

# 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.

Overview

1

**Challenge**

2

**Research Status and Innovative Method**

3

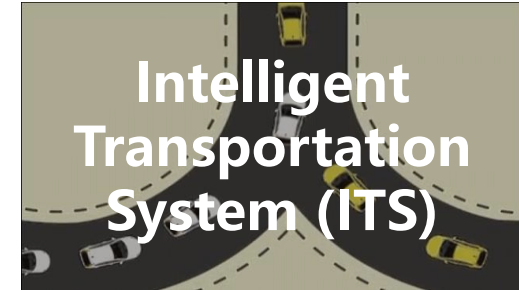
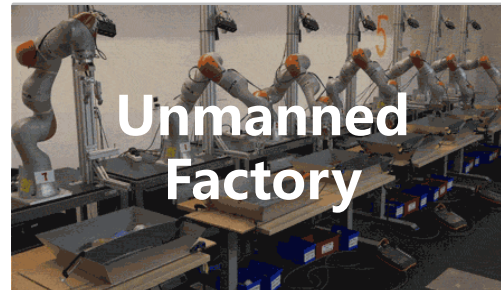
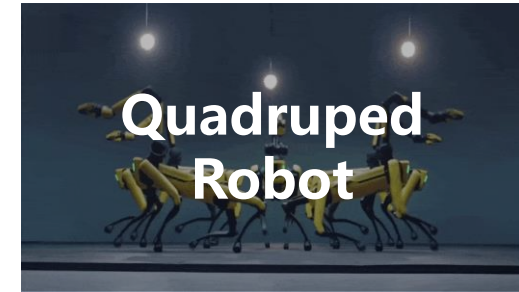
**Specific Implementation**

4

**Experiment Results**

## ■ Safe Multi-agent Reinforcement Learning Problem

- The real world generally can be viewed as a multi-agent environment. 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.

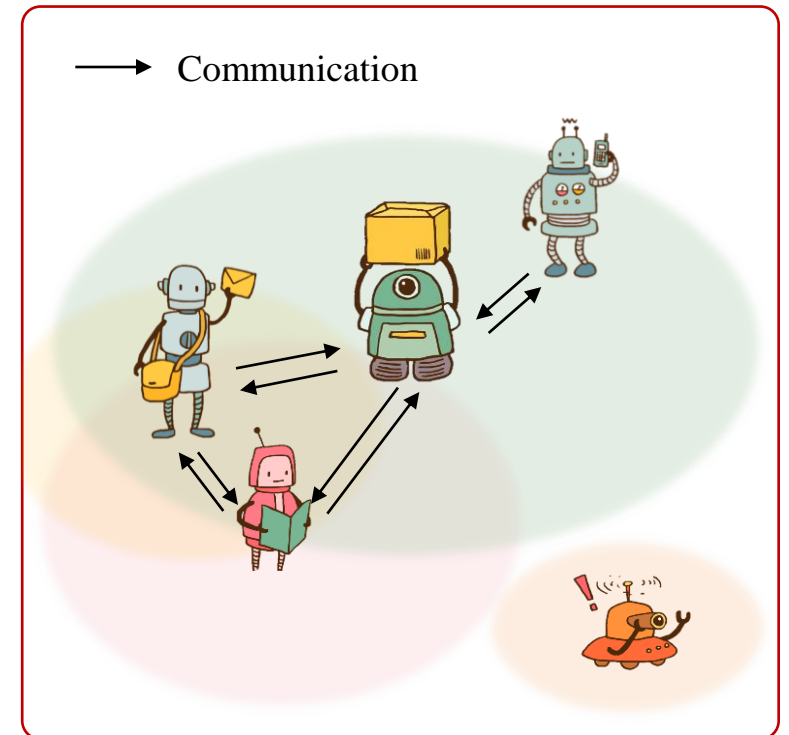


## ■ Research Status

- 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.

## ■ Our method: Scal-MAPPO-L

- We develop **a novel scalable and theoretically-justified multi-agent constrained policy optimization method**.
- Specifically, our method has the following characteristics:
  - ◆ Policy is updated by following a sequential update scheme
  - ◆ Decentralized training with local interactions
  - ◆ The safety constraints and the joint policy improvement can be met



# Specific Implementation

## Problem

How to overcome the global coupling arising from safety constraints and the exponential growth of the state-action space size on the applicability of algorithms?



## Solution

We introduce two assumptions about spatial correlation decay and updating the local policy based on the trust region method with a sequential update scheme.

## Implementation Idea

Two assumptions about spatial correlation decay are introduced.

Based on these assumptions, the maximum information loss regarding the advantage truncation is quantified.

The local policy optimization objective is provided by integrating the bounds of the trust region method and the bounds of the truncated advantage function.

A performance guarantee theorem is provided.

**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

$$\max_{i \in \mathcal{N}} \sum_{j \in \mathcal{N}} e^{\beta d(i,j)} \mathbb{1}^{ij} \leq \zeta, \quad (4)$$

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  $i \in \mathcal{N}$ ,  $s_{N_2^i} \in \mathcal{S}_{N_2^i}, s_{N_2^{i+1}}, s_{N_2^{i+2}} \in \mathcal{S}_{N_2^{i+1}}$ , one have

$$\sup_{s_{N_2^i}, s_{N_2^{i+1}}, s_{N_2^{i+2}}} \left| \pi^i(s_{N_2^i}, s_{N_2^{i+1}}) - \pi^i(s_{N_2^i}, s_{N_2^{i+2}}) \right| \leq \xi e^{-\beta n}. \quad (5)$$

**Proposition 3.3.** For any agent  $i \in \mathcal{N}$ , let the parameters  $(\eta, \phi) = \left( \frac{M^i \zeta}{1-\gamma}, e^{-\beta} \right)$ . If Assumption 2.1 and Assumption 2.2 hold, for any  $z_{N_2^i} = (s_{N_2^i}, a_{N_2^i}) \in \mathcal{S}_{N_2^i} \times \mathcal{A}_{N_2^i}$ , the exponential decay property of the advantage function holds, i.e., we have

$$\sup_{z_{N_2^i}, z_{N_2^{i+1}}, z_{N_2^{i+2}}} \left| A^i(z_{N_2^i}, z_{N_2^{i+1}}) - A^i(z_{N_2^i}, z_{N_2^{i+2}}) \right| \leq \eta \phi^n. \quad (11)$$

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

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

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

**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

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

where  $L_{j, \pi_k}^i(\bar{\pi}_k^i) = \mathbb{E}_{s_{N_2^i} \sim \rho_{\pi_k}^i, a^i \sim \bar{\pi}_k^i} \left[ A_{j, \pi_k}^i(s_{N_2^i}, a^i) \right]$ ,  $\nu_{j, \kappa}^i = \frac{2\gamma \max_{s_{N_2^i}, a^i} |A_{j, \pi_k}^i(s_{N_2^i}, a^i)|}{(1-\gamma)^2}$ ,  $(\eta'', \phi) = \left( \frac{M_j \xi}{1-\gamma} + \frac{(2+\xi)\zeta}{1-\gamma\xi}, e^{-\beta} \right)$ , and  $M_j$  is a constant.

**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:

$$\begin{aligned} \bar{\pi}_k^i &= \arg \max_{\pi_k^i \in \Pi_k^i} \left( L_{\pi_k^i}^i(\pi_k^i) - \eta' \phi^n - \nu_{j, \kappa}^i D_{\text{KL}}^{\max}(\pi_k^i, \bar{\pi}_k^i) \right), \\ \text{s.t. } \{ \bar{\pi}_k^i \in \Pi_k^i \mid D_{\text{KL}}^{\max}(\pi_k^i, \bar{\pi}_k^i) \leq \delta_k^i, \text{ and} \\ J_j^i(\pi_k) + L_{j, \pi_k}^i(\bar{\pi}_k^i) + \eta'' \phi^n + \nu_{j, \kappa}^i D_{\text{KL}}^{\max}(\pi_k^i, \bar{\pi}_k^i) \leq c_j^i - \nu_{j, \kappa}^i \sum_{b=1}^{i-1} D_{\text{KL}}^{\max}(\pi_b^i, \bar{\pi}_k^i) \}, \end{aligned} \quad (16)$$

where  $(\eta', \phi) = \left( \frac{M^i \zeta}{1-\gamma} + \frac{(2+\xi)\zeta}{1-\gamma\xi}, e^{-\beta} \right)$ ,  $(\eta'', \phi) = \left( \frac{M^i \zeta}{1-\gamma} + \frac{(2+\xi)\zeta}{1-\gamma\xi}, e^{-\beta} \right)$ ,  $\nu_{j, \kappa}^i = \frac{2\gamma \max_{s_{N_2^i}, a^i} |A_{j, \pi_k}^i(s_{N_2^i}, a^i)|}{(1-\gamma)^2}$ ,  $\delta_k^i = \min \left\{ \min_{a \leq i-1} \min_{1 \leq j \leq m^a} \frac{\bar{\pi}_j^a - L_{j, \pi_k}^a(\bar{\pi}_k^a) - \eta'' \phi^n}{\nu_{j, \kappa}^a}, \min_{a \geq i+1} \min_{1 \leq j \leq m^a} \frac{\bar{\pi}_j^a}{\nu_{j, \kappa}^a} \right\}$ ,  $\bar{\pi}_k^i = c_j^i - J_j^i(\pi_k^i) - \nu_{j, \kappa}^i \sum_{b=1}^{i-1} D_{\text{KL}}^{\max}(\pi_b^i, \bar{\pi}_k^i)$ .

**Challenge 1:** How to quantify the information loss regarding the advantage truncation?

**Challenge 2:** How to ensure the local policy updates are not overly conservative?

**Challenge 3:** How to prove that the method can improve reward performance and satisfy safe constraints?

# Experiment Results

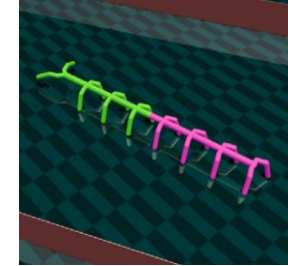
## Main results (under Safe MAMuJoCo environment)

With the truncation parameter  $\kappa \geq 3$ , we can observe that the performance of Scal-MAPPO-L improves considerably and can approach or even outperform MAPPO-L in some environments.

Safe 4 × 2 Ant task



Safe 2 × 4 Manyagent



Safe p1p-couple HalfCh

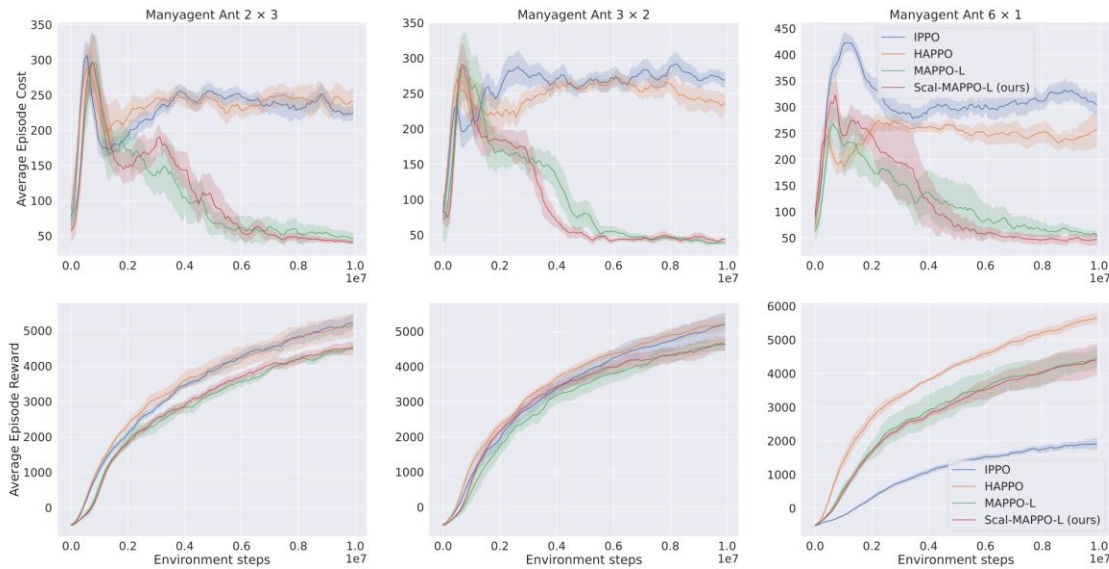


Figure 1: Performance comparisons in terms of cost and reward on three Safe ManyAgent Ant tasks. Each column subfigure represents a different task, and we plot the cost curves (the lower the better) in the upper row and the reward curves (the higher the better) in the bottom row for each task.

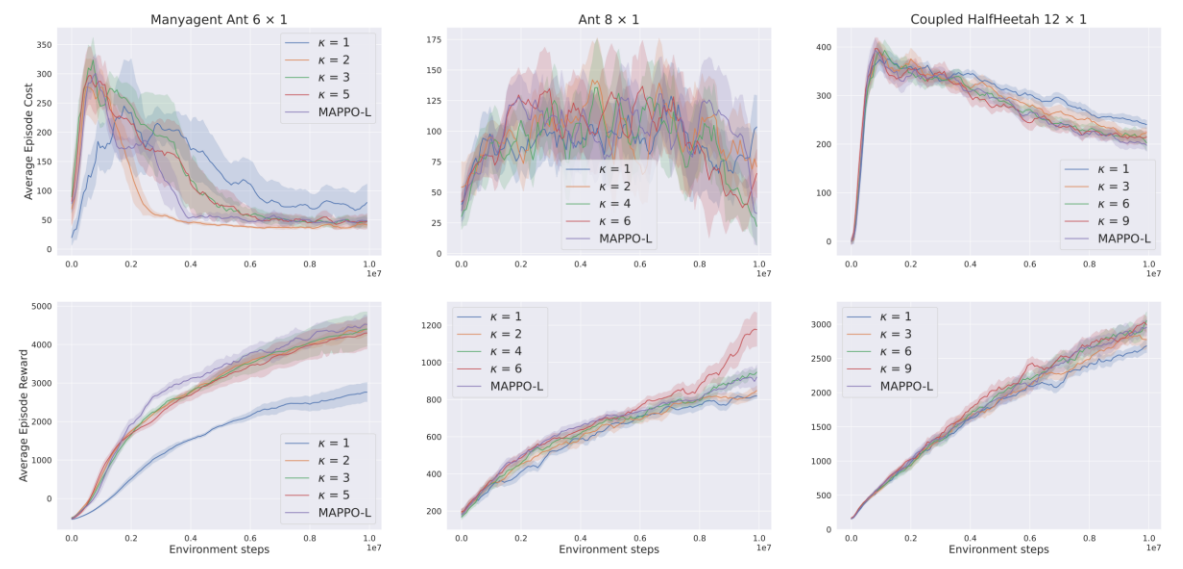
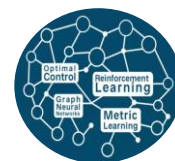


Figure 2: Performance comparisons in terms of cost and reward on Safe ManyAgent Ant task, Safe Ant task, and Safe Coupled HalfCheetah task. In each task, the performance of Scal-MAPPO-L with different  $\kappa$  and MAPPO-L are demonstrated.

**Lijun Zhang, Lin Li, Wei Wei\*, Huizhong Song, Yaodong Yang, Jiye Liang.** Scalable Constrained Policy Optimization for Safe Multi-agent Reinforcement Learning. *NeurIPS 2024*.

# Thanks for Listening



**DM&KD Lab**

—数据挖掘与知识发现实验室—