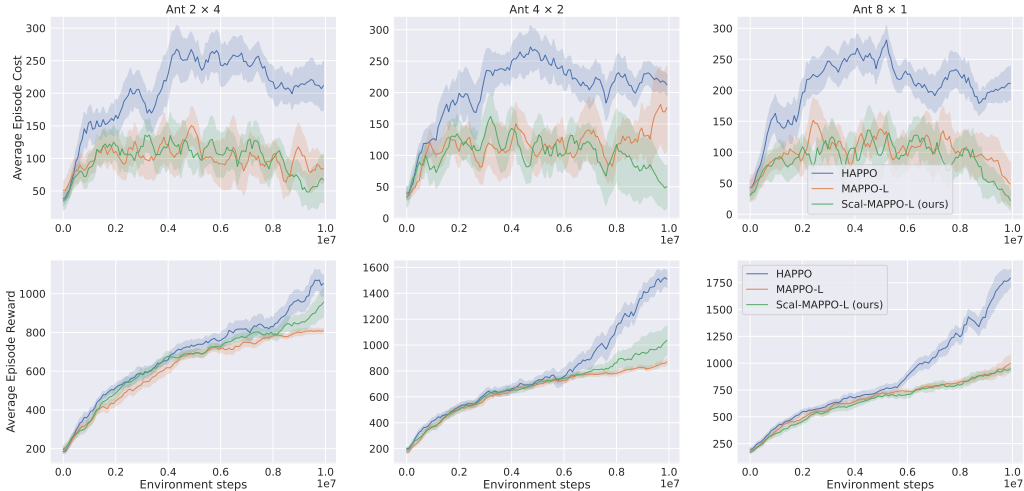


Figure 3: Performance comparisons in terms of cost and reward on three Safe Ant-v2 tasks. Each column subfigure represents a different task, and we plot the cost curves (the lower the better) in the upper row and the reward curves (the higher the better) in the bottom row for each task.

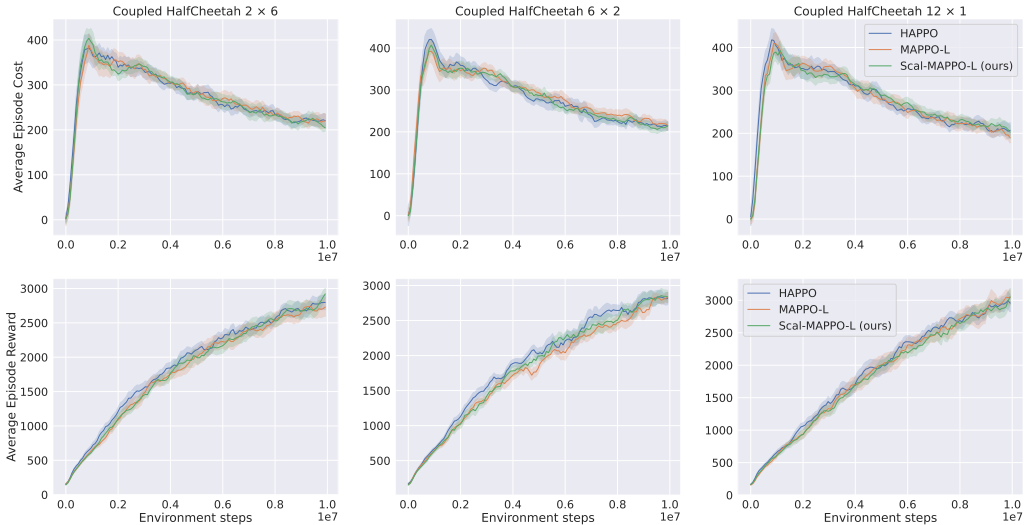


Figure 4: Performance comparisons in terms of cost and reward on three Safe Coupled HalfCheetah tasks.

*Remark D.1.* It is worth pointing out that, in our code, unlike the original, the global state consists of a patchwork of each agent’s ID and the  $\kappa$ -hop information rather than a long state vector. This is the main reason of the difference in performance from the original paper. As we consider decentralized learning, the agents in the experiments do not use parameter-sharing. In all experiments, the network architectures and common hyperparameters of our algorithm and MAPPO-L are the same for a fair

comparison. All reported results are averaged over three or more random seeds, and the curves are smooth over time.

## D.2 The computer resources and computational complexity

**Computation resources:** We executed our code on a computer with NVIDIA GeForce RTX 4090 (GPU) and Intel Core i9-13900K (CPU).

**Computational complexity:** The computational complexity of Scal-MAPPO-L (ours) is  $O(TNMHP)$ , where  $T$  denotes the number of steps,  $N$  denotes the number of agents,  $M$  denotes the number of constraints,  $H$  denotes the number of PPO-Epoch, and  $P$  denotes the number of policy parameters.

Besides, we test the running time of Scal-MAPPO-L on Safe Manyagent Ant  $6 \times 1$ , Safe Ant  $8 \times 1$ , and Safe Coupled HalfCheetah  $12 \times 1$ . The running steps are  $1 \times 10^7$  in each environment. When the parameter  $\kappa$  is maximized, the algorithm’s average wall-clock times are  $8.43h$ ,  $9.28h$ , and  $11.65h$ , respectively. It is worth noting that the wall-clock times do not significantly down when  $\kappa$  gradually decreases. This may be due to the fact that we have yet to consider the process of sending and receiving information realistically. However, based on the successful research conducted in the field of communication [51, 52], it is evident that algorithms requiring less communication undoubtedly have an advantage in terms of reducing communication burden and enhancing applicability.

## E The discussion of limitations and impacts

### E.1 Limitations

This paper is centered on theoretical analysis and also contains practical algorithms and simulation verification. The main results in the paper characterize the proposed method’s performance in terms of safety constraints and joint policy improvement. Below, we discuss the limitations of the proposed approach for both theory and experiment aspects as follows:

- 1) Our theoretical results are based on the two assumptions about the spatial decay of correlation for the dynamics and the policies in Assumption 2.1 and Assumption 2.2. Our conclusions may be useless when such assumptions do not hold, e.g., the decisions of each agent are non-negligibly related to the decisions of all other agents. However, fortunately, existing works [46, 47] have shown that many real-world situations satisfy both assumptions, so our study is still important and meaningful.
- 2) Our experiments show that Scal-MAPPO-L, with communication between a small number of neighbors, outperforms MAPPO-L in some cases, which we would like to see because it implies fewer communication requirements. However, we have yet to develop an equilibrium relationship between the amount of communication and the performance, which we will focus on next.

### E.2 Broader impacts

This paper presents work that aims to advance the field of RL, especially safe MARL. Our work has many positive societal impacts, such as providing a theoretical foundation for scalable Safe MARL, none of which we feel must be specifically highlighted. There are no negative societal impacts on our work.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Both the abstract and introduction include the claims made in the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations of the proposed approach for both theory and experiment aspects, which can be seen in Section E.1.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: This paper provide the full set of assumptions and a complete (and correct) proof. Specifically, we describe the two assumptions about the spatial correlation of the transition dynamics and policies in subsection 2.2 and provide complete proofs of the paper’s propositions, corollaries, and theorems in Appendix C.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide the pseudocode of Algorithm 1 in Appendix C.8 and provide a detailed description of the experiments in Section 4 and Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce

the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The paper provide open access to the data and code in supplemental material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: This paper states the experimental setup in Section 4.1. A more detailed setup can be found in Appendix D.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: In this paper, all reported results are averaged over three or more random seeds. The results are presented by error line plots and the curves are smooth over time.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper provide information on the computer resources in Appendix D.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper discusses both potential positive societal impacts and negative societal impacts of the work performed in Appendix E.2.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite the original paper that produced the code package.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We provide details of new assets in supplemental material..

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.



- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.