

# Safe Multi-Agent Reinforcement Learning via Distributional Cost Critic and Maximum Entropy Optimization

Qiwei Liu<sup>1,2</sup>, Ye Yuan<sup>1</sup>, Lingyue Zhang<sup>1</sup>, Kaitian Chen<sup>1</sup>, Yunkai Lv<sup>1,3</sup>, Sheng Gao<sup>1</sup>, Huaicheng Yan<sup>1,4\*</sup>

<sup>1</sup>The Key Laboratory of Smart Manufacturing in Energy Chemical Process of the Ministry of Education East China University of Science and Technology (ECUST), Shanghai, China

<sup>2</sup>Shanghai Institute of Intelligent Science and Technology, Tongji University

<sup>3</sup>Key Laboratory of System Control and Information Processing, Shanghai Jiao Tong University

<sup>4</sup>Faculty of Artificial Intelligence, Shanghai University of Electric Power

{qw\_liu,yye,22012221,ktchen}@mail.ecust.edu.cn, {yunkailv,sheng\_gao,hcyan}@ecust.edu.cn

## Abstract

Deploying multi-agent reinforcement learning (MARL) in safety-critical systems faces significant challenges due to insufficient agent exploration and inadequate safety constraint guarantees. Current approaches are constrained by two fundamental limitations: inefficient exploration leading to suboptimal policies, and expected-cost-based constraint frameworks failing to ensure full-process safety. To address these challenges, this paper proposes a novel safety-aware maximum entropy (MaxEnt) MARL framework using Conditional Value-at-Risk (CVaR) as a joint safety metric, which quantifies constraint satisfaction under worst-case scenarios for multi-agent systems. Moreover, we develop the Worst-Case Multi-Agent Soft Actor-Critic (WCMASAC) algorithm, incorporating sequential update mechanisms and maximum entropy optimization for heterogeneous agents, enhanced with distributed safety critics. Theoretically, we establish the monotonic improvement property, guaranteed constraint satisfaction, and convergence to a quantum response equilibrium for WCMASAC. Extensive experiments on safety gymnasium-based benchmarks demonstrate that WCMASAC outperforms state-of-the-art baselines in both task reward acquisition and safety constraint violation reduction, while exhibiting superior exploration efficiency and risk-aware control capabilities.

Extended version —

<https://github.com/YeY-YYe/WCMASAC>

## 1 Introduction

As a core framework for decentralized decision-making, MARL encounters substantial difficulties in scenarios demanding rigorous safety guarantees. While traditional MARL algorithms optimize cooperative efficiency through long-term return maximization, their neglect of explicit safety constraints renders them incapable of ensuring compliance with safety boundary constraints in the state-action space of agent joint actions. This limitation substantially restricts the practical applications of MARL. In particular, safe

reinforcement learning (RL) for a single agent has achieved safe exploration capabilities in complex environments (Liu et al. 2022; Yang et al. 2021; Ha et al. 2021; Li et al. 2024; Yang et al. 2023; Zhang, Vuong, and Ross 2020).

However, the non-stationary nature of multi-agent systems introduces fundamental theoretical challenges that hinder the extension of such methods to multi-agent systems. Firstly, partial observability impedes individual agents from modeling the global risk distribution of joint actions. Secondly, the interactions require each agent to pay attention to its own actions and to consider the impact of other agents' decisions on safety (Gu et al. 2024). To address these challenges, several pioneering results on safe MARL have been published (Zhang et al. 2020; Lu et al. 2021; Liu et al. 2021; Gu et al. 2023; Zhang et al. 2024; Du, Gou, and Cai 2025). Dec-PG (Lu et al. 2021) integrates decentralized policy gradient with primal-duality methods to identify saddle points between task rewards and safety costs. However, its reliance on a parameter-sharing mechanism among agents may lead to suboptimal convergence. CMIX (Liu et al. 2021) extends QMIX (Rashid et al. 2020b) by incorporating peak and average constraints. Although all agents empirically satisfy these constraints in experiments, the method still lacks theoretical safety guarantees. Furthermore, it depends on value decomposition methods that are inherently limited by restrictive underlying assumptions. Based on HARL (Zhong et al. 2024), MACPO and MAPPO-Lagrangian (Gu et al. 2023) integrate constrained optimization with advantage decomposition theory (Kuba et al. 2021), ensuring both multi-agent safety constraint satisfaction and monotonic policy improvement properties, thereby establishing a novel theoretical framework for safe MARL.

Although the MACPO framework is theoretically sound and has demonstrated promising performance, it still faces several critical challenges. First, the MACPO and MAPPO-Lagrangian algorithm, which ensures agent safety during training and execution by constraining the expectation of cost, lacks quantitative modeling capabilities for tail risks associated with safety constraints. Such constraints based on the expectation of cost can regulate the average rate of constraint violations but fail to suppress instantaneous safety violations that may occur during execution (Yang et al. 2021,

\*Corresponding Author.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

2023). Risk-sensitive RL approaches inspired by return distribution learning offer a potential solution to this challenge (Dabney et al. 2018a; Qiu et al. 2021; Hu et al. 2022; Sun et al. 2023; Shen et al. 2023). These methods enhance agents’ sensitivity to risky behaviors and demonstrate superior performance in long-term returns by learning return distributions or designing risk-sensitive functions. Regrettably, existing distributional MARL methods primarily focus on the evolution of return probability distributions while failing to explicitly consider state-action joint safety constraints in multi-agent systems, thereby being unable to guarantee execution safety for agents.

Additionally, safety-constrained RL methods based on the MACPO framework may encounter insufficient exploration issues (Liu et al. 2023; Kuba et al. 2022; Sutton, Barto et al. 1998). While policy optimization approaches like MACPO ensure the safety feasibility of iterative policies, their conservative safety constraint mechanisms significantly restrict the exploration scope of the joint action space. In sparse safety feedback scenarios, this phenomenon triggers premature policy convergence, causing systems to settle for local safe optima at the expense of global performance. A potential solution involves enabling agents to learn stochastic behaviors for effective exploration of the reward space (Haarnoja et al. 2018). This creates dynamic game conflicts between individual entropy maximization objectives and global safety constraints, potentially leading some agents to exceed local safety thresholds in pursuit of enhanced exploration. While MaxEnt safe RL has advanced in single-agent settings (Yang et al. 2021), its extension to multi-agent systems remains largely unexplored, particularly in learning stochastic policies under safety constraints. Recent work by (Liu et al. 2023) proposes a MaxEnt heterogeneous agent reinforcement learning framework (MaxEnt-HARL), which ensures convergence of stochastic policies in unconstrained scenarios through decoupling of advantage functions. However, this approach fails to address dynamically evolving feasibility boundaries induced by safety constraints, limiting its applicability to safety-critical multi-agent coordination tasks.

Inspired by the above discussion, the principal contributions of this work are fourfold:

1. We establish the first theoretically-grounded MaxEnt safe MARL framework under distributional RL principles, providing a novel paradigm for stochastic policy optimization in constrained environments. Specifically, we innovatively employ mixed CVaR as a joint safety constraint metric for multi-agent systems, constructing a risk-sensitive safety mechanism by analyzing cost tail distribution.
2. We theoretically demonstrate the existence of equilibrium policies for constrained MaxEnt MARL optimization via a dual-decoupling design (risk-aware and advantage decomposition), while establishing theoretical guarantees for monotonic policy improvement and convergence to generalized Nash equilibrium.
3. Building upon this theoretical foundation, we propose WCMASAC, the first multi-agent soft actor-critic algo-

rithm incorporating distributional safety critics, which achieves risk-aware joint policy optimization through worst-case value estimation.

4. Comprehensive experiments in safe Multi-MuJoCo cooperative control scenarios validate WCMASAC’s superior performance, demonstrating higher sample efficiency and fewer constraint violations compared to state-of-the-art baselines (Ji et al. 2023), thereby offering a reliable solution for safety-critical multi-agent deployments.

## 2 Preliminaries and Problem Formulation

### 2.1 Constrained Cooperative Markov Game

Similarly to (Gu et al. 2022), a fully cooperative safe MARL described as a decentralized constrained Markov game is considered, which is denoted as a tuple  $\langle \mathcal{N}, \mathcal{S}, \mathcal{A}, \mathbb{P}, R, \mathbf{C}, \mathbf{d} \rangle$ .  $\mathcal{N} = \{1, 2, \dots, n\}$  is the set of agents.  $\mathcal{S}$  is the state space.  $\mathcal{A}$  denotes the joint action space and there exists  $|\mathcal{A}| < \infty$ , where  $\mathcal{A} = \prod_{i=1}^n \mathcal{A}_i$  is the product of the action spaces of all agents. Let  $\mathbf{a}(t) = (a_1(t), \dots, a_n(t)) \in \mathcal{A}$  be the joint action, and the joint policy can be expressed as  $\pi(\mathbf{a}(t)|s(t)) = \prod_{i=1}^n \pi_i(a_i(t)|s(t))$ .

The probabilistic state transition function is presented as  $\mathbb{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ . When the system employs an action  $\mathbf{a}(t)$  in state  $s(t)$ , the probability of transferring to the next moment state  $s(t+1)$  is  $\mathbb{P}(s(t+1)|s(t), \mathbf{a}(t))$ , and agents receive the reward  $R(s(t), \mathbf{a}(t))$ . The primary goal of each agent is to maximize the expectation of standard cumulative rewards  $J(\pi) = \mathbb{E}_{s_{0:\infty} \sim \mathbb{P}, \mathbf{a}_{0:\infty} \sim \pi} [\sum_{t=1}^{\infty} R(s(t), \mathbf{a}(t))]$ .

Each agent’s action may adversely affect the overall safety of the whole multi-agent system. Therefore, the safety constraints are modeled as joint safety constraints under the actions of all agents, and  $\mathbf{C} = \{J_{i_j}\}_{1 \leq j \leq h}^{i \in \mathcal{N}}$  can be considered as the set of cost functions  $J_{i_j}$  for agent  $i$ , where  $h$  is the number of constraints.

Then, agent  $i \in \mathcal{N}$  also has the following accumulated safety cost:

$$C_{i_j}(\pi) = \mathbb{E}_{s_{0:\infty} \sim \mathbb{P}, \mathbf{a}_{0:\infty} \sim \pi} [\sum_{t=0}^{\infty} J_{i_j}(s(t), a_i(t), \mathbf{a}_{-i}(t))], \quad (1)$$

which should satisfy  $C_{i_j}(\pi) \leq d_{i_j}$ .  $\mathbf{a}_{-i}(t)$  denotes the actions of all the other agents except for the agent  $i$ . The set of constraint thresholds is given as  $\mathbf{d} = \{d_{i_j}\}_{1 \leq j \leq h}^{i \in \mathcal{N}}$ .

### 2.2 Safety Based on Joint CVaR

Existing safe RL methods typically ensure system safety by constraining the expected cost, where agents optimize policies to maximize the expected return while ensuring that the expected cost remains below a predefined threshold (Gu et al. 2023, 2022; Zhang et al. 2024). However, long-term safety costs exhibit stochastic characteristics under the combined effects of policy stochasticity, agent interactions inducing environmental non-stationarity, and other

factors. This renders safety constraint frameworks based on expected cost inadequate for addressing potential risks arising from heavy-tailed distributions in long-term safety costs, making them unsuitable for direct application in safety-critical systems (Yang et al. 2021).

Therefore, inspired by distributional RL (Shen et al. 2023) and WCSAC (Yang et al. 2021), we adopt the Conditional Value-at-Risk (CVaR) (Rockafellar, Uryasev et al. 2000) safety measure to replace the expected cost as the safety constraint, with the essential definitions provided as follows.

**Definition 1.** (CVaR $_{\omega}$  with Risk level  $\omega$ ). Let  $C_{i_j} := C_{i_j}(s(t), a_i(t), \mathbf{a}_{-i}(t))$  for convenience. Then for any safe cost  $C_{i_j}(s(t), a_i(t), \mathbf{a}_{-i}(t))$ , its risk function CVaR $_{\omega}$  under policy  $\pi$  is expressed as  $\phi_{1-\omega}^{\pi}(C_{i_j}) := \mathbb{E}_{\mathbf{a}_{0:\infty} \sim \pi} [C_{i_j} | C_{i_j} \geq F_{C_{i_j}}^{-1}(1 - \omega)]$  where  $F_{C_{i_j}}^{-1}$  denotes the inverse function of  $C_{i_j}$ 's cumulative distribution function and the scalar  $\omega \in (0, 1]$  denotes the level of risk.  $\omega$  closer to 0 implies stronger safety, and vice versa represents a risk-neutral case.

**Definition 2.** (CVaR $_{\omega}$ -based safety constraint). Given a risk level  $\omega$  and the safety threshold value  $d_{i_j}$ , if the following constraint holds for the system under the policy  $\pi$ :

$$\phi_{1-\omega}^{\pi}(C_{i_j}) \leq d_{i_j} \forall t \in [0, \infty), \forall i \in \mathcal{N}, 1 \leq j \leq h. \quad (2)$$

The joint policy  $\pi$  can be considered as a worst-case safe policy.

In most fully cooperative multi-agent scenarios, a local safety constraint violation by an individual agent may induce global systemic risks. For instance, abrupt torque variations in a single robotic joint may lead to overall imbalance, or unexpected acceleration/deceleration behaviors of a member in a formation may trigger cascading collisions. Therefore, this paper comprehensively considers heavy-tailed distributional safety risks and assumes, without loss of generality, that any individual agent's safety constraint violation will adversely affect the system's global safety. Specifically, we assume that the system's safety cost adheres to the following Risk-sensitive Individual-Global-Max (RIGM) criterion proposed in (Shen et al. 2023). For any joint safety cost  $\phi_{1-\omega}^{\pi}(C_{i_j}(s(t), a_i(t), \mathbf{a}_{-i}(t)))$ , the following condition exists:

$$\begin{aligned} & \arg \max_{\mathbf{a}} \phi_{1-\omega}^{\pi}[C_{i_j}(s(t), a_i(t), \mathbf{a}_{-i}(t))] \\ & = (\arg \max_{a_1} \phi_{1-\omega}^{\pi}[C_{i_j}(s(t), a_1(t))], \\ & \quad \dots, \arg \max_{a_n} \phi_{1-\omega}^{\pi}[C_{i_j}(s(t), a_n(t))]) \end{aligned} \quad (3)$$

where  $\mathbf{a} = (a_1, \dots, a_n)$  is the joint action.

### 3 Method

This section proposes a worst-case MaxEnt safe RL framework, aiming to address two critical limitations of state-of-the-art safe MARL methods: (1) suboptimal convergence caused by insufficient exploration and (2) the inability to guarantee safety constraint satisfaction during stochastic policy learning. First, we formally derive the worst-case

MaxEnt MARL objective function and characterize its generalized Nash equilibrium solution in Section 3.1. Next, Section 3.2 provides a rigorous theoretical analysis of the proposed framework, focusing on monotonic policy improvement and guarantees of Nash convergence. Finally, Section 3.3 details the specific algorithmic workflow of WC-MASAC.

#### 3.1 Worst-Case Safe MaxEnt MARL

Based on probabilistic graphical inference (Levine 2018), the MaxEnt objective function without safety constraints was successfully derived in (Liu et al. 2023). Building upon (Levine 2018; Liu et al. 2023) and (2), the constrained optimization objective for safe MaxEnt MARL can be formulated as

$$\begin{aligned} \max_{\pi} J(\pi) = & \mathbb{E}_{s_{0:\infty} \sim \mathbb{P}, \mathbf{a}_{0:\infty} \sim \pi} \left[ \sum_{t=1}^{\infty} \gamma^t \left( R(s(t), \mathbf{a}(t)) \right. \right. \\ & \left. \left. + \beta \sum_{i=1}^n \mathcal{H}(\pi_i(\cdot | s(t))) \right) \right] \end{aligned} \quad (4)$$

$$\text{s.t. } \phi_{1-\omega}^{\pi}[C_{i_j}(s, a_i, \mathbf{a}_{-i})] \leq d_{i_j}, \forall i \in \mathcal{N}, 1 \leq j \leq h. \quad (5)$$

where  $\gamma \in (0, 1]$  is the discounted factor.  $\mathcal{H}$  is the entropy term of  $\pi_i$ , and  $\beta$  is the temperature-like constant that regulates the degree of exploration randomness.

When the safety constraint and the entropy term are considered in the objective function, all agents will consider the comprehensive benefit of the safety cost and the reward, assign the mass probability, and eventually converge to the Quantum Response Equilibrium (QRE) (Liu et al. 2023; McKelvey and Palfrey 1995). The following Theorem 1 provides an expression for the equilibrium policy of the game.

**Theorem 1.** For all agents  $i \in \mathcal{N}$ , the system policy converges to the QRE equilibrium policy if no agents can increase the system reward payoff by unilaterally changing their policy. In this case, define the Q-function as

$$\begin{aligned} Q_{\pi}(s(t), \mathbf{a}(t)) = & R(s(t), \mathbf{a}(t)) + \mathbb{E}_{s_{0:\infty} \sim \mathbb{P}, \mathbf{a}_{0:\infty} \sim \pi} \\ & \left[ \sum_{l=0}^{\infty} \gamma^l \left( R(s(t+l), \mathbf{a}(t+l)) + \beta \sum_{i=1}^n \mathcal{H}(\pi_i(\cdot | s(t))) \right) \right] \end{aligned} \quad (6)$$

Then, the QRE equilibrium policy  $\pi_i^*$  of agent  $i$  is expressed as

$$\pi_i^*(a_i | s) = \frac{\Omega(a_i)}{\sum_{b_i \in \mathcal{A}_i} \Omega(b_i)}, \quad (7)$$

in which

$$\begin{aligned} \Omega(a_i) = & \exp \left( \beta^{-1} \left( \mathbb{E} \mathbf{a}_{-i} \sim \pi^{-i} [Q_{\pi}(s, a_i, \mathbf{a}_{-i})] \right. \right. \\ & \left. \left. - \sum_{j=1}^h \frac{\mu^{i_j}}{\omega} \mathbb{E} \mathbf{a}_{-i \sim \pi_i^*} [\hat{\phi}_{1-\omega}^{\pi_i}(C_{i_j}(s, a_i, \mathbf{a}_{-i}))] \right) \right). \end{aligned} \quad (8)$$



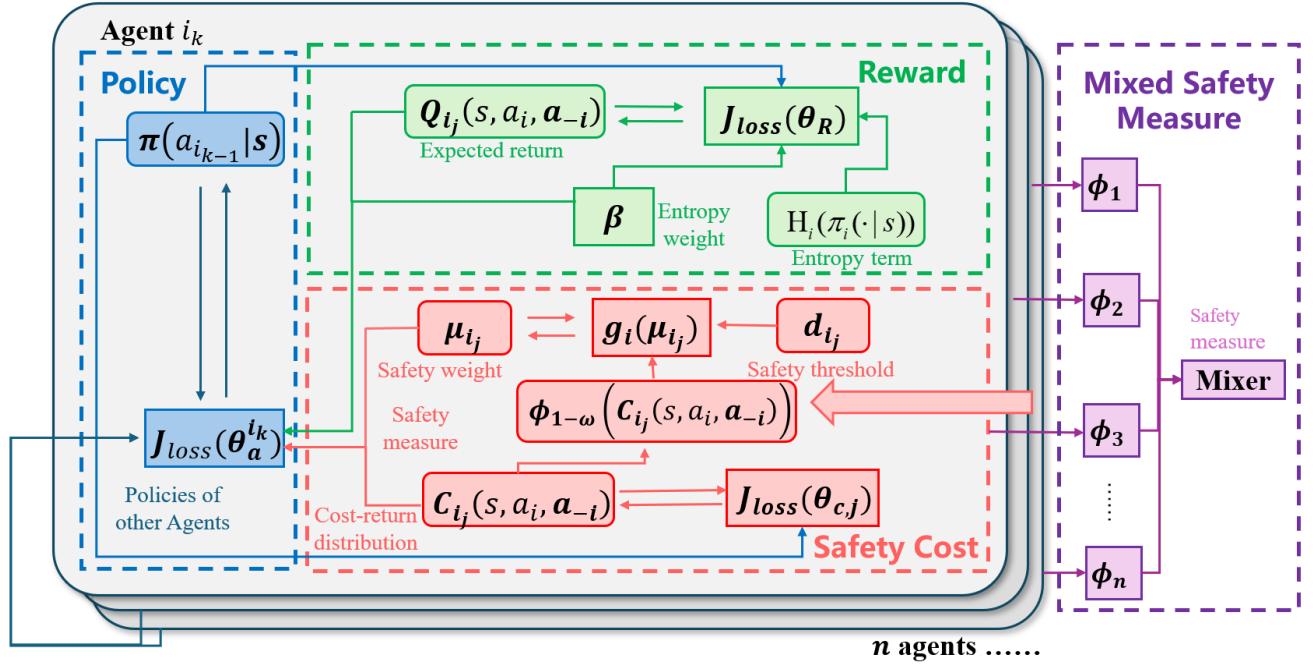


Figure 1: The framework of the WCMASAC.

Then the joint policy update process is

$$\begin{aligned} \pi^{new} = \arg \min_{\pi \in \mathcal{A}} D_{KL} \left( \pi(\cdot | s) \parallel \left( \exp \left( \beta^{-1} \left( Q_{\pi^{old}} \right) \right) \right. \right. \\ \left. \left. (s, \mathbf{a}_{i_{1:k-1}}, a_{i_k}) - \sum_{i=1}^n \sum_{j=1}^n \frac{\mu_{i_j}}{\omega} \hat{\phi}(C_{i_j}(s, \cdot))) \right) \right) \\ \cdot (Z_{\pi^{old}}(s))^{-1} \end{aligned} \quad (14)$$

The proof can be seen in the Appendix B. Theorem 2 states that an individual agent can induce a joint soft policy by sequentially optimizing the KL divergence locally. Based on Theorem 2, the following lemma and theorem theoretically verify the monotonic improvement and QRE convergence of the policy, respectively.

**Lemma 3. (Monotonic Improvement)** Let  $i_{1:k}$  be a permutation. Then for  $\forall k = 1, \dots, n$ , if agent  $i_k$  acquires a new policy  $\pi_{k_i}^{new}$  from  $\pi_{k_i}^{old} \in \mathcal{A}_{i_k}$  through formula (13), then  $Q_{\pi^{new}}(s, a) \geq Q_{\pi^{old}}(s, a)$  will hold and ensure that the safety constraints are satisfied.

Proof can be found in the Appendix C. Lemma 3 indicates that the agents can achieve monotonic improvement in performance while maintaining safety during policy update. Then the following Theorem shows that agents can perform policy evaluation and policy improvement alternately to satisfy the safety constraints and converge to QRE.

**Theorem 3. (QRE Convergence)** If all agents perform policy evaluation and policy improvement alternately, the joint policy  $\pi$  eventually converges to the QRE equilibrium policy and satisfies the safe constraint.

Proof can be found in the Appendix D. Theorem 3 provides a theoretical guarantee that agents monotonically improve their policies under safety constraints and converge to QRE, after which all agents do not unilaterally improve their own policies.

### 3.3 Practical Implementation of WCMASAC

In practical applications, we employ neural network-based approximators to estimate the centralized soft Q-function for rewards and decentralized policies for individual agents, while utilizing a hybrid approach combining quantile regression with a multi-head attention mechanism to approximate the risk measure of the centralized safety cost. Specifically, let  $\theta_r$ ,  $\theta_c$ , and  $\theta_a^{i_k}$  denote the network parameters of the reward critic, cost critic, and  $i_k$ 's actor, respectively. These neural networks are optimized alternately via stochastic gradient descent to ensure coordinated learning of reward maximization, safety constraint satisfaction, and policy decentralization. Let  $s := s(t)$ ,  $\mathbf{a} := \mathbf{a}(t)$  and  $s' = s(t+1)$ ,  $a' = a(t+1)$  for convenience, then the centralized soft Q-function is optimized by minimizing the Bellman residual as

$$\begin{aligned} J_{loss}(\theta_r) = \mathbb{E}_{(s, \mathbf{a}) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_{\theta_r}(s, \mathbf{a}) - (R(s, \mathbf{a}) \right. \right. \\ \left. \left. + \gamma \mathbb{E}_{s' \sim \mathbb{P}, \mathbf{a}' \sim \pi} [Q_{\theta_r}(s', \mathbf{a}') + \beta \sum_{i=1}^n \mathcal{H}(\pi_i(a_i' | s'))]) \right)^2 \right] \end{aligned} \quad (5)$$

Similar to the Quantile Regression (QR) loss (Dabney et al. 2018b), the  $j$ th parameterized joint cost critic  $C_{\theta_{c,j}}$ , which gathers all agents cost  $C_{i_j}$ , achieves parameter optimization by minimizing the Huber loss between the safe cost distri-

bution and the target distribution as

$$J_{loss}(\theta_{c,j}) = H \left( \mathbb{E}_{(s,\mathbf{a}) \sim \mathcal{D}} \left[ C_{\theta_{c,j}}(s, \mathbf{a}) - (c_j + \gamma \mathbb{E}_{s' \sim \mathbb{P}, \mathbf{a}' \sim \pi} [C_{\theta_{c,j}}(s', \mathbf{a}')]) \right] \right) \quad (16)$$

where  $H$  means the Huber loss (Huber 1992) and  $c_j$  is the current joint cost. Then, all the agents minimize the expected KL-divergence (2) in a sequential update manner to improve their policies, which can derive the loss of the  $i_k$ 's actor as follows:

$$J_{loss}(\theta_a^{i_k}) = \mathbb{E}_{(s,\mathbf{a}) \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{a}_{i_1:k-1} \sim \pi_{\theta_a^{i_1:k-1}}, a_{i_k} \sim \pi_{\theta_a^{i_k}}} \left[ Q_{\theta_r}(s, \mathbf{a}_{i_1:k-1}, a_{i_k}) - \lambda_{i_k} \left( \phi_{1-\omega} [C_{\theta_{c,j}}(s, \mathbf{a}_{i_1:k-1}, a_{i_k})] - \alpha \log \pi_{\theta_a^{i_k}}(a_{i_k}^{i_k} | s) \right) \right] \right] \quad (17)$$

where  $\hat{\theta}_a^{i_1:k-1}$  means the improved parameters of agents  $i_1:k-1$ . The framework can be seen in Figure 1, and the procedure as WCMASAC can be seen in the Appendix E.

## 4 Experimental Setting and Results

### 4.1 Experimental Setting

To thoroughly evaluate the efficacy and robustness of our proposed WCMASAC algorithm, we conducted a series of experiments using the Safety-Gymnasium benchmark (Ji et al. 2023) in this section. This benchmark provides several challenging multi-agent environments specifically designed for safe RL research.

We conduct comparative experiments on six velocity-constrained tasks: 2x4Ant, 4x2Ant, 2x3HalfCheetah, 6x1HalfCheetah, 3x1Hopper, and 2x3Walker2D. These experiments evaluate our method against SOTA safe MARL algorithms, including MACPO, MAPPO-Lagrangian (Gu et al. 2023), and the SOTA MARL baseline HASAC (Liu et al. 2023). The tasks simulate scenarios where agents must maximize their movement speed to achieve high rewards while strictly complying with safety constraints, particularly velocity thresholds. This setup rigorously examines the trade-off between reward optimization and safety-critical constraint satisfaction, ensuring agents avoid hazardous velocity violations during operation.

Specifically, the experiments incorporate two key performance metrics, average episode reward and average episode cost, which are evaluated during the testing phase. The first metric quantifies the overall performance of the agents' tasks, while the second measures the severity of violations of the safety constraints per episode. Each agent is required to maximize cumulative rewards while ensuring that the episode cost does not exceed predefined safety thresholds. Furthermore, agents are encouraged to minimize the accumulated cost throughout the operation, thereby balancing reward-seeking behavior with strict adherence to safety-critical constraints. This dual-objective evaluation framework rigorously assesses the trade-off between task efficiency and safety compliance in multi-agent systems.

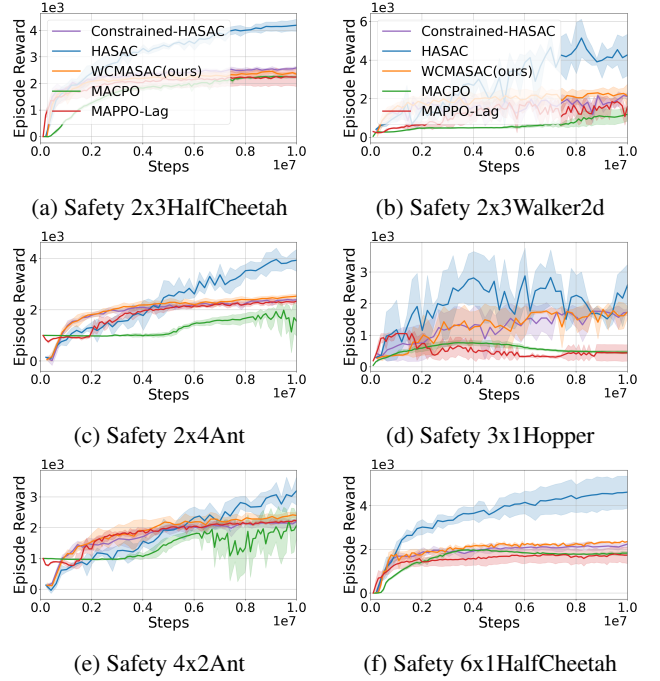


Figure 2: Reward comparisons on selected tasks of multiple benchmarks.

Moreover, Constrained-HASAC is developed to ensure fair comparisons. Unlike WCMASAC, it does not constrain the CVaR. Instead, it directly constrains the joint expectation of all agents.

All algorithms are implemented in a standard framework. Details of network architectures, hyperparameters, training procedures, and specific choices of optimizers, learning rates, and batch sizes are provided in the Appendix F.1.

### 4.2 Experimental Results

**Reward Comparison:** As shown in the comparison of reward with SOTA safe MARL algorithms and the unconstrained SOTA MARL algorithms in Figure 2, our proposed WCMASAC and Constrained-HASAC show superior task reward performance compared to MACPO and MAPPO-Lag in all the six safety environments, though trailing HASAC. Notably, HASAC exhibits multiple risk-prone behaviors due to the absence of safety constraints (see Figure 3). Furthermore, WCMASAC achieves faster learning convergence rates, as its intrinsic exploration-encouraging mechanism enables more efficient exploration of the reward landscape, thereby accelerating the acquisition of optimal behaviors.

**Cost Comparison:** Figure 3 presents the cumulative safety costs and constraint violation statistics of different algorithms under identical environments. First, due to the absence of safety constraints, the SOTA MARL algorithm HASAC exhibits cumulative costs that exceed the safety threshold by a substantial margin and are significantly higher than those of safety-aware MARL baselines. This indicates that HASAC's policy execution compromises system

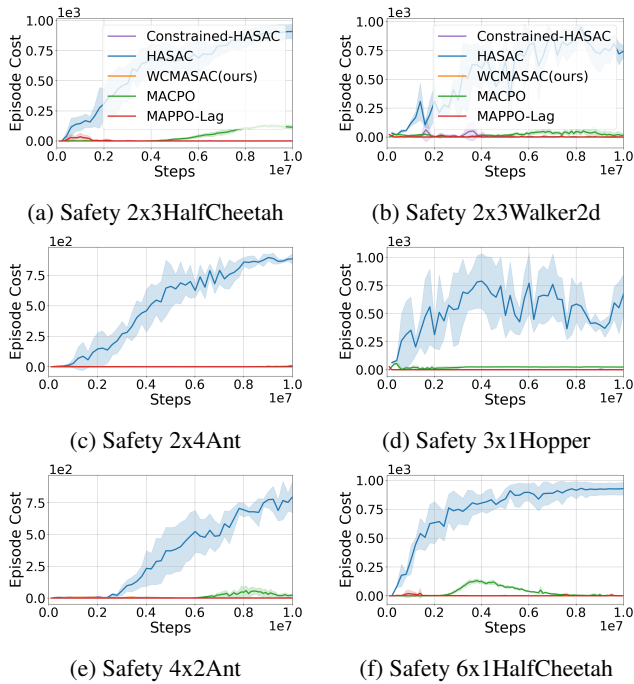


Figure 3: Cost comparisons on selected tasks of multiple benchmarks.

safety through hazardous behaviors, confirming the critical necessity of explicit safety constraints in real-world deployments.

Compared to MACPO, the proposed WCMASAC achieves lower costs in all environments while achieving higher task rewards with reduced execution risks, as evidenced in Figure 2-3. Notably, MACPO violates safety constraints in most environments, except the `Safety2x4AntVelocity` scenario, indicating dangerous agent behaviors during task execution and its failure to guarantee full-process safety. In contrast, WCMASAC and Constrained-HASAC ensure system safety throughout all operational phases in every environment, demonstrating that modeling tail risk distributions of agent costs effectively reduces the probability of selecting high-risk actions.

Figure 4 provides further safety performance comparisons between all the algorithms with safety constrains: MAPPO-Lag, MACPO, Constrained-HASAC and WCMASAC. Overall, WCMASAC and Constrained-HASAC demonstrate comparable safety performance to MAPPO-Lag but with a much lower probability of constraint violations. Specifically, from Figure 4a and Figure 4f we can see that WCMASAC and Constrained-HASAC achieve lower costs than MAPPO-Lag in the `2x3HalfCheetah` and `6x1HalfCheetah` environments, which means that our algorithms have better safety performance. In `2x3Ant` (Figure 4c), `3x1Hopper` (Figure 4d), and `4x2Ant` (Figure 4e) environments, WCMASAC exhibits slightly higher costs than MAPPO-Lag due to exploration encouragement. Notably, in the `2x3HalfCheetah`, `2x3Walker2D`, and `3x1Hopper` environments, MAPPO-Lag occasionally vio-

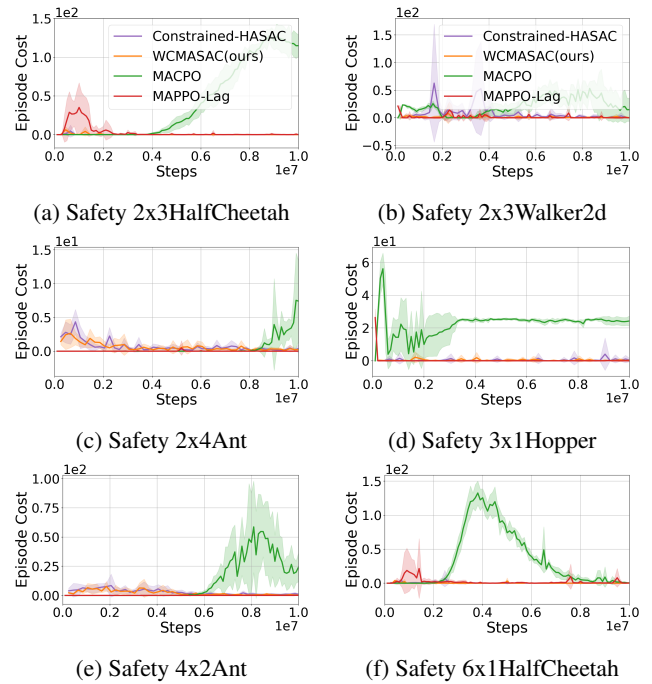


Figure 4: Cost comparisons on selected tasks of multiple safety benchmarks (without HASAC).

lates safety constraints, indicating that agents execute risk-prone behaviors that jeopardize system safety. In contrast, the proposed WCMASAC maintains strict constraint compliance across all environments, validating the effectiveness of the distributional safe critic in mitigating high-risk action selection. Moreover, WCMASAC outperforms Constrained-HASAC in most environments with lower costs in most environments (see Figure 4b, Figure 4c, Figure 4d and Figure 4e), which demonstrates the effectiveness of the CVaR constraint.

More detailed information such as safety thresholds and analysis can be found in the Appendix G.

**Analysis of risk-level parameters:** We investigate the impact of the risk-level parameter  $\omega$  on task rewards and safety costs in the system. As  $1 - \omega$  approaches 1, agents adopt more conservative policies to ensure higher safety levels. We evaluate cumulative rewards and safety performance across all environments under  $1 - \omega$  values of 0.01, 0.5, 0.99. Results reveal a declining trend in overall task rewards as  $1 - \omega$  increases, accompanied by substantial improvements in system safety (see Figure 5 above and Figure 2 in the Appendix G for detailed experiment results). This demonstrates the feasibility of adjusting  $\omega$  based on different task requirements to balance reward maximization and safety assurance.

## 5 Conclusion

This work proposes a MaxEnt safe MARL framework with a safe distributional critic, enabling multi-agent systems to learn stochastic policies for enhanced exploration while guaranteeing safety constraint compliance throughout the

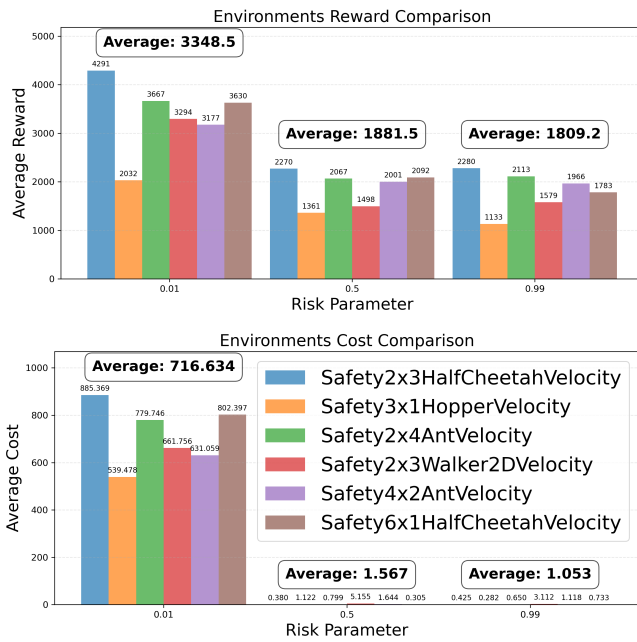


Figure 5: Average reward and average cost under different risk\_level parameter  $\omega$ .

entire operation. The proposed method effectively improves the system’s final task reward and significantly reduces the probability of agents adopting high-risk behaviors. Theoretically, the framework provides closed-form expressions for game-theoretic equilibrium policies through rigorous mathematical proofs, ensuring monotonicity and convergence during policy improvement. Furthermore, we implement the algorithm as WCMASAC and demonstrate its superior performance in both task efficiency and safety compliance within the Safe Multi-Mujoco benchmark environments. Future work will focus on enhancing the algorithm’s scalability in large-scale networks and exploring its deployment in practical multi-agent collaborative scenarios.

## Acknowledgments

This work is supported by in part by the National Natural Science Foundation of China (62333005, 62403200, 62503174, 62473133, 62533002), in part by the Shanghai International Science and Technology Cooperation Project (24510714000), in part by the Shanghai Natural Science Foundation (24ZR1416200), in part by the China University-Industry-Research Innovation Funds (2024IT032), in part by the opening project of the State Key Laboratory of Autonomous Intelligent Unmanned Systems (ZZKF2025-1-24), and in part by the Fundamental Research Funds for the Central Universities.

## References

Dabney, W.; Ostrovski, G.; Silver, D.; and Munos, R. 2018a. Implicit quantile networks for distributional reinforcement

learning. In *International conference on machine learning*, 1096–1105. PMLR.

Dabney, W.; Rowland, M.; Bellemare, M.; and Munos, R. 2018b. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Du, H.; Gou, F.; and Cai, Y. 2025. Scalable Safe Multi-Agent Reinforcement Learning for Multi-Agent System. *arXiv preprint arXiv:2501.13727*.

Gu, S.; Kuba, J. G.; Chen, Y.; Du, Y.; Yang, L.; Knoll, A.; and Yang, Y. 2023. Safe multi-agent reinforcement learning for multi-robot control. *Artificial Intelligence*, 319: 103905.

Gu, S.; Yang, L.; Du, Y.; Chen, G.; Walter, F.; Wang, J.; and Knoll, A. 2022. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*.

Gu, S.; Yang, L.; Du, Y.; Chen, G.; Walter, F.; Wang, J.; and Knoll, A. 2024. A Review of Safe Reinforcement Learning: Methods, Theories, and Applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12): 11216–11235.

Ha, S.; Xu, P.; Tan, Z.; Levine, S.; and Tan, J. 2021. Learning to Walk in the Real World with Minimal Human Effort. In *Conference on Robot Learning*, 1110–1120. PMLR.

Haarnoja, T.; Zhou, A.; Hartikainen, K.; Tucker, G.; Ha, S.; Tan, J.; Kumar, V.; Zhu, H.; Gupta, A.; Abbeel, P.; et al. 2018. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*.

Hu, J.; Sun, Y.; Chen, H.; Huang, S.; Chang, Y.; Sun, L.; et al. 2022. Distributional reward estimation for effective multi-agent deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35: 12619–12632.

Hu, T.; Luo, B.; Yang, C.; and Huang, T. 2023. MO-MIX: Multi-Objective Multi-Agent Cooperative Decision-Making With Deep Reinforcement Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10): 12098–12112.

Huber, P. J. 1992. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodology and distribution*, 492–518. Springer.

Ji, J.; Zhang, B.; Zhou, J.; Pan, X.; Huang, W.; Sun, R.; Geng, Y.; Zhong, Y.; Dai, J.; and Yang, Y. 2023. Safety gymnasium: A unified safe reinforcement learning benchmark. *Advances in Neural Information Processing Systems*, 36: 18964–18993.

Kuba, J. G.; Feng, X.; Ding, S.; Dong, H.; Wang, J.; and Yang, Y. 2022. Heterogeneous-agent mirror learning: A continuum of solutions to cooperative marl. *arXiv preprint arXiv:2208.01682*.

Kuba, J. G.; Wen, M.; Meng, L.; Zhang, H.; Mguni, D.; Wang, J.; Yang, Y.; et al. 2021. Settling the variance of multi-agent policy gradients. *Advances in Neural Information Processing Systems*, 34: 13458–13470.

Levine, S. 2018. Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review. *arXiv:1805.00909*.

- Li, Z.; Hu, C.; Wang, Y.; Yang, Y.; and Li, S. E. 2024. Safe Reinforcement Learning With Dual Robustness. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12): 10876–10890.
- Liu, C.; Geng, N.; Aggarwal, V.; Lan, T.; Yang, Y.; and Xu, M. 2021. Cmix: Deep multi-agent reinforcement learning with peak and average constraints. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part I 21*, 157–173. Springer.
- Liu, J.; Zhong, Y.; Hu, S.; Fu, H.; Fu, Q.; Chang, X.; and Yang, Y. 2023. Maximum entropy heterogeneous-agent reinforcement learning. *arXiv preprint arXiv:2306.10715*.
- Liu, Z.; Cen, Z.; Isenbaev, V.; Liu, W.; Wu, S.; Li, B.; and Zhao, D. 2022. Constrained variational policy optimization for safe reinforcement learning. In *International Conference on Machine Learning*, 13644–13668. PMLR.
- Lu, S.; Zhang, K.; Chen, T.; Başar, T.; and Horesh, L. 2021. Decentralized policy gradient descent ascent for safe multi-agent reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, 8767–8775.
- McKelvey, R. D.; and Palfrey, T. R. 1995. Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1): 6–38.
- Qiu, W.; Wang, X.; Yu, R.; Wang, R.; He, X.; An, B.; Obraztsova, S.; and Rabinovich, Z. 2021. RMIX: Learning risk-sensitive policies for cooperative reinforcement learning agents. *Advances in Neural Information Processing Systems*, 34: 23049–23062.
- Rashid, T.; Farquhar, G.; Peng, B.; and Whiteson, S. 2020a. Weighted QMIX: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 33: 10199–10210.
- Rashid, T.; Samvelyan, M.; De Witt, C. S.; Farquhar, G.; Foerster, J.; and Whiteson, S. 2020b. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178): 1–51.
- Rockafellar, R. T.; Uryasev, S.; et al. 2000. Optimization of conditional value-at-risk. *Journal of risk*, 2: 21–42.
- Shen, S.; Ma, C.; Li, C.; Liu, W.; Fu, Y.; Mei, S.; Liu, X.; and Wang, C. 2023. RiskQ: risk-sensitive multi-agent reinforcement learning value factorization. *Advances in Neural Information Processing Systems*, 36: 34791–34825.
- Sun, W.-F.; Lee, C.-K.; See, S.; and Lee, C.-Y. 2023. A unified framework for factorizing distributional value functions for multi-agent reinforcement learning. *Journal of Machine Learning Research*, 24(220): 1–32.
- Sutton, R. S.; Barto, A. G.; et al. 1998. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge.
- Tang, Y. C.; Zhang, J.; and Salakhutdinov, R. 2019. Worst cases policy gradients. *arXiv preprint arXiv:1911.03618*.
- Yang, Q.; Simão, T. D.; Tindemans, S. H.; and Spaan, M. T. 2021. WCSAC: Worst-case soft actor critic for safety-constrained reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 10639–10646.
- Yang, Q.; Simão, T. D.; Tindemans, S. H.; and Spaan, M. T. 2023. Safety-constrained reinforcement learning with a distributional safety critic. *Machine Learning*, 112(3): 859–887.
- Zhang, K.; Sun, T.; Tao, Y.; Genc, S.; Mallya, S.; and Basar, T. 2020. Robust multi-agent reinforcement learning with model uncertainty. *Advances in neural information processing systems*, 33: 10571–10583.
- Zhang, L.; Li, L.; Wei, W.; Song, H.; Yang, Y.; and Liang, J. 2024. Scalable Constrained Policy Optimization for Safe Multi-agent Reinforcement Learning. *Advances in Neural Information Processing Systems*, 37: 138698–138730.
- Zhang, T.; Li, Y.; Wang, C.; Xie, G.; and Lu, Z. 2021. Fop: Factorizing optimal joint policy of maximum-entropy multi-agent reinforcement learning. In *International conference on machine learning*, 12491–12500. PMLR.
- Zhang, Y.; Vuong, Q.; and Ross, K. 2020. First order constrained optimization in policy space. *Advances in Neural Information Processing Systems*, 33: 15338–15349.
- Zhong, Y.; Kuba, J. G.; Feng, X.; Hu, S.; Ji, J.; and Yang, Y. 2024. Heterogeneous-agent reinforcement learning. *Journal of Machine Learning Research*, 25(32): 1–67.