# Preference-based Multi-Objective Reinforcement Learning

Ni Mu*, Yao Luan*, Qing-Shan Jia†

*Abstract*—Multi-objective reinforcement learning (MORL) is a structured approach for optimizing tasks with multiple objectives. However, it often relies on pre-defined reward functions, which can be hard to design for balancing conflicting goals and may lead to oversimplification. Preferences can serve as more flexible and intuitive decision-making guidance, eliminating the need for complicated reward design. This paper introduces preference-based MORL (Pb-MORL), which formalizes the integration of preferences into the MORL framework. We theoretically prove that preferences can derive policies across the entire Pareto frontier. To guide policy optimization using preferences, our method constructs a multi-objective reward model that aligns with the given preferences. We further provide theoretical proof to show that optimizing this reward model is equivalent to training the Pareto optimal policy. Extensive experiments in benchmark multi-objective tasks, a multi-energy management task, and an autonomous driving task on a multi-line highway show that our method performs competitively, surpassing the oracle method, which uses the ground truth reward function. This highlights its potential for practical applications in complex real-world systems.

*Note to Practitioners*—Decision-making problems with multiple conflicting objectives are common in real-world applications, e.g., energy management must balance system lifespan, charge-discharge cycles, and energy procurement costs; autonomous driving vehicles must balance safety, speed, and passenger comfort. While multi-objective reinforcement learning (MORL) is an effective framework for these problems, its dependence on pre-defined reward functions can limit its application in complex situations, as designing a reward function often fails to capture the full complexity of the task fully. This paper introduces preference-based MORL (Pb-MORL), which utilizes user preference data to optimize policies, thereby eliminating the complexity of reward design. Specifically, we construct a multi-objective reward model that aligns with user preferences and demonstrate that optimizing this model can derive Pareto optimal solutions. Pb-MORL is effective, easy to deploy, and is expected to be applied in complex systems, e.g., multi-energy management through preference feedback and adaptive autonomous driving policies for diverse situations.

*Index Terms*—Reinforcement learning, Multi-objective optimization, Preference-based optimization, Pareto efficiency.

## I. INTRODUCTION

Multi-objective optimization is pervasive in real-world applications [1]–[3]. For example, in an energy system, the goal

N. Mu, Y. Luan and Q. Jia are with the Center for Intelligent and Networked System (CFINS), Department of Automation, Beijing National Research Center for Information Science and Technology, Beijing Key Laboratory of Embodied Intelligence Systems, Tsinghua University, Beijing 100084, China, {mn23@mails., luany23@mails., jiaqs@} tsinghua.edu.cn. *N. Mu and Y. Luan contributed equally. †Q. Jia is the corresponding author.

is to maximize the system lifespan and minimize charge-discharge cycles while simultaneously reducing energy procurement costs [4]. Autonomous vehicles need to provide safe, fast, and comfortable rides at the same time [5]. However, representing these objectives with a single reward can be difficult and may lose important information [6]–[8]. In addition, creating a scalar reward function for each control objective is challenging and often results in oversimplification [9], [10]. Preferences, conversely, offer a more flexible and general way to model the decision-making process [11]. Humans can easily provide their preferences, pointing out which outcome they prefer, without compressing all their decision-making information into a single reward function [12]. Therefore, it is of great practical interest to study multi-objective reinforcement learning based on preference.

However, integrating multi-objective reinforcement learning (MORL) with preference-based learning presents several challenges. First, while users can express preferences between pairs of behaviors when focusing on a single objective, establishing a complete ordering among all behaviors is often difficult. This lack of a complete preference makes it hard for algorithms to assess the relative importance of different objectives. Additionally, there are often inherent conflicts between objectives, which complicate policy optimization. Furthermore, obtaining preferences for all objectives may require pairwise comparisons, which can be computationally inefficient as the number of objectives increases, leading to increased complexity in the querying process. Given these complexities, a significant gap lies in the previous work: To the best knowledge of the authors, we did not find any method addressing the above challenges of combining MORL and preference-based optimization, highlighting the need for a novel approach to multi-objective, preference-driven decision-making problems.

### A. Related Work

**Single-objective reinforcement learning with preference.** Reinforcement learning (RL) [13] has gained significant attention in recent years due to its remarkable success in solving challenging problems [14]–[19]. Traditional RL algorithms often rely on a pre-defined reward function, which serves as the guidance for policy optimization. However, designing such a reward function can be complex and even impractical [9]. However, there are two typical categories of RL problems where we face difficulty obtaining the optimal policy. The first category of RL problems has a pre-defined reward function, such as cost savings [20], reducing carbon emissions [21],
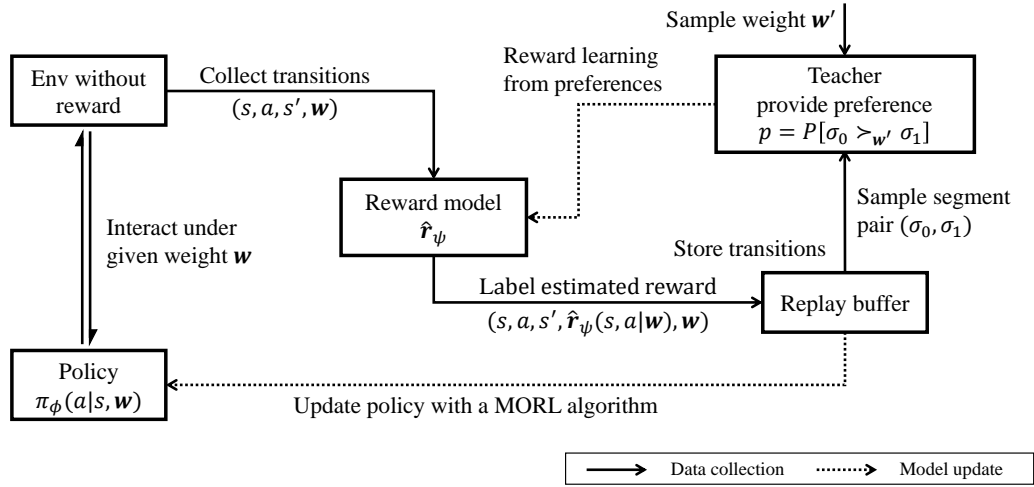
Fig. 1. A demonstration of the proposed Pb-MORL framework. An explicit multi-objective reward model $\hat{r}_\psi$ is learned using preference data. Then, the multi-objective policy $\pi_\phi(a|s, \boldsymbol{w})$ can be updated through any MORL algorithm based on the reward model. In this figure, $\boldsymbol{w}$ denotes the weight vector, $(\sigma_0, \sigma_1)$ denotes the segment pair, $p$ denotes the preference provided by the teacher. A detailed introduction to the settings and notations will be provided in the following sections.

or sparse rewards like a "+1" bonus for reaching the goal in a maze [22]. The challenge in this category is identifying the optimal policy, as task dynamics are often complex and stochastic, while sparse reward signals complicate policy learning from the reward function [23]. The second category consists of RL problems where defining the reward function is challenging, such as in robotics systems [9] and large language models [24]. In these cases, while we want the system behaviors to align with human expectations, formalizing the objective function is often difficult [9]. Preference-based reinforcement learning (PbRL) provides a solution by utilizing user feedback to guide agent behavior, making it suitable for both categories of RL problems. This approach offers preferences that may be more accessible and more naturally aligned with policy optimization than traditional reward signals. Early works of PbRL, such as [11] and [12], have shown the ability of agents to learn from simple comparisons between pairs of trajectory segments, thereby eliminating the necessity for complex reward engineering. With the development of deep learning, techniques like pre-training [6], [25], [26] and data augmentation [27] are employed to improve learning efficiency. Meta-learning approaches [28] also enable agents to adapt to new tasks based on past experiences quickly. Moreover, PbRL has been successfully applied to fine-tune large-scale language models like GPT-3 for challenging tasks, as highlighted by [29]. While PbRL omits reward engineering through leveraging user feedback, it primarily deals with single-objective optimization instead of multi-objective preference modeling.

**Multi-objective reinforcement learning with explicit reward functions.** Multi-objective reinforcement learning (MORL) is a pivotal subfield of reinforcement learning [30]–[32], focusing on decision-making problems under multiple objectives. Envelope Multi-objective Q-learning [33] extends the traditional Q-learning algorithm to the multi-objective domain and proves the convergence of its Q-learning algorithm

in tabular settings. Expected Utility Policy Gradient (EUPG) [34] and Prediction-Guided MORL (PGMORL) [35] further integrate deep learning into MORL. EUPG incorporates policy gradients to balance current and prospective returns, while PGMORL applies an evolutionary strategy to enhance the Pareto frontier. Additionally, Pareto Conditioned Networks [36] and Generalized Policy Improvement Linear Support [37] employ neural networks conditioned on target returns to predict optimal actions within deterministic settings. Despite their advancements, current MORL methods rely on predefined multi-objective reward functions, posing challenges for their application in real-world control scenarios. Extending preferences from single-objective reinforcement learning to multi-objective contexts is feasible, which is the main contribution of this paper.

### B. Main Contributions

In this paper, we introduce Preference-based Multi-Objective Reinforcement Learning (Pb-MORL), which integrates preference modeling into Multi-Objective Reinforcement Learning (MORL), as illustrated in Fig. 1. Specifically, we first establish theorems that demonstrate a teacher providing preferences can guide the learning of optimal multi-objective policies (Theorem 1, 2). Furthermore, we propose a method to construct an explicit multi-objective reward model that aligns with the teacher's preferences. Our theoretical proof (Theorem 4) shows that, in a multi-objective context, if the reward function perfectly matches the teacher's preferences, optimizing this reward is equivalent to learning the optimal policy. To implement Pb-MORL, we combine the Envelope Q Learning (EQL) method [33] with our proposed reward model. This implementation is simple yet effective, for EQL guarantees the convergence of policy optimization in multi-objective tasks. To demonstrate the effectiveness of our method, we conduct experiments in benchmark multi-objective reinforcement learning tasks. The results show that

TABLE I
NOTATION TABLE

| Symbol | Definition | Description |
|---|---|---|
| *Multi-Objective RL Elements* | | |
| $m$ | Number of objectives | Dimension of the reward vector. |
| $\boldsymbol{r}(s,a)$ | True multi-objective reward vector | $\boldsymbol{r}(s,a) \in \mathbb{R}^m$: Ground truth reward signal provided by the environment for each objective. |
| $\boldsymbol{w}$ | Weight vector | $\boldsymbol{w} \in \mathcal{W} = \{\boldsymbol{w} \in \mathbb{R}^m \mid w_i \geq 0, \sum_i w_i = 1\}$: Vector encoding the relative importance (preference) assigned to each objective. |
| $\mathcal{W}$ | Weight space | Set of all valid weight vectors. |
| $D_w$ | Prior weight distribution | Distribution over the weight space $\mathcal{W}$ from which weights are sampled during training or evaluation. In this study, we assume this distribution to be uniform over $\mathcal{W}$. |
| $\Pi$ | Policy space | Set of all possible policies. |
| $\Pi^*$ | Pareto optimal policy set | Set of policies that are not dominated by any other policy in $\Pi$ with respect to all objectives. |
| *Preference Elements* | | |
| $\sigma$ | Trajectory segment | Finite sequence of state-action pairs: $\sigma = \{s_k, a_k, \ldots, s_{k+H-1}, a_{k+H-1}\}$ of length $H$. |
| $H$ | Segment length | Number of steps in a trajectory segment $\sigma$. |
| $p$ | Preference label | $p \in \{0, 0.5, 1\}$: Human teacher's preference judgment for a pair of segments $(\sigma_0, \sigma_1)$ under weight $\boldsymbol{w}$. $p = 0$: $\sigma_0$ preferred, $p = 1$: $\sigma_1$ preferred, $p = 0.5$: no preference/indifferent. |
| $\sigma_0 \succ_{\boldsymbol{w}} \sigma_1$ | Preference relation | Segment $\sigma_0$ is strictly preferred over segment $\sigma_1$ under weight $\boldsymbol{w}$. |
| *Learned Components* | | |
| $\hat{\boldsymbol{r}}_\psi(s,a)$ | Learned multi-objective reward model | $\hat{\boldsymbol{r}}_\psi(s,a) \in \mathbb{R}^m$: Model parameterized by $\psi$, trained using preference data to approximate the underlying objectives. Used as the reward signal for MORL policy optimization. |
| $\pi_\phi(a\|s, \boldsymbol{w})$ | Parameterized policy | Stochastic policy parameterized by $\phi$, conditioned on the current state $s$ and the weight vector $\boldsymbol{w}$ (indicating the desired objective trade-off). Outputs a distribution over actions. |
| $Q_\theta(s, a, \boldsymbol{w})$ | Multi-objective Q-function | $Q_\theta(s, a, \boldsymbol{w}) \in \mathbb{R}^m$: State-action value function parameterized by $\theta$. Estimates the vector of expected discounted future rewards for each objective, starting from state $s$, taking action $a$, and following policy $\pi_\phi(\cdot\|\cdot, \boldsymbol{w})$ thereafter. |
| $J$ | Optimization objective | Scalarized expected return: $J = \mathbb{E}_{\boldsymbol{w} \sim D_w, \tau \sim (\mathcal{P}, \pi)}[\boldsymbol{w}^T \sum_t \gamma^t \hat{\boldsymbol{r}}_\psi(s_t, a_t)]$ (Eq. 7). Maximized during policy learning using the learned reward model. |

our approach achieves performance levels comparable to the oracle method, which uses the ground truth reward function to learn the optimal policy. To validate our method's applicability in real-world scenarios, we evaluate our method on both a multi-energy management task and an autonomous driving task on a multi-line highway. In both settings, the Pb-MORL algorithm outperforms the oracle, showing its potential for practical implementation in complex, real-world environments. Through this work, we aim to broaden the applications of MORL in real-world settings, by employing preferences as a more accessible and intuitive optimization guidance.

The main contributions of this paper are as follows:

- We establish theorems for preference-based optimization in multi-objective settings, demonstrating that a preference-based teacher can guide the learning of optimal multi-objective policies (Theorem 1, 2, 4).
- We introduce Pb-MORL, which develops an explicit multi-objective reward model that aligns with preference data through the construction of the Bradley-Terry model and the optimization of the cross-entropy loss function. In addition, we combine the EQL algorithm with the reward model to achieve a simple yet effective implementation of Pb-MORL.
- We conduct experiments in multi-objective benchmark tasks, a multi-energy management task, and an autonomous driving task on a multi-line highway, showing that Pb-MORL performs comparably to the oracle method

using ground truth reward functions. It demonstrates Pb-MORL's potential for real-world applications.

The remaining sections are organized as follows. In Section II, we introduce preliminaries and the problem formulation. In Section III, we present the theoretical guarantees of Pb-MORL and propose the specific algorithm for explicit reward modeling and policy optimization. In Section IV, we describe the experimental setting and discuss the experimental results. Finally, we conclude the paper in Section V.

## II. PROBLEM FORMULATION

In this section, we first introduce the multi-objective MDP and Q-learning in multi-objective settings, then formulate the Pb-MORL framework.

### A. MDP and Q-learning in Multi-Objective Settings

For single-objective settings, an MDP with discrete time, infinite-stage discounted reward, and finite or countable state and action spaces could be characterized as a tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, r, \gamma \rangle$. Here, $\mathcal{S}$ is the state space, $\mathcal{A}$ is the action space, and $P(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the one-step state transition probability of transiting from $s$ to $s'$ by taking action $a$. Besides, $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defines the immediate reward of taking action $a$ under state $s$, and $\mathbb{R}$ denotes the set of real numbers. Finally, $\gamma \in (0, 1)$ is the discount factor for balancing immediate and long-term rewards.

For multi-objective settings, the MDP framework is extended to include multiple reward functions. The reward function can be represented as a vector $\boldsymbol{r}(s,a) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}^m$, where $m$ is the number of objectives. In the case of linear reward combination, the overall reward is defined by a linear combination of these objectives, $r_{\boldsymbol{w}}(s,a) = \boldsymbol{w}^T \boldsymbol{r}(s,a)$, where $\boldsymbol{w} \in \mathcal{W}$ is the weight vector, and the weight space $\mathcal{W} = \{\boldsymbol{w} | \boldsymbol{w} \in \mathbb{R}^m, w_i \geq 0, \sum w_i = 1\}$. The goal in the multi-objective MDP is to find a policy $\pi(a|s, \boldsymbol{w})$ : $\mathcal{S} \times \mathcal{W} \times \mathcal{A} \to [0,1]$ that maximizes the inner product of the multi-dimensional discounted return and the weight vector $\boldsymbol{w}$, that is,

$$\max J = \mathbb{E}_{\substack{\boldsymbol{w} \sim D_{\boldsymbol{w}} \\ \tau \sim (P, \pi(\cdot|\cdot, \boldsymbol{w}))}} \boldsymbol{w}^T \sum_{\tau} \gamma^t \boldsymbol{r}(s_t, a_t), \qquad (1)$$

where $\tau$ denotes the trajectory, and under $D_{\boldsymbol{w}}$ is a prior weight distribution. Denote the policy space as $\Pi$.

Then, the Q-learning algorithm can be adapted to the multi-objective setting. The standard Q-Learning [13], [38] for single-objective RL is based on the Bellman optimality operator $\mathcal{B}$:

$$(\mathcal{B}Q)(s,a) := r(s,a) + \sup_{a'} \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} Q(s', a'). \qquad (2)$$

Following the previous work [33], we extend this to the MORL setting, by considering multi-objective Q-value functions $\boldsymbol{Q}(s,a,\boldsymbol{w}) : \mathcal{S} \times \mathcal{A} \times \mathcal{W} \to \mathbb{R}^m$, which estimates expected total vector rewards under state $s$, action $a$ and $m$-dimensional weight $\boldsymbol{w}$. It is important to note that the parameter $\boldsymbol{w}$ in the Q function represents that the Q value is under the policy $\pi(\cdot|\cdot, \boldsymbol{w})$, because the policies conditioning on different weight vectors $\boldsymbol{w}$ results in different behaviors, and the corresponding Q functions can vary.

We define the distance between two multi-objective Q functions $\boldsymbol{Q}_1, \boldsymbol{Q}_2$ as follows:

$$d(\boldsymbol{Q}_1, \boldsymbol{Q}_2) := \sup_{s,a,\boldsymbol{w}} \left| \boldsymbol{w}^T \left( \boldsymbol{Q}_1(s,a,\boldsymbol{w}) - \boldsymbol{Q}_2(s,a,\boldsymbol{w}) \right) \right|. \quad (3)$$

The metric $d$ forms a complete pseudo-metric space, as the identity of indiscernibles [39] does not hold.

With a little abuse of notation, we use the same $\mathcal{B}_\pi$ and $\mathcal{B}$ as in the single-objective RL to represent the Bellman operator in the multi-objective setting. Specifically, given a policy $\pi$ and sampled trajectories $\tau$, the multi-objective Bellman operator for policy evaluation $\mathcal{B}_\pi$ is defined as:

$$(\mathcal{B}_\pi \boldsymbol{Q})(s,a,\boldsymbol{w}) := \boldsymbol{r}(s,a) + \gamma \mathbb{E}_{\tau \sim (P, \pi)} \boldsymbol{Q}(s', a', \boldsymbol{w}). \quad (4)$$

To construct the multi-objective Bellman optimality operator, an optimality filter $\mathcal{H}$ for the multi-objective Q function is first defined as:

$$(\mathcal{H}\boldsymbol{Q})(s,\boldsymbol{w}) := \arg_{\boldsymbol{Q}} \sup_{a \in \mathcal{A}, \boldsymbol{w}' \in \mathcal{W}} \boldsymbol{w}^T \boldsymbol{Q}(s,a,\boldsymbol{w}'), \qquad (5)$$

where the $\arg \boldsymbol{Q}$ takes the multi-objective value corresponding to the supremum (i.e., $\boldsymbol{Q}(s,a,\boldsymbol{w}')$ ) such that $(a, \boldsymbol{w}') \in \arg \sup_{a \in \mathcal{A}, \boldsymbol{w}' \in \mathcal{W}} \boldsymbol{w}^T \boldsymbol{Q}(s,a,\boldsymbol{w}'))$. Then, the multi-objective Bellman optimality operator $\mathcal{B}$ is defined as:

$$(\mathcal{B}\boldsymbol{Q})(s,a,\boldsymbol{w}) := \boldsymbol{r}(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} (\mathcal{H}\boldsymbol{Q})(s', \boldsymbol{w}) \quad (6)$$

---

**Algorithm 1** Using the teacher to derive convex Pareto frontier, based on traversing the weight space

1: Initialize the solution set $\Pi^* = \emptyset$
2: **for** each $\boldsymbol{w} \in \mathcal{W}^{[N_w]}$ **do**
3:     **for** each $\pi_i \in \Pi$ **do**
4:         **if** $\not\exists \, \pi' \in \Pi, \pi' \neq \pi_i$ s.t. $\boldsymbol{w}^T \sum_{(s,a) \sim \sigma_i} \gamma^t \boldsymbol{r}(s_t, a_t) < \boldsymbol{w}^T \sum_{(s,a) \sim \sigma'} \gamma^t \boldsymbol{r}(s_t, a_t)$, where $\sigma_i, \sigma'$ are segments generated by $\pi_i, \pi'$, **then**
5:             $\Pi^* \leftarrow \Pi^* \cup \{\pi_i\}$
6:         **end if**
7:     **end for**
8: **end for**
9: **return** $\Pi^*$

---

Intuitively, the optimality Bellman operator $\mathcal{B}$ solves the minimum convex envelope of the current $\boldsymbol{Q}$ frontier. Previous works of MORL [33] have provided proof of the convergence of the above multi-objective Q-learning algorithm, by proving the Bellman operator $\mathcal{B}_\pi$ and $\mathcal{B}$ are both contrastive mappings under the metric $d$ defined in Eq. (3).

### B. Pb-MORL Formulation

For single-objective settings, by following the previous work [11], [12], we can define the preference in the form of tuple $(\sigma_0, \sigma_1, p)$, where segment $\sigma_0, \sigma_1$ are sequences of states and actions $\{s_k, a_k, ..., s_{k+H-1}, a_{k+H-1}\}$ with length $H$ and arbitrary starting time $k$, and $p \in \{0, 0.5, 1\}$ encodes the preference relations:

- $\sigma_0$ strictly preferred to $\sigma_1$ when $p = 0$.
- $\sigma_1$ strictly preferred to $\sigma_0$ when $p = 1$.
- Indeterminate preference (equivalence or ambiguous judgment) when $p = 0.5$.

This scheme accounts for human rating uncertainty while maintaining annotation efficiency. When $\sigma_0 = \sigma_1$ or trajectories are equally preferable, $p = 0.5$ explicitly captures the uncertainty.

For multi-objective settings, we redefine the preference as a tuple $(\sigma_0, \sigma_1, \boldsymbol{w}, p)$, where $\boldsymbol{w} \in \mathcal{W}$ is a weight vector. The preference $p \in \{0, 0.5, 1\}$ is a scalar which encodes preference relations under $\boldsymbol{w}$, defined similarly as in the single-objective settings. In fact, given any weight vector, we introduce a complete ordering of the trajectory segments. However, we employ a pairwise comparison method due to practical constraints and use partial ordering notation ($\succ$) in the following paper. Specifically, let $\sigma_0 \succ_{\boldsymbol{w}} \sigma_1$ means that trajectory segment $\sigma_0$ is preferred over $\sigma_1$ under the weight vector $\boldsymbol{w}$. Then, the preference $p$ can be written in the form of $p = \mathbb{I}(\sigma_0 \succ_{\boldsymbol{w}} \sigma_1)$, where $\mathbb{I}(\cdot)$ is an indicator function that returns 1 if the condition is true, and 0 otherwise. As mentioned earlier, the weight $\boldsymbol{w}$ represents the importance assigned to each objective within the multi-objective framework. The weight $\boldsymbol{w}$ is crucial in defining the multi-objective preference, as preferences can vary for the same trajectory pair depending on the weights.

To align the problem formulation with the RL framework, we define an explicit multi-objective reward model

$\hat{r} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^m$, where each dimension corresponds to a distinct objective. This reward model can be trained using preference data, serving as a bridge between qualitative preference and quantitative rewards. Based on this model, we propose the objective of Pb-MORL as finding a policy $\pi(a|s, \boldsymbol{w})$ conditioned on the weight vector $\boldsymbol{w}$. Specifically, the goal is to maximize the inner product between the conditioned weight and the discounted return of the reward model $\hat{r}$, under a prior weight distribution $D_{\boldsymbol{w}}$, as Eq. (1) shows:

$$\max J = \mathbb{E}_{\substack{\boldsymbol{w} \sim D_{\boldsymbol{w}} \\ \tau \sim (P, \pi(\cdot|\cdot, \boldsymbol{w}))}} \boldsymbol{w}^T \sum_\tau \gamma^t \hat{r}(s_t, a_t). \qquad (7)$$

In the form of Q function, it can also be written as:

$$\max J = \mathbb{E}_{\substack{\boldsymbol{w} \sim D_{\boldsymbol{w}} \\ (s_0, a_0) \sim (P, \pi(\cdot|\cdot, \boldsymbol{w}))}} \boldsymbol{w}^T \boldsymbol{Q}_\pi(s_0, a_0, \boldsymbol{w}|\hat{r}). \qquad (8)$$

## III. A Pb-MORL algorithm with Explicit Reward Modeling

### A. Theoretical Analysis

In this subsection, we present the theoretical foundations of the Pb-MORL framework. We demonstrate how our approach ensures convergence to Pareto-optimal policies. To ease the proof, we discretize the weight space $\mathcal{W}$ to a finite space $\mathcal{W}^{[N_w]}$ with size $N_w$, and assume that when $N_w$ is large enough, $\mathcal{W}^{[N_w]}$ could fully represent $\mathcal{W}$, then induce the same set of optimal policies. We formalize this in Assumption 4.

First, we assume the presence of preferences over pairs of trajectory segments with arbitrary finite length under an arbitrary given weight. To formalize this, we introduce the following assumption.

**Assumption 1.** *The preference $p \in \{0, 0.5, 1\}$ over a pair of trajectory segments $(\sigma_0, \sigma_1)$ exists, with arbitrary finite segment length $H$, under an arbitrary given weight $\boldsymbol{w} \in W$. These preferences satisfy **symmetry**, **consistency**, and **transitivity**, which are defined as follows.*

**Definition 1** (Symmetry). *Symmetry means that if trajectory segment $\sigma_0$ is preferred over $\sigma_1$ under a weight vector $\boldsymbol{w}$, then the opposite must also be true: $\sigma_1$ is less preferred than $\sigma_0$ under the same weight. Formally, this is written as:*

$$\sigma_0 \succ_{\boldsymbol{w}} \sigma_1 \implies \sigma_1 \prec_{\boldsymbol{w}} \sigma_0. \qquad (9)$$

*This ensures that preferences are reversible under the same weight vector.*

**Definition 2** (Consistency). *Consistency means that if $\sigma_0 \succ_{\boldsymbol{w}} \sigma_1$ holds for a given $\boldsymbol{w}$, this preference remains unchanged over time. Formally, this is expressed as:*

$$\sigma_0^{t_0} \succ_{\boldsymbol{w}} \sigma_1^{t_0} \implies \forall t > 0, \ \sigma_0^t \succ_{\boldsymbol{w}} \sigma_1^t, \qquad (10)$$

*where $\sigma^t$ denotes a trajectory segment starting from time $t$, i.e. $\sigma^t = \{s_t, a_t, \cdots, s_{t+H-1}, a_{t+H-1}\}$. Here, $\sigma^{t_0}$ and $\sigma^t$ are segments with the same state action sequence $(s, a, s', \cdots)$, but starting from the different time.*

**Definition 3** (Transitivity). *Transitivity means that if the teacher prefers $\sigma_0$ over $\sigma_1$ and $\sigma_1$ over $\sigma_2$ under the same*

weight $\boldsymbol{w}$, then the teacher must also prefer $\sigma_0$ over $\sigma_2$ under weight $\boldsymbol{w}$. Formally, this is expressed as:

$$(\sigma_0 \succ_{\boldsymbol{w}} \sigma_1) \wedge (\sigma_1 \succ_{\boldsymbol{w}} \sigma_2) \implies \sigma_0 \succ_{\boldsymbol{w}} \sigma_2. \qquad (11)$$

*This property ensures logical coherence of preferences across multiple trajectory segments. Thus, the teacher's feedback does not contradict itself when extended to multiple comparisons.*

The symmetry, consistency, and transitivity requirements in Assumption 1 align with standard preference modeling in single-objective RL [9]. Then, we assume the presence of a perfect teacher, which can provide the preference over an arbitrary pair of trajectory segments with arbitrary finite length under an arbitrarily given weight.

**Assumption 2.** *We assume the existence of a teacher who can provide the preference feedback for two arbitrary trajectory segments $(\sigma_0, \sigma_1)$, based on an arbitrary weight vector $\boldsymbol{w}$.*

In Assumption 1 and 2, we assume that the teacher can provide preferences $p \in \{0, 0.5, 1\}$ over arbitrary pairs of segments $(\sigma_0, \sigma_1)$ under a given weight $\boldsymbol{w}$, and that these preferences satisfy symmetry, consistency, and transitivity. The assumption of preference availability under given weights is based on existing single-objective preference learning works [9], [27]. This indicates that the teacher's preferences are based on stable and consistent feedback related to the task objectives. Based on Assumption 1, it is reasonable to assume that the task has an underlying true reward, which is aligned with the teacher's preferences. We formalize it in Assumption 3. This assumption helps to establish a connection between the teacher's preferences and policy optimization.

**Assumption 3.** *There exists a true reward function $\boldsymbol{r}$ for a certain multi-objective task, if there exists a teacher that can express preferences for this task. Furthermore, the value of the true weighted reward $\boldsymbol{w}^T \boldsymbol{r}$ is bounded by a constant $r_{max}$. Formally, this is written as:*

$$\max_{\boldsymbol{w}, s, a} |\boldsymbol{w}^T \boldsymbol{r}(s, a)| \leq r_{max}. \qquad (12)$$

*The above equation indicates that regardless of the chosen weight vector $\boldsymbol{w}$, the absolute value of the weighted reward will not exceed this predefined upper limit.*

Assumption 3 is a common practice in existing works [40]–[42], as most real-world problems involve bounded rewards. By doing this, Assumption 3 prevents issues such as divergence in the reward function, thereby enabling Theorem 1, as discussed in the following paper.

**Assumption 4.** *The optimal policy $\pi^*(a|s, \boldsymbol{w}_0)$ under weight $\boldsymbol{w}_0$ is also the optimal policy under weight $\boldsymbol{w} \in \{\boldsymbol{w} | \|\boldsymbol{w} - \boldsymbol{w}_0\|_\infty \leq \epsilon\}, \exists \epsilon > 0, \forall s \in \mathcal{S}, a \in \mathcal{A}, \boldsymbol{w}_0 \in \mathcal{W}.$*

Assumption 4 is based on the assumption that the value function is continuous with respect to the weight vector $\boldsymbol{w}$, which is reasonable and commonplace in industrial applications. With Assumption 4, we could discretize the weight space $\mathcal{W}$ into a finite space $\mathcal{W}^{[N_w]}$ of size $N_w = \frac{|\mathcal{W}|}{\epsilon^m} \leq \epsilon^{-m}$, i.e. divide the weight space $\mathcal{W}$ to super cubes with side length $\epsilon$. The optimal policies for weights within each super cube are

identical. Therefore, $\mathcal{W}^{[N_w]}$ could fully represent $\mathcal{W}$, as they induce the same set of optimal policies.

Under Assumption 1, 2, 3 and 4, in the following theorems, we illustrate that the entire Pareto frontier could be obtained by a simple algorithm (Algorithm 1) using preferences given different weights. Specifically, we first prove that any optimal policy in an arbitrary given weight is in the Pareto frontier in Theorem 1. Then in Theorem 2, we prove that the optimal policies in all weights could form any convex Pareto frontier. Further, for non-convex Pareto frontiers, we prove the frontier could be obtained using preferences collected in designed weights in Theorem 3.

**Theorem 1.** *Each policy in the policy set $\pi^*(a|s, \boldsymbol{w}) \in \Pi^*$ derived from Algorithm 1 is in the Pareto frontier when the segment length $H \to \infty$.*

*Proof.* We prove this theorem by contradiction. Suppose $\pi^*(a|s, \boldsymbol{w})$ is not in the Pareto frontier. Then there must exist a policy $\pi^\circ(a|s, \boldsymbol{w}) \neq \pi^*(a|s, \boldsymbol{w})$ which dominates $\pi^*(a|s, \boldsymbol{w})$. And then there must exist a weight $\boldsymbol{w}_0$ and a pair of trajectories $\tau^\circ$ and $\tau^*$ which are generated from $\pi^\circ(a|s, \boldsymbol{w})$ and $\pi^*(a|s, \boldsymbol{w})$ respectively, and $\tau^\circ \succ_{\boldsymbol{w}_0} \tau^*$. We extract segments of length $H$ from $\tau^\circ$ and $\tau^*$, denoted as $\sigma^\circ$ and $\sigma^*$ respectively. Under Assumption 1, the teacher can always output the true preference between two segments.

Let $s_t^\square$ and $a_t^\square$ denote the state and action at time $t$ in the trajectory $\tau^\square$, where $\square$ is an arbitrary symbol. With discount factor $\gamma$, the difference between the discounted total return $\sum_{t=0}^\infty \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t, a_t)$ and the truncated discounted total return $\sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t, a_t)$ is bounded, i.e. $|\sum_{t=H}^\infty \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t, a_t)| \leq \frac{\gamma^H}{1-\gamma} r_{\max}$. Let $\mathcal{R}_{\underline{t}}^{\bar{t}}(\sigma^\circ) = \sum_{t=\underline{t}}^{\bar{t}} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^\circ, a_t^\circ)$, $\mathcal{R}_{\underline{t}}^{\bar{t}}(\sigma^*) = \sum_{t=\underline{t}}^{\bar{t}} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^*, a_t^*)$, we have

$$\sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^\circ, a_t^\circ) - \sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^*, a_t^*) \geq 2 \frac{\gamma^H}{1-\gamma} r_{\max}$$

$$\Leftrightarrow \mathcal{R}_0^{H-1}(\sigma^\circ) - \mathcal{R}_0^{H-1}(\sigma^*) \geq 2 \frac{\gamma^H}{1-\gamma} r_{\max}$$

$$\geq |\mathcal{R}_H^\infty(\sigma^\circ)| + |\mathcal{R}_H^\infty(\sigma^*)| \geq \mathcal{R}_H^\infty(\sigma^*) - \mathcal{R}_H^\infty(\sigma^\circ)$$

$$\Rightarrow \mathcal{R}_0^\infty(\sigma^\circ) - \mathcal{R}_0^\infty(\sigma^*) \geq 0$$

$$\Leftrightarrow \sum_{t=0}^\infty \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^\circ, a_t^\circ) - \sum_{t=0}^\infty \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^*, a_t^*) \geq 0.$$

Therefore, a sufficient condition that the preference between two trajectories is consistent with the preference between the two segments is that $|\sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^\circ, a_t^\circ) - \sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^*, a_t^*)| \geq 2 \frac{\gamma^H}{1-\gamma} r_{\max}$. This condition can always be satisfied when $H \to \infty$, which means the true preference between two trajectories ($\tau^\circ \succ_{\boldsymbol{w}_0} \tau^*$) can always be obtained from the teacher. That contradicts Algorithm 1 which only terminates when $\nexists \pi^\circ$ s.t. $\pi^\circ \succ_{\boldsymbol{w}_0} \pi^*$ and completes the proof. $\square$

In practice, we typically select pairs of segments with distinct behaviors for human comparison, facilitating humans to provide preferences. Therefore, it is reasonable to assume

---

**Algorithm 2** Using the teacher to obtain non-convex Pareto frontier, based on insertion sort

1: **for** each policy $\pi_i \in \Pi$ **do**
2:     **for** each policy $\pi_j \in \Pi$ **do**
3:         **if** for each $\boldsymbol{w}_k \in W_I$, $\boldsymbol{w}_k^T \sum_{(s,a) \sim \sigma_i} \gamma^t \boldsymbol{r}(s_t, a_t) > \boldsymbol{w}_k^T \sum_{(s,a) \sim \sigma_j} \gamma^t \boldsymbol{r}(s_t, a_t)$, where $\sigma_i, \sigma_j$ are segments generated by $\pi_i, \pi_j$, **then**
4:             Assign $\pi_i > \pi_j$
5:         **end if**
6:     **end for**
7: **end for**
8: Use insertion sort, obtain one or multiple biggest policies

---

that there exists a minimum difference $\delta$ in discounted returns between any two segments, that is, $\exists \delta \geq 0$ such that $|\sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}_1(s_t, a_t) - \sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}_2(s_t, a_t)| \geq \delta > 0$. Under this assumption, we derive the following Corollary 1.

**Corollary 1.** *If all segment pairs are distinct enough, i.e. $\exists \delta \geq 0$ s.t. $|\sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^1, a_t^1) - \sum_{t=0}^{H-1} \gamma^t \boldsymbol{w}_0^T \boldsymbol{r}(s_t^2, a_t^2)| \geq \delta > 0 \ \forall \sigma_1, \sigma_2$, then each policy in the policy set $\pi^*(a|s, \boldsymbol{w}) \in \Pi^*$ derived from Algorithm 1 is in the Pareto frontier when the segment length $H \geq \log_\gamma \frac{\delta(1-\gamma)}{2r_{max}}$.*

**Theorem 2.** *Algorithm 1 obtains the entire convex Pareto frontier, i.e., $\Pi^*$ is the entire convex Pareto frontier.*

*Proof.* Since the Pareto frontier is convex, for each policy $\pi^*$ on the Pareto frontier, there must exist a weight $\boldsymbol{w}$ s.t. $\boldsymbol{w}^T \bar{\boldsymbol{R}}^* \geq \boldsymbol{w}^T \bar{\boldsymbol{R}}'$, where $\bar{\boldsymbol{R}}^* = \mathbb{E}_\pi \sum_{t=0}^\infty \gamma^t \boldsymbol{r}(s_t, a_t)$ is the expected total discounted return derived by $\pi^*$, and $\bar{\boldsymbol{R}}'$ is that derived by any other policies. Using Theorem 1, the optimal policy under weight $\boldsymbol{w}$ could be obtained by Algorithm 1. Therefore, by traversing $\boldsymbol{w}$, we can traverse each policy on the Pareto frontier. $\square$

**Theorem 3.** *An arbitrary Pareto frontier could be completely obtained with preferences under every weight from an identity matrix weight set $W_I = \{\boldsymbol{w}_i \mid [\boldsymbol{w}_i, \cdots, \boldsymbol{w}_m] = I, i = 1, \cdots, m\}$.*

*Proof.* We prove it by providing a constructive Algorithm 2.

If there is only one policy in the policy space $\Pi$, then it is the Pareto frontier.

If we add a policy $\pi'$ into the current policy space $\Pi$, then $\pi'$ will be compared to all $\pi \in \Pi$, specifically, compared to the current Pareto frontier $\pi \in \Pi^*$ and the non-Pareto frontier $\pi \in \Pi \setminus \Pi^*$.

If $\pi'$ is not in the Pareto frontier, then $\exists \pi^* \in \Pi^*$ s.t. $\boldsymbol{w}^T \boldsymbol{R}(\sigma') < \boldsymbol{w}^T \boldsymbol{R}(\sigma^*)$ for all $\boldsymbol{w} \in W_I$, where $\boldsymbol{R}(\sigma) = \sum_{t=0,(s_t,a_t) \sim \sigma}^H \gamma^t \boldsymbol{r}(s_t, a_t)$. Thus, through Algorithm 2, $\pi'$ won't be included in the new Pareto frontier.

If $\pi'$ is in the Pareto frontier, then $\nexists \pi^* \in \Pi^*$ s.t. $\boldsymbol{w}^T \boldsymbol{R}(\sigma') < \boldsymbol{w}^T \boldsymbol{R}(\sigma^*)$ for all $\boldsymbol{w} \in W_I$. Thus, through Algorithm 2, $\pi'$ will be included in the new Pareto frontier. That completes the proof. $\square$

While linear weighting approaches ($\boldsymbol{w}^T \boldsymbol{R}$) discover only the convex Pareto frontier as in Algorithm 2, Theorem 3

operates differently. By evaluating policies under unit vector weights $W_I$ via pairwise preference comparisons, we directly assess the Pareto dominance relationship. This allows identification of non-convex Pareto-optimal policies. We provide another proof for Theorem 3 in Appendix A.

The theoretical analysis above has demonstrated that the preference-based multi-objective reinforcement learning framework can converge to the Pareto optimal set under specific conditions, providing important guarantees on its performance. Based on these results, we will describe the detailed steps of the algorithm in the next subsection, showing how this framework can be applied to optimize multi-objective policies in practical scenarios.

### B. Multi-Objective Reward Modeling

Based on the theoretical foundations established in the previous section, we now focus on the practical implementation of Pb-MORL. In particular, we focus on constructing a multi-objective reward model that aligns with human preferences. By utilizing the preference data given by the teacher, we can develop an explicit reward model that captures the complexities of human decision-making.

Inspired by the previous work [11] in the single-objective scenario, we construct a preference predictor $P_\psi[\sigma_0 \succ \sigma_1 | \boldsymbol{w}]$, which is designed to predict the preference $p$ given the pair of segments $\sigma_0$ and $\sigma_1$ under the weight $\boldsymbol{w}$, and is parameterized by $\psi$. The preference predictor $P_\psi$ can be trained by minimizing the cross-entropy loss:

$$\mathcal{L}^{\mathrm{p}} = - \mathbb{E}_{(\sigma_0, \sigma_1, \boldsymbol{w}, p) \sim \mathcal{D}} \Big[ p(0) \log P_\psi[\sigma_0 \succ \sigma_1 | \boldsymbol{w}] + p(1) \log P_\psi[\sigma_1 \succ \sigma_0 | \boldsymbol{w}] \Big]. \quad (13)$$

Utilizing the Bradley-Terry model [11], [43], an explicit reward model $\hat{\boldsymbol{r}}_\psi$ can be constructed to predict the preference as follows:

$$P_\psi[\sigma_1 \succ \sigma_0 | \boldsymbol{w}] = \frac{\exp \sum_t \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}_\psi(s_t^1, a_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}_\psi(s_t^i, a_t^i)}. \quad (14)$$

Eq. (14) models preferences as probabilistic outcomes, thereby accommodating the inherent ambiguity found in human judgments. Specifically, it suggests that the preference is exponentially related to the reward sum over the segment. Then, the reward model $\hat{\boldsymbol{r}}_\psi$ is trained to predict the preference under the weight $\boldsymbol{w}$. Although the estimator $\hat{\boldsymbol{r}}_\psi$ is not inherently a binary classifier, the process of learning this estimator can be regarded as a binary classification, where the preferences $p$ serve as the classification labels.

In the previous discussion, we introduce how to leverage the preference data to construct a reward model. Theoretically, when the reward model $\boldsymbol{r}$ aligns perfectly with the teacher's preferences, we can directly optimize this model to derive the optimal policy. To formalize this relationship, we present the following theorem:

**Theorem 4.** *If the reward model $\hat{\boldsymbol{r}}$ is perfectly aligned with the teacher's preferences, that is, for segments $(\sigma_0, \sigma_1)$ with arbitrary length $H$,*

$$\sigma_0 \succ_{\boldsymbol{w}} \sigma_1 \iff$$
$$\sum_{(s,a) \sim \sigma_0} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t, a_t) > \sum_{(s,a) \sim \sigma_1} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t, a_t). \quad (15)$$

*Since the segment length $H$ can be arbitrarily long, the above equation is equivalent to*

$$\pi_0 \succ_{\boldsymbol{w}} \pi_1 \iff$$
$$\mathbb{E}_{\tau \sim \pi_0} \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t, a_t) > \mathbb{E}_{\tau \sim \pi_1} \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t, a_t). \quad (16)$$

*Then, under a given weight vector $\boldsymbol{w}$, maximizing the discounted return*

$$J(\pi) = \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t, a_t) \quad (17)$$

*is equivalent to selecting the optimal policy $\pi^*(\cdot|\cdot, \boldsymbol{w})$.*

*Proof.* For contradiction, assume that there exists another policy $\pi'$ that performs better than $\pi^*$ under the weight vector $\boldsymbol{w}$, i.e.,

$$\sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t', a_t') > \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t^*, a_t^*), \quad (18)$$

where $(s_t', a_t')$ and $(s_t^*, a_t^*)$ are from the trajectories generated by policies $\pi'$ and $\pi^*$, respectively. In this case, the teacher would prefer the trajectory of $\pi'$ over that of $\pi^*$.

However, since the reward model $\boldsymbol{r}$ is perfectly aligned with the teacher, we have:

$$\pi' \prec_{\boldsymbol{w}} \pi^* \implies$$
$$\sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t', a_t') < \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T \hat{\boldsymbol{r}}(s_t^*, a_t^*), \quad (19)$$

which contradicts the fact that $\pi^*$ is the optimal policy under the reward model $\hat{\boldsymbol{r}}$. Therefore, the assumption is false, and the theorem holds. $\square$

### C. MORL based on Multi-Objective Reward Model

Having established the construction method of the reward model $\hat{\boldsymbol{r}}_\psi$, we now focus on implementing the Pb-MORL algorithm. Specifically, we leverage the learned reward model as a substitute for the traditional reward function, enabling the direct application of existing MORL techniques.

In typical MORL training process, the algorithm collects transitions $(s, a, s', \boldsymbol{r}, \boldsymbol{w})$, which are composed of state $s$, action $a$, next state $s'$, multi-objective reward $\boldsymbol{r}$ and the weight vector $\boldsymbol{w}$. These transitions are then utilized to update value functions and policies. In contrast, our method collects transitions where the reward $\boldsymbol{r}$ is replaced with the predicted reward from the model, $\hat{\boldsymbol{r}}_\psi$. This allows us to align the policy with preference data by minimizing the loss as Eq. (13). This leads to a straightforward implementation of Pb-MORL. We first train the multi-objective reward model, followed by conducting MORL training based on this reward model.
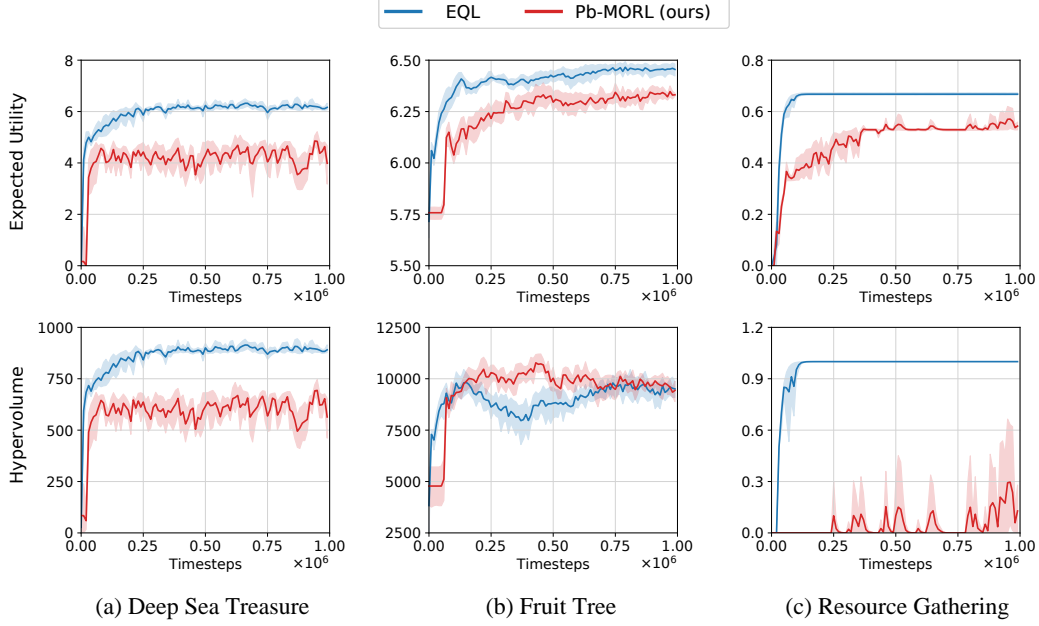
Fig. 2. The training curves of the expected utility and hypervolume on three multi-objective benchmark tasks. The experiments are conducted on 5 random seeds. Blue: the oracle method (EQL). Red: our method.

However, directly using the reward model to train an MORL agent may potentially result in inefficient policy learning:

1) Insufficient amount of preference data: To train a high-quality reward model, a substantial amount of preference data may need to be collected beforehand.
2) Imprecise reward model: When the amount of preference data is insufficient, the reward model may overfit the limited training data, resulting in an imprecise reward model and consequently leading to suboptimal policy performance.

Below are two techniques that can help improve sample efficiency and performance.

- Continuous preference collection: Continuously gather preference data during training, which can enrich the training data of the reward model.
- Relabeling: Relabel historical data with the updated reward model, which can increase the sample efficiency of preference and transition data.

Based on the above techniques, we present Algorithm 3, which is a variant of the Pb-MORL approach discussed above. By integrating the Envelope multi-objective Q-learning (EQL) [33] into our learning process, Algorithm 3 achieves a simple yet effective approach for policy optimization. Specifically, in lines 3-13, the agent interacts with the environment to collect transition data. In lines 14-22, the multi-objective reward model is updated continuously during policy training. In line 23, the rewards of the transition data in the replay buffer are relabeled. In lines 26-29, the Q-function and policy are updated using the EQL method.

## IV. EXPERIMENTAL RESULTS

### A. Setups

In this section, we conduct several experiments to evaluate the effectiveness of the proposed method. We test Pb-MORL on several benchmark multi-objective tasks [33], [35], [44] to demonstrate its effectiveness across diverse multi-objective settings. Additionally, we evaluate Pb-MORL on a custom task for multi-energy storage system charging and discharging as well as an autonomous driving task on a multi-line highway, showing its potential for real-world industrial applications.

**Construction of the multi-objective teacher.** Similar to prior PbRL works [9], [27], [45], in order to systemically evaluate the performance, we construct a "scripted teacher", which provides preferences $p$ between two trajectory segments $\sigma_0, \sigma_1$ under certain weight $w$ according to the task's ground truth reward function. The following paragraph formalizes the process used by the scripted teacher in multi-objective RL.

Let $r_{gt}$ denote the task's ground truth reward function. Then, for the segment pair $(\sigma_0, \sigma_1)$, the scripted teacher first computes the discounted reward sum:

$$\boldsymbol{R}_i = \sum_{t=0}^{H-1} \gamma^t \boldsymbol{r}_{gt}(s_t^i, a_t^i) \quad i = 0, 1, \tag{20}$$

where $\gamma$ is the discount factor, $t$ is the time step, $s_t^i, a_t^i$ represents the state and action of segment $\sigma_i$ in time step $t$, and $H$ is the segment length. Next, the teacher computes the weighted inner product with the given weight vector $w$:

$$R_i = \boldsymbol{w}^T \boldsymbol{R}_i = \boldsymbol{w}^T \left( \sum_{t=0}^{H-1} \gamma^t \boldsymbol{r}_{gt}(s_t^i, a_t^i) \right) \quad i = 0, 1, \tag{21}$$

**Algorithm 3** Pb-MORL algorithm using the EQL method

---

**Input:** Frequency of teacher feedback $K$, number of sampled segment $N_s$, number of sample weights $N_w$ for reward learning, timesteps for learning start $T_0$

**Output:** Multi-objective reward model $\hat{\boldsymbol{r}}_\psi$, multi-objective Q function $\boldsymbol{Q}_\theta$, policy $\pi_\phi(\cdot|\cdot, \boldsymbol{w})$

1: Initialize parameter vectors $\psi, \theta, \phi$
2: **for** each iteration **do**
3:    **for** each environment step $t$ **do**
4:      Obtain current state $s_t$
5:      **if** global step $< T_0$ **then**
6:        Randomly sample action $a_t$
7:      **else**
8:        Obtain action $a_t \sim \pi_\phi(\cdot|s_t, \boldsymbol{w})$ under a weight $\boldsymbol{w}$
9:      **end if**
10:      Obtain transition $(s_t, a_t, s_{t+1}, \boldsymbol{w})$
11:      Obtain multi-objective reward $\hat{\boldsymbol{r}}(s_t, a_t)$
12:      Add $(s_t, a_t, s_{t+1}, \hat{\boldsymbol{r}}(s_t, a_t), \boldsymbol{w})$ into replay buffer $\mathcal{D}$
13:    **end for**
14:    **if** iteration mod K == 0 **then**
15:      Sample $N_s$ query $(\sigma_0, \sigma_1) \sim \mathcal{D}$
16:      Sample $N_w$ weights $\boldsymbol{w}$
17:      Query overseer for preference $p$ for all queries $(\sigma_0, \sigma_1)$ under all $\boldsymbol{w}$
18:      Store all $N_s \times N_w$ $(\sigma_0, \sigma_1, \boldsymbol{w}, p)$ to buffer $\mathcal{D}_p$
19:      **for** each gradient step **do**
20:        Sample minibatch $(\sigma_0, \sigma_1, \boldsymbol{w}, p) \sim \mathcal{D}_p$
21:        Optimize Eq. (13) to update reward model $\hat{\boldsymbol{r}}_\psi$
22:      **end for**
23:      Relabel entire replay buffer $\mathcal{D}$ using $\hat{\boldsymbol{r}}_\psi$
24:    **end if**
25:    **for** each gradient step **do**
26:      Sample a minibatch from replay buffer $\mathcal{D}$
27:      Update Q function $\boldsymbol{Q}_\theta$ by minimizing $|\boldsymbol{Q} - B\boldsymbol{Q}|$ under $(s, a, s', \boldsymbol{r}, \boldsymbol{w}) \sim \mathcal{D}$, as Eq. (6)
28:      Update the Q-learning policy $\pi_\phi(\cdot|\cdot, \boldsymbol{w})$ by maximizing $\boldsymbol{w}^T \boldsymbol{Q}(s, \pi_\phi(\cdot|s, \boldsymbol{w}), \boldsymbol{w})$ under $(s, a, s', \boldsymbol{r}, \boldsymbol{w}) \sim \mathcal{D}$
29:    **end for**
30: **end for**

---

The scripted teacher compares $R_0$ and $R_1$ to determine which segment performs better:

$$p = \begin{cases} 1, & \text{if } R_0 > R_1, \\ 0.5, & \text{if } R_0 = R_1, \\ 0, & \text{if } R_0 < R_1. \end{cases} \quad (22)$$

Since the scripted teacher's preferences directly correspond to the task's ground truth reward, the algorithms can be quantitatively evaluated using the ground truth reward function.

**Evaluation metrics.** We use two metrics to evaluate the empirical performance on each task:

1) **Expected Utility (EU)** [30]: This metric measures the average utility under randomly sampled weights. Let $\boldsymbol{w}$ be a weight vector randomly sampled from the uniform distribution in $\mathcal{W}$ space. Let $U(\pi, \boldsymbol{w})$ represent the

---

TABLE II
HYPERPARAMETER SETTINGS FOR Pb-MORL

| Hyperparameter | Value |
|---|---|
| Preference frequency $K$ | 500 |
| Number of sampled segment $N_s$ | 300 |
| Number of sample weights $N_w$ | 10 |
| Discount factor $\gamma$ | 0.99 |
| Batch size | 256 |
| Learning rate | $3 \times 10^{-4}$ |
| Training timesteps | $1 \times 10^6$ |
| Number of Q network hidden layers | 2 |
| Number of hidden units per layer | 128 |
| Q target update $\tau$ | $1 \times 10^{-4}$ |
| Optimizer | Adam |

utility function of policy $\pi(\cdot|\cdot, \boldsymbol{w})$ under the weight $\boldsymbol{w}$, which is usually the inner product of the discounted return and the weight $\boldsymbol{w}$. The expected utility $\text{EU}(\pi)$ is then defined as:

$$\text{EU}(\pi) = \mathbb{E}_{\boldsymbol{w}} U(\pi, \boldsymbol{w}). \quad (23)$$

Expected Utility is crucial for evaluation, as it comprehensively measures a policy's overall performance across objectives. Unlike hypervolume [46], which focuses on boundary solutions, EU evaluates the policy's average behavior over the entire weight space. Thus, it serves as a more relevant indicator of general performance in many multi-objective tasks.

2) **Hypervolume (HV)** [46] : Given an approximate Pareto Frontier set $\tilde{\mathbf{F}}$ of multi-objective return and a reference point $\boldsymbol{R}_{\text{ref}}$, the hypervolume metric is defined as:

$$\text{HV}(\tilde{\mathbf{F}}, \boldsymbol{R}_{\text{ref}}) = \bigcup_{\boldsymbol{R} \in \tilde{\mathbf{F}}} \text{volume}(\boldsymbol{R}_{\text{ref}}, \boldsymbol{R}), \quad (24)$$

where $\text{volume}(\boldsymbol{R}_{\text{ref}}, \boldsymbol{R})$ is the volume of the hypercube spanned by the reference vector $\boldsymbol{R}_{\text{ref}}$ and the vector $\boldsymbol{R}$. The reference point here is typically an estimation of the worst possible return for all objectives.

**Experimental details.** For implementation details, for line 12 of Algorithm 3, we use a query-policy aligned replay buffer to maintain an accurate reward model in the near-policy region [47]. For line 17, we use the scripted teacher mentioned earlier to generate preference data. The detailed hyperparameter settings of Pb-MORL are shown in Table II. For the baseline, we use EQL [33] as an oracle method, which leverages the ground truth reward function for policy learning.

### B. Experimental Results on Multi-Objective Benchmark Tasks

**Tasks.** We evaluate our method on three multi-objective benchmark tasks [30], each presenting distinct challenges, such as balancing time-cost and total reward or optimizing independently across multiple objectives:

- **Deep Sea Treasure (DST)** [48]: An agent controls a submarine in a 10×11 grid to discover treasures, balancing time and treasure value. The grid contains 10 treasures, with the value increasing with the distance from the starting point $s_0 = (0, 0)$. The multi-objective reward
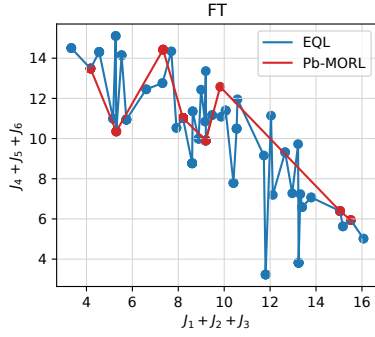
Fig. 3. Visualization of the estimated Pareto frontier of two methods in FT. Note that the actual Pareto frontier in FT has 6 dimensions, we add up the first 3 and last 3 dimensions of rewards for illustration.
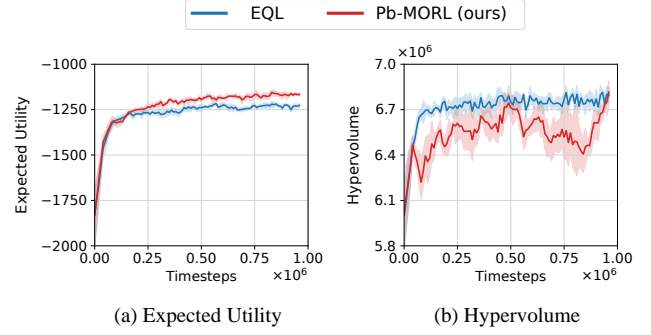


(a) Expected Utility      (b) Hypervolume

Fig. 4. The training curves of the expected utility and hypervolume on the energy system task, averaging over 5 random seeds. Blue: the oracle method (EQL). Red: our method.

$r(s, a)$ has two dimensions: $r_1(s, a)$ for treasure value and $r_2(s, a)$ for time cost, decreasing by 1 for each step.

- **Fruit Tree (FT) [33]**: A full binary tree provides a six-dimensional reward $r \in \mathbb{R}^6$ at each leaf, representing nutritional components: {PROTEIN, CARBS, FATS, VITAMINS, MINERALS, WATER}. The agent maximizes utility for a given weight by selecting the optimal path from root to leaf while choosing between the left and right subtrees.

- **Resource Gathering (RG) [49]**: An agent collects the gold or gem in a 5×5 grid while evading two enemies. Encountering an enemy in the same cell poses a 10% risk of death. The multi-objective reward $r(s, a)$ has three dimensions: $r_1(s, a) = -1$ if being killed, $r_2(s, a) = +1$ if safely returning home with gold, $r_3(s, a) = +1$ if returning with gem.

To justify the selection of $H$ values for each task, we analyze their characteristics. For the DST task, where episode length varies and rewards accumulate over time, we choose $H = 7$ to capture the cumulative effects. In the FT task, with a fixed episode length of 6, $H = 6$ is suitable to encompass the full trajectory. For the RG task, which features sparse rewards and early episode termination, we select $H = 10$ to ensure enough steps are available to differentiate between policies.

Figure 2 presents the expected utility and hypervolume results for the three tasks. In the DST task, our method performs comparably to the oracle in expected utility, demonstrating consistent utility improvement over time. For the FT task, our method matches the oracle in expected utility and surpasses it in hypervolume, indicating effective use of preference for enhancing the Pareto frontier quality. In the RG task, while our method's utility approaches optimal performance, the hypervolume results are less favorable. This may be because the returns of RG are limited to 0 or ±1, and the learned reward model exhibits imprecision in capturing edge-cases under these sparse rewards, which restricts hypervolume growth.

To demonstrate the high quality of the Pareto frontier we learned, we visualized the Pareto frontier learned by EQL and our method in the FT task, as shown in Fig. 3. Our method captures the key factors of the Pareto frontier of the oracle method, showing its effectiveness for application in practice.

### C. Experimental Results on the Custom Energy Task

**The multi-energy management task.** To assess the potential of Pb-MORL for real-world industry applications, we designed a custom multi-objective task for multi-energy storage, simulating the charging and discharging of an energy storage system. The agent controls discharge and charge levels to satisfy external energy demands while balancing cost savings and system lifespan.

- **State space**: The state space includes four scalar values: Current stored energy $s_{\text{storage}}$ (kWh), current energy generated from renewable sources $s_{\text{new}}$ (kWh), external energy demand $s_{\text{demand}}$ (kWh), and the electricity market price $s_{\text{price}}$ (monetary units). Thus, the state vector can be represented as:

$$s = [s_{\text{storage}}, s_{\text{new}}, s_{\text{demand}}, s_{\text{price}}]. \quad (25)$$

- **Action space**: The action $a$ is a scalar indicating the discharge level. Positive values represent energy discharged to meet external demand, while negative values indicate energy charged from renewable sources or the grid.

- **Transition**: After a state transition, the new storage level is calculated as:

$$s_{\text{storage},t+1} = \min(s_{\text{storage}}^{\text{max}}, (s_{\text{storage},t} - a_t)^+), \quad (26)$$

where $s_{\text{storage}}^{\text{max}}$ is the maximum capacity of the energy storage, and $(\cdot)^+$ denote $\max(\cdot, 0)$.

- **Reward function**: The reward is a two-dimensional vector, where the first dimension $r_1(s_t, a_t)$ penalizes the electricity purchasing cost. At each time step, the system may purchase energy to satisfy the external energy demand and charge the storage. The amount of energy bought for charging is

$$g_{\text{charge}} = \begin{cases} (-a - (s_{\text{new}} - s_{\text{demand}})^+)^+ & a < 0 \\ (a - s_{\text{storage}})^+ & a \geq 0 \end{cases}, \quad (27)$$

and that for external demand is

$$g_{\text{demand}} = ((s_{\text{demand}} - s_{\text{new}})^+ - (a)^+)^+, \quad (28)$$

$r_1(s_t, a_t)$ is calculated as $r_1(s_t, a_t) = s_{\text{price}} \times (g_{\text{demand}} + g_{\text{charge}})$. The second dimension $r_2(s_t, a_t)$ indicates a penalty for discharging: $r_2(s_t, a_t) = -1$ when energy is

discharged and 0 otherwise. This design aims to reduce discharges, thus prolonging the system's lifespan.

In this task, rewards are cumulative, with a maximum episode length of 50. Since the agent does not face failures leading to early termination, a sufficiently long $H$ is critical for capturing long-term policy performance. Setting $H = 10$ allows for comprehensive observation of cumulative returns, enabling effective differentiation among policy performances during optimization.

Figure 4 presents the experimental results of the multi-energy management task. Our method surpasses the oracle method in expected utility. This can be attributed to the task's inherent randomness and complex transition dynamics, which make it challenging to directly optimize task rewards like electricity costs or the lifespan loss of system charging. Instead, preference provides a more flexible way to guide policy optimization, allowing the policy to adapt to system complexities more effectively. Additionally, our method matches the oracle in the hypervolume metric.

The weight-conditioned policy $\pi(a|s, w)$ learned by Pb-MORL allows dynamic adaptation to changing user preferences. For instance, in energy management, operators may prioritize cost reduction during peak pricing periods ($w_1 \uparrow$) and system longevity during high-stress operations ($w_2 \uparrow$). Since Pb-MORL trains a single policy conditioned on arbitrary weights $w \in \mathcal{W}$, adapting to such changes only requires modifying the input weight vector $w$ at deployment, eliminating the need for policy retraining. Similarly, in autonomous driving in Section IV-D, safety weights can be adjusted during adverse weather by simply updating $w$. This adaptability minimizes operational overhead and enables Pb-MORL to respond instantly to evolving objectives, making it well-suited for dynamic environments with non-stationary preferences.
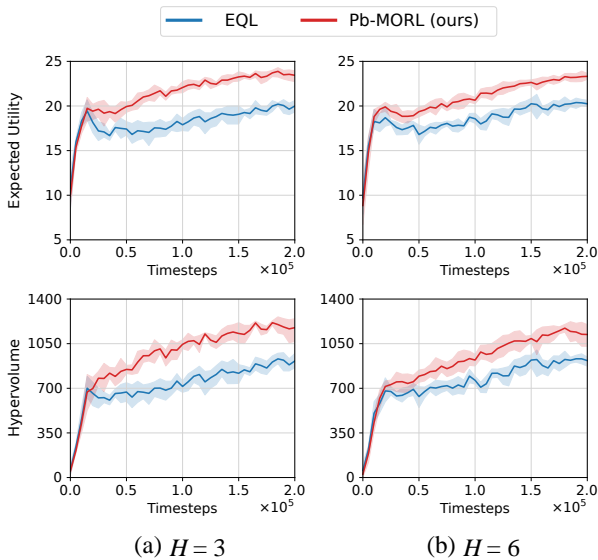


Fig. 5. The training curves of the expected utility and hypervolume on the highway task, averaging over 5 random seeds. Blue: the oracle method (EQL). Red: our method.

### D. Experimental Results on the Multi-Lane Highway Task

**The multi-lane highway task.** To validate the effectiveness of our approach in real-world complex control scenarios, we evaluate it in a multi-lane highway task [30], [50]. In this task, the agent navigates a three-lane highway while driving as quickly as possible, avoiding collisions and prioritizing positioning in the rightmost lane. This setting comprehensively tests the agent's ability to perform in dynamic and multi-faceted environments.

- **State space**: The state is represented by a $V \times 5$ matrix that includes the coordinates and speeds of the ego vehicle and $V - 1$ surrounding vehicles. Each line consists of [presence of the vehicle, $x$ coordinate, $y$ coordinate, $x$ velocity, $y$ velocity].
- **Action space**: Actions are categorized discretely as follows: lane change to the left (0), maintaining the current state (1), lane change to the right (2), acceleration (3), and deceleration (4). These actions are integrated with a lower-level controller for speed and steering.
- **Transition**: The vehicle kinematic is modeled using a simplified kinematic bicycle model, which assumes the left and right wheels function as a single wheel. It regards the front wheel as the primary steering control while omitting sliding effects, thus enabling a more straightforward representation of vehicle dynamics. This model formulation captures the essential dynamics of real-world vehicle behavior, enhancing the simulation's fidelity. The following equations describe the vehicle's motion:

$$\begin{cases} \dot{x} & = v \cos(\theta + \beta), \\ \dot{y} & = v \sin(\theta + \beta), \\ \dot{v} & = a, \\ \dot{\theta} & = \frac{v}{l_r} \tan(\delta), \\ \beta & = \arctan\left(\frac{l_r}{l_f + l_r} \tan(\delta)\right). \end{cases}$$

The surrounding vehicles are controlled by the Intelligent Driver Model (IDM) [51] and the Minimum Overall Braking Distance (MOBIL) model [52].

- **Reward function**: The reward function comprises a three-dimensional vector. The first element represents speed reward, calculated as $\frac{v - v_{\min}}{v_{\max} - v_{\min}}$, where $v$ is the current speed of the ego vehicle, and $v_{\min}$ and $v_{\max}$ denote the minimum and maximum allowable speeds, respectively. The second element indicates lane position reward, as $1$ if the ego vehicle is in the rightmost lane and $0$ otherwise. The third element is a collision penalty, assigned $-1$ upon collision and $0$ otherwise.

We selected $H = 3$ and $H = 6$ for our evaluation. $H = 3$ corresponds to a low-level behavior over a 5-second period, while $H = 6$ represents longer driving behavior, allowing for a more comprehensive assessment of the agent's performance. By evaluating our method with both time horizons, we obtain a more nuanced understanding of its capabilities across different driving scenarios. Following prior works [53], [54], we trained the agent for 200,000 steps.

As shown in Figure 5, our method surpasses the oracle method regarding both expected utility and hypervolume.

In contrast to EQL, which experiences a significant performance decline after an initial increase followed by a slow recovery, our approach maintains stable performance without such setbacks. The initial drop in EQL may be because the agent becomes overly focused on immediate goals, such as speed, resulting in aggressive policies that often neglect safety. Consequently, when the negative impact of collision occurs, the agent must re-explore to find more stable policies. In contrast, Pb-MORL addresses this overfitting through preference-driven reward learning. This method emphasizes the relative benefits of multiple objectives (e.g., "safe overtaking $\succ$ aggressive overtaking") rather than focusing on absolute value differences, enabling the reward model to learn the trade-offs of specific scenarios. Additionally, continuous preference feedback helps to recalibrate the reward model in three phases: balancing objectives, refining scene-specific policies and optimizing for long-tail risks. This approach effectively prevents oscillations caused by conflicting targets.

In summary, Sections IV-C and IV-D demonstrate that our preference-guided policy outperforms the oracle across multiple metrics and adapts more effectively to complex systems than direct task reward optimization. Additionally, the resulting multi-objective policies exhibit strong interpretability, clearly showing how weight vectors impact policy behavior. These findings underscore the potential of our approach for optimizing complex real-world systems.

## V. CONCLUSION

This paper presents the preference-based multi-objective reinforcement learning (Pb-MORL) algorithm, which leverages preference data to overcome the limitations of complicated reward design. Our contributions include a theoretical proof of optimality, showing the Pb-MORL framework can guide the learning of Pareto-optimal policies. In addition, we construct an explicit multi-objective reward model that directly aligns with user preferences, enabling more intuitive decision-making in complex scenarios. Extensive experiments demonstrate the effectiveness and interpretability of Pb-MORL in optimizing various types of multi-objective tasks. Through this work, we highlight the potential of preference-based frameworks in enhancing multi-objective optimization.

Future research can explore several directions. First, we recognize that some assumptions in our work may not hold in practical scenarios. Specifically, for the symmetry, consistency, and transitivity requirements in Assumption 1, we can explore non-transitive cases through pairwise ranking methods [55] and utilize preference aggregation strategies to address violations of other properties. Second, Assumption 2 could be relaxed through active query strategies [26] that optimize comparison requests. Third, to expand its utility in complex systems, we aim to apply Pb-MORL to various domains, such as financial investment and smart manufacturing. Notably, we provide an alternative perspective in Appendix B, discussing the motivation for Pb-MORL, highlighting the impact of human subjectivity in preference data on learning quality in traditional PbRL.

## APPENDIX A
### ADDITIONAL PROOF

*Another proof of Thm.3.* We prove it by providing a constructive Algorithm 4.

If policy $\pi_i$ is never added into $\Pi^*$, there must exists a policy $\pi_a$ $a < i$ and $\boldsymbol{w}_k \in W_I$ s.t. $\boldsymbol{w}_k^T \boldsymbol{r}(\sigma_a) > \boldsymbol{w}_k^T \boldsymbol{r}(\sigma_j)$. Therefore $\pi_i$ is dominated by $\pi_a$, thus not being in the Pareto frontier.

If policy $\pi_i$ is added into $\Pi^*$, and removed when traversing $\pi_b$, there must exists a $\boldsymbol{w}_k \in W_I$ s.t. $\boldsymbol{w}_k^T \boldsymbol{r}(\sigma_b) > \boldsymbol{w}_k^T \boldsymbol{r}(\sigma_j)$. Therefore $\pi_i$ is dominated by $\pi_b$, thus not being in the Pareto frontier.

If policy $\pi_i$ is added into $\Pi^*$, and remains in the $\Pi^*$. Assume there is $\pi_c$ dominates $\pi_i$. If $c > i$, when traversing $\pi_c$, $\pi_i$ must be in $\Pi^*$ and be removed, which contradicts the algorithm. If $c < i$, $\pi_c$ or its dominator must be in $\Pi^*$ when traversing $\pi_i$, and $\pi_i$ must be removed, which also contradicts the algorithm. Therefore the assumption is false.

Summarizing these all concludes the proof.    $\square$

---

**Algorithm 4** Using the teacher to obtain non-convex Pareto frontier, based on insertion sort

---

1: Initialize the estimated Pareto frontier $\Pi^* = \emptyset$
2: **for** each policy $\pi_i \in \Pi$ **do**
3:     **for** each policy $\pi_j \in \Pi^*$ **do**
4:         **if** $\boldsymbol{w}_k^T \boldsymbol{R}(\sigma_i) > \boldsymbol{w}_k^T \boldsymbol{R}(\sigma_j)$ for each $\boldsymbol{w}_k \in W_I$ **then**
5:             Remove $\pi_j(a|s)$ from $\Pi^*$
6:         **else if** $\boldsymbol{w}_k^T \boldsymbol{R}(\sigma_i) < \boldsymbol{w}_k^T \boldsymbol{R}(\sigma_j)$ for each $\boldsymbol{w}_k \in W_I$ **then**
7:             break
8:         **end if**
9:     **end for**
10:     Add $\pi_i$ into $\Pi^*$
11: **end for**

---

## APPENDIX B
### A DIFFERENT PERSPECTIVE ON PB-MORL MOTIVATION

From a different perspective, our motivation for proposing Pb-MORL comes from the common issue of subjective bias in preference data, which exists in traditional PbRL and any application that relies on human expert preferences. For example, in autonomous driving, some people prefer aggressive driving that prioritizes efficiency, while others prefer safer, more cautious driving. Similarly, in large language models using RLHF (Reinforcement Learning from Human Feedback), the preferences assigned by different annotators often conflict due to personal biases. This subjectivity not only affects learning

quality but also makes it difficult for models to accommodate diverse needs.

However, we believe that these scenarios can be reframed as multi-objective optimization problems. To address this, we propose Pb-MORL, which models the problem as a multi-objective optimization task. This approach learns a multi-objective policy that can effectively handle different weights $w$, enabling better use of preference data, improving learning quality, and meeting the diverse needs of individuals.

## REFERENCES

[1] J.-H. Cho, Y. Wang, R. Chen, K. S. Chan, and A. Swami, "A survey on modeling and optimizing multi-objective systems," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 3, pp. 1867–1901, 2017.

[2] Z. Liu, X. Zhang, and B. Jiang, "Active learning with fairness-aware clustering for fair classification considering multiple sensitive attributes," *Information Sciences*, vol. 647, p. 119521, 2023.

[3] X. Guan, C. Song, Y.-C. Ho, and Q. Zhao, "Constrained ordinal optimization—a feasibility model based approach," *Discrete Event Dynamic Systems*, vol. 16, no. 2, pp. 279–299, 2006.

[4] H. R. Baghaee, M. Mirsalim, G. B. Gharehpetian, and H. Talebi, "Reliability/cost-based multi-objective pareto optimal design of stand-alone wind/pv/fc generation microgrid system," *Energy*, vol. 115, pp. 1022–1041, 2016.

[5] X. He and C. Lv, "Toward personalized decision making for autonomous vehicles: a constrained multi-objective reinforcement learning technique," *Transportation research part C: emerging technologies*, vol. 156, p. 104352, 2023.

[6] K. Lee, L. Smith, A. Dragan, and P. Abbeel, "B-pref: Benchmarking preference-based reinforcement learning," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.

[7] F. Memarian, W. Goo, R. Lioutikov, S. Niekum, and U. Topcu, "Self-supervised online reward shaping in sparse-reward environments," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 2369–2375.

[8] P. Mannion, S. Devlin, J. Duggan, and E. Howley, "Reward shaping for knowledge-based multi-objective multi-agent reinforcement learning," *The Knowledge Engineering Review*, vol. 33, p. e23, 2018.

[9] K. Lee, L. M. Smith, and P. Abbeel, "Pebble: Feedback-efficient interactive reinforcement learning via relabeling experience and unsupervised pre-training," in *International Conference on Machine Learning*. PMLR, 2021, pp. 6152–6163.

[10] D. Dewey, "Reinforcement learning and the reward engineering principle," in *2014 AAAI Spring Symposium Series*, 2014.

[11] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," *Advances in neural information processing systems*, vol. 30, 2017.

[12] A. Wilson, A. Fern, and P. Tadepalli, "A bayesian approach for policy learning from trajectory preference queries," *Advances in neural information processing systems*, vol. 25, 2012.

[13] R. S. Sutton, A. G. Barto *et al.*, *Reinforcement learning: An introduction*. MIT press Cambridge, 1998, vol. 1, no. 1.

[14] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.

[15] C. Chang, N. Mu, J. Wu, L. Pan, and H. Xu, "E-mapp: Efficient multi-agent reinforcement learning with parallel program guidance," in *Advances in Neural Information Processing Systems*, 2022.

[16] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, "Mastering the game of go without human knowledge," *nature*, vol. 550, no. 7676, pp. 354–359, 2017.

[17] Y. Luan and Q.-S. Jia, "Simplify twin crane scheduling in railway yard by spatial task assignment," in *2023 China Automation Congress (CAC)*. IEEE, 2023, pp. 3034–3039.

[18] N. Mu, X. Hu, and Q.-S. Jia, "Large-scale data center cooling control via sample-efficient reinforcement learning," in *2024 IEEE 20th International Conference on Automation Science and Engineering (CASE)*. IEEE, 2024.

[19] Y. Luan, Q.-S. Jia, Y. Xing, Z. Li, and T. Wang, "An efficient real-time railway container yard management method based on partial decoupling," *IEEE Transactions on Automation Science and Engineering*, vol. 22, pp. 14 183–14 200, 2025.

[20] N. Mu, X. Hu, and Q.-S. Jia, "Integrating mechanism and data: Reinforcement learning based on multi-fidelity model for data center cooling control," in *2023 China Automation Congress (CAC)*. IEEE, 2023, pp. 5283–5288.

[21] X. Guo, X. Zhang, and X. Zhang, "Incentive-oriented power-carbon emissions trading-tradable green certificate integrated market mechanisms using multi-agent deep reinforcement learning," *Applied Energy*, vol. 357, p. 122458, 2024.

[22] G. Brockman, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.

[23] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, "Exploration by random network distillation," in *International Conference on Learning Representations*, 2018.

[24] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[25] B. Ibarz, J. Leike, T. Pohlen, G. Irving, S. Legg, and D. Amodei, "Reward learning from human preferences and demonstrations in atari," *Advances in neural information processing systems*, vol. 31, 2018.

[26] N. Mu, Y. Luan, Y. Yang, and Q.-S. Jia, "S-epoa: Overcoming the indistinguishability of segments with skill-driven preference-based reinforcement learning," *arXiv preprint arXiv:2408.12130*, 2024.

[27] J. Park, Y. Seo, J. Shin, H. Lee, P. Abbeel, and K. Lee, "Surf: Semi-supervised reward learning with data augmentation for feedback-efficient preference-based reinforcement learning," in *International Conference on Learning Representations*, 2022.

[28] D. J. Hejna III and D. Sadigh, "Few-shot preference learning for human-in-the-loop rl," in *Conference on Robot Learning*. PMLR, 2023, pp. 2014–2025.

[29] N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano, "Learning to summarize with human feedback," *Advances in Neural Information Processing Systems*, vol. 33, pp. 3008–3021, 2020.

[30] F. Felten, L. N. Alegre, A. Nowé, A. L. Bazzan, E.-G. Talbi, G. Danoy, and B. C. da Silva, "A toolkit for reliable benchmarking and research in multi-objective reinforcement learning," in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, 2023, pp. 23 671–23 700.

[31] C. F. Hayes, R. Rădulescu, E. Bargiacchi, J. Källström, M. Macfarlane, M. Reymond, T. Verstraeten, L. M. Zintgraf, R. Dazeley, F. Heintz *et al.*, "A practical guide to multi-objective reinforcement learning and planning," *Autonomous Agents and Multi-Agent Systems*, vol. 36, no. 1, p. 26, 2022.

[32] K. Li and H. Guo, "Human-in-the-loop policy optimization for preference-based multi-objective reinforcement learning," *arXiv preprint arXiv:2401.02160*, 2024.

[33] R. Yang, X. Sun, and K. Narasimhan, "A generalized algorithm for multi-objective reinforcement learning and policy adaptation," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019, pp. 14 636–14 647.

[34] D. M. Roijers, D. Steckelmacher, and A. Nowé, "Multi-objective reinforcement learning for the expected utility of the return," in *Proceedings of the Adaptive and Learning Agents workshop at FAIM*, vol. 2018, 2018.

[35] J. Xu, Y. Tian, P. Ma, D. Rus, S. Sueda, and W. Matusik, "Prediction-guided multi-objective reinforcement learning for continuous robot control," in *International conference on machine learning*. PMLR, 2020, pp. 10 607–10 616.

[36] M. Reymond, E. Bargiacchi, and A. Nowé, "Pareto conditioned networks," in *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 2022, pp. 1110–1118.

[37] L. N. Alegre, A. L. Bazzan, D. M. Roijers, A. Nowé, and B. C. da Silva, "Sample-efficient multi-objective learning via generalized policy improvement prioritization," in *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, 2023, pp. 2003–2012.

[38] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.

[39] A. W. Naylor and G. R. Sell, *Linear operator theory in engineering and science*. Springer Science & Business Media, 1982.

[40] F. S. Melo, "Convergence of q-learning: A simple proof," *Institute Of Systems and Robotics, Tech. Rep*, pp. 1–4, 2001.

[41] K. Li and Q.-S. Jia, "Decentralized multi-agent reinforcement learning: An off-policy method," 2021.

[42] K. Li, X. Jin, Q.-S. Jia, D. Ren, and H. Xia, "An ocba-based method for efficient sample collection in reinforcement learning," *IEEE Transactions on Automation Science and Engineering*, vol. 21, no. 3, pp. 3615–3626, 2024.

[43] R. A. Bradley and M. E. Terry, "Rank analysis of incomplete block designs: I. the method of paired comparisons," *Biometrika*, vol. 39, no. 3/4, pp. 324–345, 1952.

[44] N. Mu, Y. Luan, and Q.-S. Jia, "Preference-based multi-objective reinforcement learning with explicit reward modeling," in *2024 China Automation Congress (CAC)*. IEEE, 2024.

[45] N. Mu, H. Hu, X. Hu, Y. Yang, B. Xu, and Q.-S. Jia, "Clarify: Contrastive preference reinforcement learning for untangling ambiguous queries," in *Proceedings of the 42th International Conference on Machine Learning*, 2025.

[46] E. Zitzler, *Evolutionary algorithms for multiobjective optimization: Methods and applications*. Shaker Ithaca, 1999, vol. 63.

[47] X. Hu, J. Li, X. Zhan, Q.-S. Jia, and Y.-Q. Zhang, "Query-policy misalignment in preference-based reinforcement learning," in *The Twelfth International Conference on Learning Representations*, 2024.

[48] P. Vamplew, R. Dazeley, A. Berry, R. Issabekov, and E. Dekker, "Empirical evaluation methods for multiobjective reinforcement learning algorithms," *Machine learning*, vol. 84, pp. 51–80, 2011.

[49] L. Barrett and S. Narayanan, "Learning all optimal policies with multiple criteria," in *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 41–47.

[50] E. Leurent, "An environment for autonomous driving decision-making," https://github.com/eleurent/highway-env, 2018.

[51] M. Treiber, A. Hennecke, and D. Helbing, "Congested traffic states in empirical observations and microscopic simulations," *Physical review E*, vol. 62, no. 2, p. 1805, 2000.

[52] A. Kesting, M. Treiber, and D. Helbing, "General lane-changing model mobil for car-following models," *Transportation Research Record*, vol. 1999, no. 1, pp. 86–94, 2007.

[53] L. Forneris, A. Pighetti, L. Lazzaroni, F. Bellotti, A. Capello, M. Cossu, and R. Berta, "Implementing deep reinforcement learning (drl)-based driving styles for non-player vehicles," *International Journal of Serious Games*, vol. 10, no. 4, pp. 153–170, 2023.

[54] H. Tian, K. Reddy, Y. Feng, M. Quddus, Y. Demiris, and P. Angeloudis, "Enhancing autonomous vehicle training with language model integration and critical scenario generation," *arXiv preprint arXiv:2404.08570*, 2024.

[55] H. Choi, S. Jung, H. Ahn, and T. Moon, "Listwise reward estimation for offline preference-based reinforcement learning," in *Forty-first International Conference on Machine Learning*, 2024.