SAFE MULTI-OBJECTIVE REINFORCEMENT LEARN-ING VIA MULTI-PARTY PARETO NEGOTIATION

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

033 034

035

036

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Safe multi-objective reinforcement learning (Safe MORL) seeks to optimize performance while satisfying safety constraints. Existing methods face two key challenges: (i) incorporating safety as additional objectives enlarges the objective space, requiring more solutions to uniformly cover the Pareto front and maintain adaptability under changing preferences; (ii) strictly enforcing safety constraints is feasible for single or compatible constraints, but conflicting constraints prevent flexible, preference-aware trade-offs. To address these challenges, we cast Safe MORL within a multi-party negotiation framework that treats safety as an external regulatory perspective, enabling the search for a consensus-based multi-party Pareto-optimal set. We propose a multi-party Pareto negotiation (MPPN) strategy built on NSGA-II, which employs a negotiation threshold ε to represent the acceptable solution range for each party. During evolutionary search, ε is dynamically adjusted to maintain a sufficiently large negotiated solution set, progressively steering the population toward the $(\varepsilon_{\text{efficiency}}, \varepsilon_{\text{safety}})$ -negotiated common Pareto set. The framework preserves user preferences over conflicting safety constraints without introducing additional objectives and flexibly adapts to emergent scenarios through progressively guided ($\varepsilon_{\text{efficiency}}, \varepsilon_{\text{safety}}$). Experiments on a MuJoCo benchmark show that our approach outperforms state-of-the-art methods in both constrained and unconstrained MORL, as measured by multi-party hypervolume and sparsity metrics, while supporting preference-aware policy selection across stakeholders.

1 Introduction

Multi-objective reinforcement learning (MORL) addresses decision-making problems with multiple, often conflicting objectives (Dulac-Arnold et al., 2021). Since no single policy can be optimal across all objectives simultaneously, existing approaches are typically divided into two categories. Single-policy MORL (Chen et al., 2021; Skalse et al., 2022; Kyriakis & Deshmukh, 2022) reduces the multi-objective problem to a scalar one by applying predefined weights, allowing standard RL algorithms to be used directly. However, this scalarization produces a policy tailored to a fixed preference, limiting adaptability across tasks. Multi-policy MORL (Yang et al., 2019; Chen et al., 2019; Xu et al., 2020; Hayes et al., 2022), on the other hand, aims to approximate the Pareto front (PF) by learning a set of non-dominated policies, thereby supporting diverse objective preferences and enabling flexible policy selection in practice.

While MORL has shown promising progress, incorporating safety considerations introduces new challenges. Safe MORL (Huang et al., 2022) aims to optimize multiple objectives while ensuring that agents adhere to safety requirements, thereby preventing hazardous behaviors during training and deployment. One natural formulation is to treat safety as an additional objective alongside performance goals. This enables explicit exploration of safety—performance trade-offs but increases the dimensionality of the objective space, leading to an exponential growth of the Pareto set and making full coverage intractable. Alternatively, safety can be enforced as a hard constraint, giving rise to constrained MORL (CMORL) (Huang et al., 2022; Lin et al., 2024; Gu et al., 2025), which restricts the search to policies that satisfy predefined safety conditions. This avoids dimensionality explosion and directly guarantees safe behavior, but struggles with conflicting or overly strict constraints and lacks adaptability in dynamic environments.

In practice, however, safety is not always absolute. Real-world decision making often requires negotiating between efficiency and safety, where tolerating minor violations of certain safety constraints may yield significant performance gains. For example, as illustrated in Figure 1, cargo-handling robots aim to maximize movement speed and payload capacity while maintaining body stability and limiting energy consumption. Allowing slight violations in stability or energy usage can enable faster transport of larger loads, which may be desirable in time-sensitive scenarios. Such flexibility is difficult to achieve with CMORL, since hard constraints restrict the feasible solution set and often eliminate practically useful trade-offs. Likewise, objective-based MORL approaches may struggle to capture balanced solutions between efficiency and safety when the inclusion of multiple safety objectives causes the Pareto set to expand excessively.

To address these limitations, we reconceptualize Safe MORL as a multi-party negotiation problem, where the safety objectives and efficiency objectives are treated as separate multi-objective decision parties rather than as additional objectives in a single objective space. This formulation enables the search for a common Pareto set that balances efficiency and safety while resolving potential conflicts among safety constraints. Building on this idea, we develop a negotiation-driven evolutionary framework, MPPN-MORL, which integrates multi-party Pareto negotiation into policy search without increasing the dimensionality of the objective space. Our approach flexi-

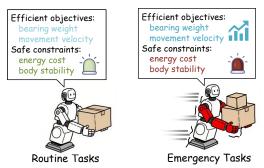


Figure 1: Trade-off between efficiency and safety. Safety constraints are relaxed in emergency tasks to prioritize efficiency objectives.

bly adapts to user-specified preferences over both performance and safety, preserves diversity in the solution set, and promotes fairness across parties. Extensive experiments on a multi-objective Mu-JoCo benchmark demonstrate that MPPN-MORL achieves superior trade-offs between efficiency and safety compared to existing MORL and CMORL methods, while effectively handling conflicting safety constraints and supporting preference-aware policy deployment.

2 PRELIMINARIES

2.1 MULTI-OBJECTIVE DECISION-MAKING

A multi-objective decision-making (MODM) problem involves optimizing multiple, potentially conflicting objectives. Formally, it can be formulated as

$$\min_{\pi_{\theta}} \mathbf{F}_{\pi_{\theta}} = \min_{\pi_{\theta}} \left[f_1(\pi_{\theta}), \dots, f_m(\pi_{\theta}) \right], \tag{1}$$

where π_{θ} denotes a parameterized policy, and $f_i(\pi_{\theta})$ is the expected performance of π_{θ} with respect to the *i*-th objective, for $i=1,\ldots,m$. Unlike single-objective settings, which seek a unique optimal policy, MODM problems typically yield a set of Pareto-optimal solutions, each reflecting a different trade-off among objectives.

2.2 Constrained Multi-objective Reinforcement Learning

Safe MORL extends MODM by enforcing safety constraints, restricting the set of admissible policies. Constrained MORL (CMORL) formalizes this idea using explicit cost functions. Specifically, CMORL introduces p additional cost functions c_{m+1},\ldots,c_{m+p} , each mapping a state-action pair (s,a) to a scalar cost. For a policy π , the expected cumulative cost under the (m+i)-th function is denoted as $c_{m+i}(\pi)$, which must satisfy

$$c_{m+i}(\pi) \le d_i, \quad \forall i = 1, \dots, p, \tag{2}$$

where d_i is a predefined safety threshold. The objective in CMORL is to optimize the vector-valued function $\mathbf{F}(\pi)$ representing performance across m objectives, while ensuring that π lies within the safe policy set:

$$\Pi_{\text{safe}} = \{ \pi \in \Pi \mid c_i(\pi) \le d_i, \forall i = m+1, \dots, m+p \},$$
(3)

from which the agent identifies a set of Pareto-optimal policies.

Gu et al. (2025) extend the Pareto frontier concept to safety-constrained MDPs. A policy $\pi \in \Pi_{safe}$ is safe Pareto-optimal if no other policy in Π_{safe} strictly improves all objectives without violating any constraint. Thus, the central goal of CMORL is to efficiently find such policies, balancing objective performance with constraint satisfaction.

3 Метнор

In this section, we introduce the modeling approach for MPMORL and present the MPPN-MORL algorithm, which is capable of selecting appropriate Pareto-optimal policy sets based on the preferences of multiple parties. We first describe the modeling framework for MPMORL, and then elaborate on MPPN-MORL from two key aspects. The detail of MPPN-MORL is shown in Algorithm 2.

3.1 MULTI-PARTY MULTI-OBJECTIVE DECISION-MAKING

Multi-party multi-objective decision-making (MPMODM) models scenarios with multiple decision-makers (DMs), where each DM optimizes its own set of objectives and at least one DM faces multiple, potentially conflicting goals. Such scenarios arise naturally in decentralized systems, multidepartmental planning, and cooperative multi-agent environments where each party holds distinct priorities.

In sequential decision-making under uncertainty, MPMODM can be formulated as a multi-party multi-objective Markov decision process (MPMOMDP). In this work, we focus on two parties: the *safety side* and the *efficiency side*. Formally, the problem is defined by the tuple

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, T, \mu, \Gamma, \mathcal{R} \rangle \tag{4}$$

where S is the state space, \mathcal{A} is the action space, $T(s'\mid s,a)$ denotes the state transition probability, and μ is the initial state distribution. The discount factors are represented by $\Gamma=\gamma^1,\gamma^2$, where $\gamma^k=[\gamma_1^k,\ldots,\gamma_m^k]\in[0,1]^m$ denotes the discount vector for party k. The reward function is defined as

$$\mathcal{R}(s, a, s') = \left[R^1(s, a, s'), R^2(s, a, s') \right], \tag{5}$$

with $R^k(s, a, s') = [r_1^k(s, a, s'), \dots, r_m^k(s, a, s')]^{\top}$ representing the m-dimensional reward vector for party k.

A policy $\pi_{\theta}: \mathcal{S} \to \mathcal{A}$ guides the agent's behavior. For each party $k \in \{1, 2\}$, its performance is evaluated using a vector of expected discounted returns:

$$J_{i,\pi_{\theta}}^{k} = \mathbf{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{\infty} (\gamma_{i}^{k})^{t} r_{i}^{k}(s_{t}, a_{t}, s_{t+1}) \right], \quad i = 1, \dots, m.$$
 (6)

Let $\mathbf{J}_{\pi_{\theta}}^{k} = [J_{1,\pi_{\theta}}^{k},\ldots,J_{m,\pi_{\theta}}^{k}]^{\top}$ be the m-dimensional return vector for party k. The goal of MP-MORL is to identify a set of policies $\{\pi_{\theta}\}$ that approximates the joint two-party Pareto front, which balances trade-offs between the safety side and the efficiency side.

Definition 3.1 (Pareto Dominance). Given two solutions $X, Y \in \mathcal{X}$ and the objective set F_k of the k-th DM, X is said to *Pareto dominate* Y with respect to DM k, denoted as $X \prec_k Y$, if $f_{k,i}(X) \leq f_{k,i}(Y)$ for all $i \in \{1, \ldots, m_k\}$ and there exists at least one j such that $f_{k,j}(X) < f_{k,j}(Y)$.

Definition 3.2 (Multi-Party Pareto Dominance). Given two solutions $X,Y \in \mathcal{X}$, X is said to *multi-party Pareto dominate* Y, denoted as $X \prec_{\mathsf{MP}} Y$, if $X \prec_k Y$ holds in the local objective space of each DM k.

Different from traditional MORL, where Pareto optimality is defined with respect to a centralized objective space, MPMORL introduce a perspective-dependent notion of optimality. A solution regarded as globally Pareto optimal may appear suboptimal from the standpoint of an individual DM with unique preferences. MPMORL requires identifying solutions that not only balance multiple objectives but also ensure diversity and fairness among all participating DMs.

Example: Consider a robotic cargo transportation task in which a robot delivers goods from a workstation to a designated target area. In this scenario, two DMs focus on different aspects of the robot's policy: the efficiency party emphasizes transportation speed and payload capacity, whereas the safety party prioritizes energy consumption and body stability.

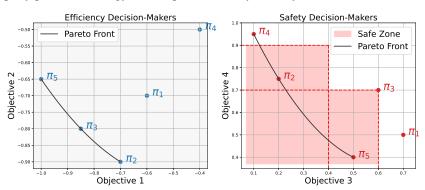


Figure 2: Performance of different policies under MPMORL modeling and CMORL modeling.

As shown in Figure 2, when modeled as a conventional MORL problem, each policy constitutes a Pareto-optimal solution. However, individual DMs may have diverging preferences: the efficiency party favors polices π_1 , π_2 and π_5 , while the safety party favors polices π_2 , π_3 , and π_4 . Although all policies are Pareto-optimal in the multi-objective sense, under the multi-party perspective, policies π_1 , π_3 and π_4 are dominated by policies π_2 and π_5 . Therefore, the multi-party Pareto front includes π_2 and π_5 , eliminating solutions that appear optimal in the centralized view.

When facing conflicting safety constraints, CMORL must perform optimization within a fixed constraint space, thus failing to find all Pareto solutions. While feasible, it lacks the diversity required for negotiation among DMs and cannot effectively resolve conflicts between objectives.

3.2 Multi-Party Pareto Negotiation-Based Non-Dominated Sorting

In MPMODM, the goal is generally to identify a common Pareto set that satisfies all parties. However, such common solutions are often limited. To enlarge the set of negotiable policies while respecting individual preferences, one or both parties may relax their acceptance criteria. We model this process as a bargaining game, where both parties start from an initial compromise level ε and iteratively negotiate toward their reference thresholds ($\varepsilon_{\rm efficiency}, \varepsilon_{\rm safety}$). During this negotiation, the framework ensures that solutions maintain both high quality and uniform coverage, providing a balanced compromise between efficiency and safety.

To address multi-objective decision-making in multi-party scenarios, we propose a *multi-party Pareto negotiation* (MPPN) framework that extends the classical Pareto concept. The detailed procedure of this algorithm is presented in Algorithm 1. The key idea is to relax the strict dominance relation by introducing an ε -compromise degree with respect to a shared reference solution, allowing each DM to accept solutions that are not strictly superior but still fall within an acceptable margin of improvement relative to this baseline.

Specifically, for each DM, ε -dominance is evaluated against a predefined reference reward vector \mathbf{r}_{ref} . A candidate solution \mathbf{r} is said to ε -dominate the reference if, within the DM's objective subspace, it satisfies $\mathbf{r} \succeq_{\varepsilon} \mathbf{r}_{ref}$, meaning that its performance is no worse than $\varepsilon \cdot \mathbf{r}_{ref}$ across all relevant objectives (accounting for optimization direction via element-wise scaling) and strictly better in at least one. This mechanism prevents excessive rejection of solutions due to minor differences and provides a negotiation margin centered on a common target, enabling more practical and stable preference modeling under conflicting interests. As shown in Figure 3, by reducing the value of ($\varepsilon_{\text{efficiency}}, \varepsilon_{\text{safety}}$), safety constraints can be appropriately relaxed in exchange for tightening policies toward higher-performance regions.

To retain appropriate policies during the evolutionary process, we have designed a multi-party non-dominated sorting (MPNDS) method under ε -dominance. This MPNDS approach constructs a ranking dictionary for each DM, which records the Pareto front level of each individual from the perspective of that DM. For ε -dominance, the Pareto level is defined as the minimum policy value that can

217218219220221222

223224225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240241242

243

244

245

246

247

248249

250

251

252

253

254

255

256257258

259260

261

262

263

264265

266

267

268

269

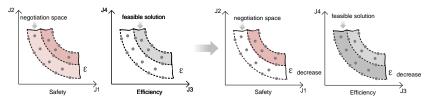


Figure 3: The reward space under ε -dominance varies in accordance with the outcomes of the negotiation process.

Algorithm 1 Multi-party ε -dominance Sorting

```
Input: Candidate set \mathcal{C} with rewards \{\mathbf{r}_i\}, reference solution \mathbf{r}_{ref}, compromise vector \varepsilon = [\varepsilon_1, \varepsilon_2],
objective partitions \{J_1, J_2\}, joint threshold \tau
Output: Joint \varepsilon-front \mathcal{F}_{\text{joint}}, local fronts \{\mathcal{F}_1, \mathcal{F}_2\}
  1: Initialize \mathcal{F}_{\mathrm{joint}} \leftarrow \emptyset, \mathcal{F}_1 \leftarrow \emptyset, \mathcal{F}_2 \leftarrow \emptyset
  2: for each solution \mathbf{r}_i \in \mathcal{C} do
  3:
             if \mathbf{r}_i \in \mathcal{E}-dominates \mathbf{r}_{ref} w.r.t. J_1 and J_2 then
  4:
                  Add \mathbf{r}_i to \mathcal{F}_{\text{joint}}
  5:
             else if \mathbf{r}_i \in \text{-dominates } \mathbf{r}_{\text{ref}} w.r.t. J_1 or J_2 then
                  Add \mathbf{r}_i to \mathcal{F}_1 or \mathcal{F}_2
  7:
             end if
  8: end for
  9: if |\mathcal{F}_{\text{joint}}| \geq \tau then
             Shrink \varepsilon \leftarrow \max(\varepsilon * \Delta \varepsilon, \varepsilon_{\min})
10:
11: end if
12: return \mathcal{F}_{\text{joint}}, \mathcal{F}_1, \mathcal{F}_2
```

 ε -dominate the reference policy. We then calculate the sum of the levels of each individual across all DMs, and stratify and sort the individuals in ascending order based on the aggregated total level value. This strategy effectively mitigates the dominant impact of overly strict preferences from a single DM on the overall ranking, and better integrates the perspectives of all participating parties. By adopting a level summation mechanism to aggregate multi-party negotiation opinions, MPPN can identify solutions that more fairly reflect the diverse and even potentially conflicting objectives of various stakeholders.

By dynamically tightening ε , the algorithm transitions from broad exploration to focused exploitation. Initially large to allow diverse policies near \mathbf{r}_{ref} , ε decays only when enough solutions satisfy the condition for all DMs. If one party fails to improve, ε for the other party is temporarily relaxed to explore better policies.

Overall, the multi-party ε -nondominated sorting first identifies locally preferred solutions within each DM and then integrates them to form a global ranking. This yields a set of ε -nondominated solutions that balance conflicting objectives while maintaining diversity and fairness.

3.3 MULTI-PARTY PARETO NEGOTIATION FOR SAFE MORL

To address conflicts among objectives from multiple parties, we propose MPPN-MORL, which incorporates a multi-party negotiation mechanism into evolutionary search. Inspired by NSGA-II(Storn & Price, 1997), the algorithm replaces genetic operators with differential evolution for efficiency and substitutes standard Pareto dominance with an ε -dominance criterion to enable negotiation.

MPPN-MORL initializes a population of candidate policies, evaluates their multi-objective rewards, and assigns party-specific compromise parameters ε to guide negotiation. In each generation, offspring are generated via differential evolution and combined with parents. Solutions are compared against a reference under the negotiation-based dominance criterion. Jointly dominant solutions tighten ε to enforce stricter optimality, while if none exist, each party updates its own dominant set, preserving individual preferences.

Population diversity is maintained by prioritizing ε -dominated solutions and filling remaining slots based on crowding distance. This iterative process continues until termination, yielding a negotiation-based Pareto front that balances cooperation and competition, reflecting both individual preferences and mutual consensus.

4 EXPERIMENTS

4.1 EVALUATION METRICS

In MORL, the most commonly used evaluation metric is the hypervolume (HV) and Sparsity (SP) (Xu et al., 2020; Basaklar et al., 2023; Hu & Luo, 2024; Liu et al., 2025). To evaluate the overall performance in MPMORL scenarios, we employ the Multi-Party Hypervolume (MPHV) and Multi-Party Sparsity (MPSP). Assume that for each DM, an approximated Pareto front L_k is obtained in an m_k -dimensional objective space, where $k \in \{1, \ldots, K\}$ indexes the parties and M_k is the number of solutions in L_k . Let $r_k \in \mathbb{R}^{m_k}$ be the reference point for the k-th party. The HV for L_k is defined as:

$$HV(L_k) = \delta(H(L_k, r_k)), \tag{7}$$

where

$$H(L_k, r_k) = \{ w \in \mathbb{R}^{m_k} \mid \exists j, \ r_k \le w \le L_{k,j} \},$$
(8)

 $L_{k,j}$ is the j-th solution in L_k , and $\delta(\cdot)$ denotes the Lebesgue measure in \mathbb{R}^{m_k} . The relation \preceq is the weak Pareto dominance operator, meaning that for two vectors $a,b\in\mathbb{R}^{m_k}$, $a\preceq b$ holds if and only if $a_i\leq b_i$ for all objectives i. HV measures the volume of the region dominated by the approximated Pareto set L_k and bounded by the reference point r_k , where a larger HV indicates better convergence and diversity properties of the approximation.

By introducing the negotiation thresholds ($\varepsilon_{\text{efficiency}}$, $\varepsilon_{\text{safety}}$), MPHV aggregates the HV of all parties with preference weights, reflecting the overall performance of the approximated Pareto sets. Its calculation formula is expressed as follows:

$$MPHV = (1 - \varepsilon_{\text{efficiency}}) \cdot HV(L_{\text{efficiency}}) + (1 - \varepsilon_{\text{safety}}) \cdot HV(L_{\text{safety}}). \tag{9}$$

The SP metric is further introduced to evaluate the distribution of solutions along the approximated Pareto front. Unlike HV, which focuses on convergence and overall coverage of the objective space, SP emphasizes the evenness of solution spacing, reflecting how well the algorithm maintains diversity across objectives. Formally, let $L = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$ be the approximated Pareto front in an m-dimensional objective space, where M is the number of solutions. For each objective dimension $k \in \{1, \dots, m\}$, the solutions are sorted in descending order by their k-th objective value. The sparsity is then computed as:

$$SP(L) = \frac{1}{M-1} \sum_{k=1}^{m} \sum_{j=1}^{M-1} (z_{j,k} - z_{j+1,k})^2,$$
(10)

where $z_{j,k}$ denotes the k-th objective value of the j-th solution after sorting. A lower SP value indicates that the solutions are more evenly distributed along the Pareto front. Therefore, SP serves as a complementary indicator to HV, as it directly measures the diversity of solutions rather than the dominated volume.

Analogous to MPHV, we extend SP to the multi-party setting by defining the Multi-Party Sparsity (MPSP). Specifically, MPSP aggregates the sparsity values of all parties under negotiation thresholds, capturing the overall evenness of solution distribution across different parties. Its formulation is given as:

$$MPSP = (1 - \varepsilon_{\text{efficiency}}) \cdot SP(L_{\text{efficiency}}) + (1 - \varepsilon_{\text{safety}}) \cdot SP(L_{\text{safety}}). \tag{11}$$

This metric reflects the overall quality of the Pareto approximations across all parties. A higher MPHV indicates that the solutions perform well on average for individual parties, maintaining good convergence and diversity. In contrast, a lower MPSP value signifies that the algorithm achieves a well-spread set of solutions for each party and avoids clustering or large gaps between adjacent solutions.

4.2 Environment Settings

Based on the MuJoCo (Todorov et al., 2012) and MO-MuJoCo (Xu et al., 2020) benchmark, we developed a MPMO MuJoCo benchmark to evaluate the performance of the proposed algorithms within the MuJoCo framework. This benchmark consists of six continuous robotic locomotion control tasks: mp-HalfCheetah, mp-Walker, mp-Hopper, mp-Pusher, mp-Swimmer, and mp-Humanoid. Each task involves two decision-making parties, namely the safety party and the efficiency party, where each party is associated with two distinct objectives.

The definitions of objectives and reward formulations for the other experimental environments are detailed in the supplementary material.

4.3 BASELINES

To demonstrate the advantages of the MPMORL formulation and to evaluate the effectiveness of the proposed MPPN-MORL algorithm, we conducted experiments against leading methods from both domains. For CMORL, we first adopted LP3 (Huang et al., 2022) as a baseline algorithm. We also adopted the state-of-the-art algorithm CR-MOPO (Gu et al., 2025). For MORL, we employed PGMORL (Xu et al., 2020), advanced approach designed for continuous state-action spaces. Notably, CR-MOPO-S (Gu et al., 2025), which reformulates the safety constraint in CR-MOPO as an additional objective, can also be viewed as a MORL algorithm. Furthermore, to evaluate the ef-

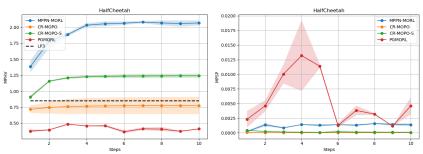


Figure 4: The MPHV and MPSP curve for the MP-HalfCheetah environment. The shaded region represents the standard deviation across six independent experimental runs.

fectiveness of the MPPN-MORL algorithm, we performed an ablation study in which the MPPN component is removed, and only the MPNDS (Liu et al., 2020) component is employed during the evolutionary process; this variant is referred to as MPPN-ablated.

Further details regarding the algorithmic procedures and parameter settings of the baseline methods are provided in the supplementary material.

4.4 RESULTS

We evaluate the proposed methods on the developed MPMO continuous control benchmark, MPMO-MuJoCo. Figure 4 illustrates the MPHV and MPSP curves during training for all methods in the MP-HalfCheetah environment, while Table 1 reports the evaluation results across all environments. The MPPN-MORL algorithm employs an initial negotiation vector of (0.5,0.5).

On the MPHV metric, MPNN-MORL performs slightly worse than CR-MOPO only in the mp-Hopper environment, while achieving the best performance across all other environments. This result validates the effectiveness of the proposed algorithm in balancing the interests of multiple parties. However, on the MPSP metric, the CR-MOPO-S algorithm achieves the best performance in three environments, and the CR-MOPO algorithm achieves the best performance in two environments, which can be attributed to their gradient-based optimization that yields a large number of dense policies. It is worth noting that MPNN-MORL exhibits relatively weaker performance on the MPSP metric. This phenomenon is mainly attributed to the negotiation mechanism of the algorithm, which places greater emphasis on global convergence during optimization, resulting in a sparser distribution of solutions along the Pareto front and consequently higher MPSP values. This indicates that MPNN-MORL has certain limitations in terms of solution distribution.

Table 1: Experimental results of MP-MuJoCo environments. Each algorithm was independently executed six times under identical experimental conditions, reporting the mean \pm standard deviation.LP3 outputs a single policy and thus cannot calculate MPSP.

Environments	Metrics	PGMORL	LP3	CR-MOPO	CR-MOPO-S	MPNN-MORL
MP-HalfCheetah-v4	MPHV MPSP(10 ⁻²)	0.411±0.006 0.458±0.146	0.852 N/A	0.778±0.134 0.001 ± 0.001	1.241±0.032 0.007±0.008	2.067 ± 0.040 0.137±0.021
MP-Walker-v4	MPHV	0.000±0.000	0.000	1.518±0.013	0.294±0.037	1.797±0.220
	MPSP(10 ⁻²)	3.136±1.343	N/A	0.195±0.019	0.246±0.070	0.102±0.056
MP-Hopper-v4	MPHV	0.273±0.273	0.000	1.235±0.034	1.376±0.094	1.190±0.006
	MPSP(10 ⁻²)	0.828±0.488	N/A	0.012±0.018	0.191±0.032	0.670±0.050
MP-Pusher-v4	MPHV MPSP(10 ⁻²)	0.142±0.029 6.093±2.731	0.063 N/A	0.398±0.041 0.015±0.002	0.753 ± 0.062 0.014 ± 0.003	$\begin{array}{c} 0.824 {\pm} 0.001 \\ 0.004 {\pm} 0.003 \end{array}$
MP-Swimmer-v4	MPHV	0.011±0.000	0.839	0.944±0.035	0.995±0.008	1.322±0.009
	MPSP(10 ⁻²)	0.074±0.009	N/A	0.032 ± 0.011	0.099±0.005	0.764±0.004
MP-Humanoid-v4	MPHV	1.720±0.330	0.000	1.508±0.275	1.847±0.194	2.761 ± 0.361
	MPSP(10 ⁻²)	0.381±0.030	N/A	0.002±0.001	0.000 ± 0.000	0.694±0.032

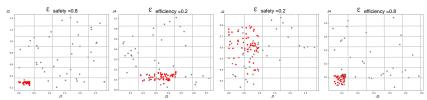


Figure 5: The Pareto policy sets ultimately obtained with different initial negotiation vectors. The red point means the multi-party Pareto policy.

Across all environments, PGMORL and LP3 exhibits inferior performance, which may be attributed to the difficulty of the predictive model in accurately guiding the policy when the number of objectives is large. CR-MOPO-S consistently outperforms CR-MOPO, indicating that enforcing safety as a hard constraint limits policy exploration.

Figure 5 compares the Pareto policy sets ultimately obtained with different initial negotiation vectors. It can be observed that by relaxing the constraints of the safety party, significant improvements can be achieved in performance objectives.

By removing the MPPN component, we obtain the multi-objective evolutionary reinforcement learning algorithm MPNDSRL. Figure 6 depicts the Pareto fronts obtained by MPPN-MORL and MPNDSRL for the two parties in the MP-Halfcheetah environment. MPPN-MORL achieves a better balance between safety and efficiency, and finds a sufficient number of policies while achieving better performance. In contrast, MPNDSRL only finds a very small number of common solutions, which proves that our method can still achieve excellent results when common solutions are scarce.

5 RELATED WORK

5.1 MORL AND CMORL

MORL tackles tasks with multiple conflicting objectives and mainly includes single-policy and multi-policy approaches. Single-policy methods scalarize multiple rewards into a single objective and apply standard RL to maximize it (Roijers et al., 2013), but they rely on expert-defined preference weights (Van Moffaert et al., 2013; Abdolmaleki et al., 2020) that may vary with real-world conditions. Multi-policy methods approximate the Pareto front by learning a set of policies under different preferences (Roijers et al., 2014; Mossalam et al., 2016; Zuluaga et al., 2016). Typical methods include PGMORL (Xu et al., 2020), which improves efficiency through predictive models and PPO updates but risks local minima; PD-MORL (Basaklar et al., 2023) obtains a unified network covering the entire preference space through single-round training; PA2D-MORL (Hu & Luo, 2024), which uses Pareto ascent directions for automatic optimization and better coverage; and

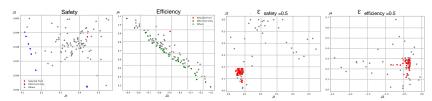


Figure 6: MPNDS algorithm without the MPPN mechanism compares with the MPPN-MORL.

PSL-MORL (Liu et al., 2025), which employs hypernetworks to generate preference-conditioned policies compatible with single-objective RL.

Despite these advances, existing MORL methods optimize multiple objectives only for a single agent and cannot model multi-party interactions or conflicts. Consequently, they fail to capture the negotiation dynamics and collective trade-offs essential in multi-stakeholder scenarios, leading to suboptimal solutions.

CMORL further incorporates safety requirements into multi-objective optimization. LP3 (Huang et al., 2022) jointly learns preferences and policies by treating task rewards and constraint costs as independent objectives. PDOA (Lin et al., 2024) supports offline adaptation under unknown preferences and safety thresholds by learning diverse policies and conservatively estimating preference weights to mitigate violation risks. CR-MOPO (Gu et al., 2025) integrates conflict-aware gradients and hard constraint corrections to ensure safety while efficiently approximating the Pareto front.

Nevertheless, CMORL still focuses on single-agent optimization, lacking mechanisms to model interactions and negotiations among multiple parties. In contrast, MPMORL explicitly captures multiagent interactions and negotiation dynamics, offering superior modeling capability and adaptability in complex multi-stakeholder environments.

5.2 MPMOP

MPMOPs aim to identify mutually optimal solutions for multiple DMs with diverse and often conflicting objectives, a critical challenge in many real-world scenarios. To address this, researchers have developed various MPMOEAs by extending existing MOEA frameworks with ranking and selection mechanisms for multi-party settings. OptMPNDS (Liu et al., 2020) ranks solutions by their worst dominance level across all DMs, while OptMPNDS2 (She et al., 2021) refines this by treating dominance levels from each DM as new objectives and applying a second non-dominated sorting for finer evaluation. A theoretical analysis (Sun et al., 2025) further revealed the inefficiency of traditional MOEAs, especially for NP-hard problems.

Despite the success of MPMOEAs in solving MPMOPs, no prior studies have integrated them into RL. The proposed MPPN-MORL addresses this gap by reformulating MORL as an MPMOP, thereby establishing the first link between these two research domains.

6 Conclusion

This paper reformulates MORL with safety constraints as a MPMORL problem and proposes an evolutionary algorithm, MPPN-MORL, based on a multi-party Pareto negotiation mechanism. In the proposed approach, efficiency and safety are treated as two independent decision-making parties. During the evolutionary process, the algorithm maintains individual Pareto fronts for each party and integrates them into a shared solution set through a NBMPNDS procedure. This design effectively alleviates the complexity caused by the proliferation of objectives in traditional MORL. Unlike CMORL, which enforces safety as a hard constraint and strictly limits exploration, MPPN-MORL dynamically adjusts the trade-off between safety and efficiency, producing high-quality compromise solutions. Experimental results demonstrate that across six mp-MuJoCo environments, MPPN-MORL consistently achieves the highest MeanHV and SP metrics, significantly outperforming state-of-the-art MORL and CMORL methods, while exhibiting superior balance and diversity in strategies when handling conflicts between safety and efficiency.

7 REPRODUCIBILITY STATEMENT

We have taken extensive efforts to ensure the reproducibility of our work. The proposed algorithms and benchmark implementations have been anonymously submitted as supplementary materials and will be publicly released upon publication. The benchmark environment setups are detailed. These resources collectively enable independent verification and reproduction of our reported results.

REFERENCES

- Abbas Abdolmaleki, Sandy Huang, Leonard Hasenclever, Michael Neunert, Francis Song, Martina Zambelli, Murilo Martins, Nicolas Heess, Raia Hadsell, and Martin Riedmiller. A distributional view on multi-objective policy optimization. In *International conference on machine learning*, pp. 11–22. PMLR, 2020.
- Toygun Basaklar, Suat Gumussoy, and Umit Y Ogras. Pd-morl: Preference-driven multi-objective reinforcement learning algorithm. The Eleventh International Conference on Learning Representations, https..., 2023.
- SenPeng Chen, Jia Wu, and XiYuan Liu. Emorl: Effective multi-objective reinforcement learning method for hyperparameter optimization. *Engineering Applications of Artificial Intelligence*, 104: 104315, 2021.
- Xi Chen, Ali Ghadirzadeh, Mårten Björkman, and Patric Jensfelt. Meta-learning for multi-objective reinforcement learning. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 977–983. IEEE, 2019.
- Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and TAMT Meyarivan. A fast and elitist multiobjective genetic algorithm: Nsga-ii. *IEEE transactions on evolutionary computation*, 6(2): 182–197, 2002.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, 110(9):2419–2468, 2021.
- Shangding Gu, Bilgehan Sel, Yuhao Ding, Lu Wang, Qingwei Lin, Alois Knoll, and Ming Jin. Safe and balanced: A framework for constrained multi-objective reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz, et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.
- Tianmeng Hu and Biao Luo. Pa2d-morl: Pareto ascent directional decomposition based multiobjective reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 12547–12555, 2024.
- Sandy Huang, Abbas Abdolmaleki, Giulia Vezzani, Philemon Brakel, Daniel J Mankowitz, Michael Neunert, Steven Bohez, Yuval Tassa, Nicolas Heess, Martin Riedmiller, et al. A constrained multi-objective reinforcement learning framework. In *Conference on Robot Learning*, pp. 883–893. PMLR, 2022.
- Panagiotis Kyriakis and Jyotirmoy Deshmukh. Pareto policy adaptation. In *International Conference on Learning Representations*, volume 2022, 2022.
- Qian Lin, Zongkai Liu, Danying Mo, and Chao Yu. An offline adaptation framework for constrained multi-objective reinforcement learning. Advances in Neural Information Processing Systems, 37: 140292–140319, 2024.
- Erlong Liu, Yu-Chang Wu, Xiaobin Huang, Chengrui Gao, Ren-Jian Wang, Ke Xue, and Chao Qian. Pareto set learning for multi-objective reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 18789–18797, 2025.

540	Wenjie Liu, Wenjian Luo, Xin Lin, Miqing Li, and Shengxiang Yang. Evolutionary approach to				
541	multiparty multiobjective optimization problems with common pareto optimal solutions. In 2020				
542	IEEE Congress on Evolutionary Computation (CEC), pp. 1–9. IEEE, 2020.				
543					
544	Hossam Mossalam, Yannis M Assael, Diederik M Roijers, and Shimon Whiteson. Multi-objective				
545	deep reinforcement learning. arXiv preprint arXiv:1610.02707, 2016.				
546					
547	Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of multi-				
548	objective sequential decision-making. <i>Journal of Artificial Intelligence Research</i> , 48:67–11				
549	2013.				
550					

- Diederik M Roijers, Shimon Whiteson, and Frans A Oliehoek. Linear support for multi-objective coordination graphs. In *AAMAS'14: PROCEEDINGS OF THE 2014 INTERNATIONAL CONFERENCE ON AUTONOMOUS AGENTS & MULTIAGENT SYSTEMS*, volume 2, pp. 1297–1304. IFAAMAS/ACM, 2014.
- Zeneng She, Wenjian Luo, Yatong Chang, Xin Lin, and Ying Tan. A new evolutionary approach to multiparty multiobjective optimization. In *International Conference on Swarm Intelligence*, pp. 58–69. Springer, 2021.
- Joar Skalse, Lewis Hammond, Charlie Griffin, and Alessandro Abate. Lexicographic multi-objective reinforcement learning. *arXiv* preprint arXiv:2212.13769, 2022.
- Rainer Storn and Kenneth Price. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization*, 11(4):341–359, 1997.
- Yuetong Sun, Peilan Xu, and Wenjian Luo. Runtime analysis of evolutionary algorithms for multiparty multiobjective optimization. *arXiv* preprint arXiv:2501.16336, 2025.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. IEEE, 2012.
- Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforcement learning: Novel design techniques. In 2013 IEEE symposium on adaptive dynamic programming and reinforcement learning (ADPRL), pp. 191–199. IEEE, 2013.
- Jie Xu, Yunsheng Tian, Pingchuan Ma, Daniela Rus, Shinjiro Sueda, and Wojciech Matusik. Prediction-guided multi-objective reinforcement learning for continuous robot control. In *International conference on machine learning*, pp. 10607–10616. PMLR, 2020.
- Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective reinforcement learning and policy adaptation. *Advances in neural information processing systems*, 32, 2019.
- Marcela Zuluaga, Andreas Krause, and Markus Püschel. e-pal: An active learning approach to the multi-objective optimization problem. *Journal of Machine Learning Research*, 17(104):1–32, 2016.

A THE USE OF LARGE LANGUAGE MODELS

We used a large language model to assist in polishing the writing and checking for grammatical issues.

B MPPN-SAFEMORL

In Algorithm 2, we describe the whole procedure of our proposed MPPN-SafeMORL algorithm.

C EXPERIMENT SETUP DETAILS

For each episode, the reward (or cost) for each objective is computed as the average of the corresponding per-step values over all time steps within that episode.

MP-Halfcheetah: In the MP-HalfCheetah environment, the safety party seeks to minimize energy consumption and maintain the stability of the robot's height, whereas the efficiency party aims to maximize forward velocity while mitigating excessive oscillations. Energy consumption is quantified as the squared norm of the action vector:

$$C_e^i = -\alpha_a \|a_{\text{cheetah}}^i\|^2, \tag{12}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and $a_{\rm cheetah}^i$ represents the action vector applied to the HalfCheetah at time step i.

Height stability is evaluated by the deviation of the robot's height from a target value:

$$C_h^i = \left| H_{\text{cheetah}}^i - H_{\text{target}}^i \right|, \tag{13}$$

where C_h^i is the height stability cost at time step i, $H_{\rm cheetah}^i$ is the actual torso height of the HalfCheetah at time step i, and $H_{\rm target}^i$ is the predefined target height.

Forward velocity is represented by the absolute value of the robot's horizontal velocity:

$$R_x^i = \left| V_x^i \right|,\tag{14}$$

where R_x^i denotes the forward velocity reward at time step i, and V_x^i is the horizontal velocity of the HalfCheetah at time step i.

Excessive oscillation is penalized using the absolute value of the robot's vertical velocity:

$$R_y^i = -\left|V_y^i\right|,\tag{15}$$

where R_y^i is the oscillation penalty at time step i, and V_y^i represents the vertical velocity of the HalfCheetah at time step i.

MP-Hopper: In the MP-Hopper environment, the safety party aims to minimize the robot's angular deviation around the z-axis and reduce energy consumption, while the efficiency party seeks to maximize forward velocity in the x-direction while minimizing vertical oscillations in the y-direction.

The angular deviation around the z-axis is quantified by the absolute value of the robot's z-axis angle:

$$C_z^i = \left| \Theta_z^i \right|, \tag{16}$$

where C_z^i denotes the angular deviation cost at time step i, and Θ_z^i is the robot's orientation angle around the z-axis at time step i.

Energy consumption is measured as the squared norm of the action vector:

$$C_e^i = -\alpha_a |a_{\text{hopper}}^i|^2, \tag{17}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and a_{hopper}^i represents the action vector applied to the Hopper at time step i.

Forward velocity is represented by the absolute value of the robot's horizontal velocity in the x-direction:

$$R_x^i = \left| V_x^i \right|,\tag{18}$$

where R_x^i denotes the forward velocity reward at time step i, and V_x^i is the horizontal velocity of the Hopper at time step i.

Vertical oscillation is penalized using the absolute value of the robot's velocity in the y-direction:

$$R_y^i = -\left|V_y^i\right|,\tag{19}$$

where R_y^i is the oscillation penalty at time step i, and V_y^i represents the vertical velocity of the Hopper at time step i.

MP-Walker: In the MP-Walker environment, the safety party aims to minimize the absolute height of the robot's head and reduce the degree of body posture deviation, while the efficiency party seeks to maximize forward velocity in the x-direction while minimizing energy consumption.

The head height cost is quantified by the absolute value of the robot's head height:

$$C_z^i = \left| Z_{\text{head}}^i \right|, \tag{20}$$

where C_z^i denotes the head height cost at time step i, and $Z_{\rm head}^i$ is the vertical height of the Walker's head at time step i.

The body posture cost is measured by the absolute deviation of the robot's posture:

$$C_n^i = \left| P_{\text{walker}}^i \right|,\tag{21}$$

where C_p^i denotes the posture deviation cost at time step i, and P_{walker}^i represents the robot's body inclination angle at time step i.

Forward velocity is represented by the absolute value of the robot's horizontal velocity in the x-direction:

$$R_x^i = \left| V_x^i \right|,\tag{22}$$

where R_x^i denotes the forward velocity reward at time step i, and V_x^i is the horizontal velocity of the Walker at time step i.

Energy consumption is quantified as the squared norm of the action vector:

$$C_e^i = -\alpha_a |a_{\text{walker}}^i|^2, \tag{23}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and a_{walker}^i represents the action vector applied to the Walker at time step i.

MP-Swimmer: In the MP-Swimmer environment, the safety party aims to minimize energy consumption and reduce the degree of body oscillation, while the efficiency party seeks to maximize forward velocity in the x-direction while minimizing vertical velocity in the y-direction.

Energy consumption is quantified as the squared norm of the action vector:

$$C_e^i = -\alpha_a |a_{\text{swimmer}}^i|^2, \tag{24}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and a_{swimmer}^i represents the action vector applied to the Swimmer at time step i.

Body oscillation is measured by the absolute value of the robot's angular velocity:

$$C_o^i = \left| \Omega_{\text{swimmer}}^i \right|, \tag{25}$$

where C_o^i denotes the body oscillation cost at time step i, and $\Omega_{\mathrm{swimmer}}^i$ is the angular velocity of the Swimmer at time step i.

Forward velocity is represented by the absolute value of the robot's horizontal velocity in the x-direction:

$$R_x^i = |V_x^i|, (26)$$

where R_x^i denotes the forward velocity reward at time step i, and V_x^i is the horizontal velocity of the Swimmer at time step i.

Vertical velocity is penalized using the absolute value of the robot's velocity in the y-direction:

$$R_y^i = -\left|V_y^i\right|,\tag{27}$$

where R_y^i is the vertical velocity penalty at time step i, and V_y^i represents the vertical velocity of the Swimmer at time step i.

MP-Pusher: In the MP-Pusher environment, the safety party aims to minimize energy consumption and reduce the velocity of the robot's end-effector, while the efficiency party seeks to minimize the distance between the actuator and the object as well as the distance between the object and the target position.

Energy consumption is quantified as the squared norm of the action vector:

$$C_e^i = -\alpha_a |a_{\text{pusher}}^i|^2, \tag{28}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and a_{pusher}^i represents the action vector applied to the Pusher at time step i.

The end-effector velocity cost is measured by the absolute value of the end-effector's velocity:

$$C_v^i = \left| V_{\text{end}}^i \right|, \tag{29}$$

where C_v^i denotes the end-effector velocity cost at time step i, and V_{end}^i is the velocity of the Pusher's end-effector at time step i.

The actuator-to-object distance is evaluated as the Euclidean distance between the actuator and the object:

$$R_{ao}^{i} = -\left|P_{\text{actuator}}^{i} - P_{\text{object}}^{i}\right|,\tag{30}$$

where R_{ao}^i denotes the actuator-to-object distance reward at time step i, $P_{\rm actuator}^i$ is the position of the actuator at time step i, and $P_{\rm object}^i$ is the position of the object at time step i.

The object-to-target distance is evaluated as the Euclidean distance between the object and the target position:

$$R_{ot}^{i} = -\left|P_{\text{object}}^{i} - P_{\text{target}}^{i}\right|,\tag{31}$$

where R_{ot}^i denotes the object-to-target distance reward at time step i, P_{target}^i is the predefined target position, and P_{object}^i is the object's position at time step i.

MP-Humanoid: In the MP-Humanoid environment, the safety party aims to minimize energy consumption and reduce contact impact, while the efficiency party seeks to maximize forward velocity in the x-direction and enhance the humanoid's health reward.

Energy consumption is quantified as the squared norm of the action vector:

$$C_e^i = -\alpha_a |a_{\text{humanoid}}^i|^2, \tag{32}$$

where C_e^i denotes the energy consumption at time step i, α_a is a scaling coefficient, and a_{humanoid}^i represents the action vector applied to the Humanoid at time step i.

Contact impact is measured by the magnitude of the external contact forces exerted on the humanoid:

$$C_c^i = \left| F_{\text{contact}}^i \right|, \tag{33}$$

where C_c^i denotes the contact impact cost at time step i, and $F_{\rm contact}^i$ represents the contact force vector applied to the Humanoid at time step i.

Forward velocity is represented by the absolute value of the humanoid's horizontal velocity in the x-direction:

$$R_x^i = \left| V_x^i \right|,\tag{34}$$

where R_x^i denotes the forward velocity reward at time step i, and V_x^i is the horizontal velocity of the Humanoid at time step i.

The health reward is quantified by the humanoid's uprightness and stability:

$$R_h^i = H_{\text{humanoid}}^i, \tag{35}$$

where R_h^i denotes the health reward at time step i, and H_{humanoid}^i is the environment-defined health indicator of the Humanoid at time step i.

```
758
759
760
761
762
763
764
765
766
767
768
            Algorithm 2 Multi-Party NSDE with \varepsilon-Dominance and Priority Selection
769
            Input: Number of iterations T, population size N, mutation factor F, crossover rate CR, initial
770
            tolerance \varepsilon_{\text{init}}, decay rate \Delta \varepsilon, objective partitions \{O_1, O_2\}
771
            Output: Final \varepsilon-dominant Pareto front \mathcal{F}
             1: Initialize empty population \mathcal{P} \leftarrow \emptyset
772
773
             2: for i = 1 to N do
774
                     Initialize random policy \pi_i with parameters \theta_i
             4:
                     Evaluate \pi_i to obtain reward vector \mathbf{r}_i
775
             5:
                     Add (\pi_i, \theta_i, \mathbf{r}_i) to \mathcal{P}
776
             6: end for
777
             7: Initialize tolerance vector \varepsilon = [\varepsilon_{\text{init}}, \varepsilon_{\text{init}}]
778
             8: Compute reference solution \mathbf{r}_{ref}
779
             9: for t = 1 to T do
780
            10:
                     \mathcal{Q} \leftarrow \emptyset
781
            11:
                     for i = 1 to N do
782
            12:
                         Generate offspring parameters \theta_{\text{trial}} using DE mutation and crossover with (F, CR)
783
                         Evaluate offspring to obtain \mathbf{r}_{\mathrm{trial}}
            13:
            14:
                         Add offspring (\pi_{\text{trial}}, \theta_{\text{trial}}, \mathbf{r}_{\text{trial}}) to Q
784
            15:
                     end for
785
                     Combine populations: C \leftarrow P \cup Q
            16:
786
                     Extract rewards \mathcal{R} = \{\mathbf{r}_j\}_{j \in \mathcal{C}}
            17:
787
            18:
                     Identify \varepsilon-dominant fronts:
788
                          1.
                                  Find joint solutions \mathbf{r}_j such that \mathbf{r}_j \varepsilon-dominates \mathbf{r}_{ref} for both DMs.
789
                                  If none, update each DM's front separately based on its own objective set O_k.
790
            19:
                     If joint \varepsilon-dominant solutions are found: update \varepsilon \leftarrow \varepsilon - \Delta \varepsilon
791
            20:
                     Build next population \mathcal{P}_{t+1}:
792
                                  Add all individuals from \varepsilon-dominant fronts to \mathcal{P}_{t+1}
793
                                  If |\mathcal{P}_{t+1}| < N, fill the remaining slots by applying crowding distance selection on
794
                                  the rest of C
                     \mathcal{P} \leftarrow \mathcal{P}_{t+1}
795
            21:
            22: end for
796
            23: Extract final front \mathcal{F} from \mathcal{P} based on \varepsilon-dominance
797
            24: return \mathcal{F}
798
799
800
```