
Preference Controllable Reinforcement Learning with Advanced Multi-Objective Optimization

Yucheng Yang¹ Tianyi Zhou² Mykola Pechenizkiy¹ Meng Fang^{3,1}

Abstract

Practical reinforcement learning (RL) usually requires agents to be optimized for multiple potentially conflicting criteria, e.g. speed vs. safety. Although Multi-Objective RL (MORL) algorithms have been studied in previous works, their trained agents often cover limited Pareto optimal solutions and they lack precise controllability of the delicate trade-off among multiple objectives. Hence, the resulting agent is not versatile in aligning with customized requests from different users. To bridge the gap, we develop the ‘‘Preference controllable (PC) RL’’ framework, which trains a preference-conditioned meta-policy that takes user preference as input controlling the generated trajectories within the preference region on the Pareto frontier. PCRL is compatible with advanced Multi-Objective Optimization (MOO) algorithms that are rarely seen in previous MORL approaches. We also proposed a novel preference-regularized MOO algorithm specifically for PCRL. We provide a comprehensive theoretical analysis to justify its convergence and preference controllability. We evaluate PCRL with different MOO algorithms against state-of-the-art MORL baselines in various challenging environments with up to six objectives. In these experiments, our proposed method exhibits significantly better controllability than existing approaches and can generate Pareto solutions with better diversity and utilities.

1. Introduction

Multi-Objective Reinforcement Learning (MORL) has attracted growing interests in applications of training sequential decision-making agents that satisfy multiple objectives. In practice, optimizing for potentially multiple criteria often involves managing trade-offs between them. For example, speed vs. safety or distance vs. energy for robotic control tasks. Many previous MORL methods (Yang et al., 2019; Abels et al., 2019; Xu et al., 2020; Lu et al., 2023; Alegre et al., 2023) tried to address the trade-off issue by optimizing a linearly scalarized objective, which sums up multiple objectives with preference weights. However, the Linear Scalarization (LS) approach’s solutions are limited to a subset of the Pareto optimal solutions, as shown in Fig. 1. As a result, LS solutions have limited optimality and are often not well aligned with the trade-off preferences. However, the Linear Scalarization (LS) approach restricts solutions to a subset of the Pareto-optimal solutions, as illustrated in Fig. 1. Consequently, LS solutions often exhibit limited optimality and may not align well with trade-off preferences. More recently, Lu et al. (2023) sought to enhance LS through reward augmentation, but this can lead to information loss from the original problem. Additionally, Basaklar et al. (2023) proposed optimizing cosine similarity for better preference alignment; however, their approach lacks a mechanism to resolve conflicts between similarity gradients and objective gradients.

This motivates us to develop a novel MORL framework, **Preference Control (PC) RL**, to train a preference controllable agent for user’s preference trade-offs with advanced Multi-Objective Optimization (MOO) algorithms. By leveraging well-established MOO methods (Lin et al., 2019; Mahapatra & Rajan, 2020), our approach overcomes the limitations of LS and enables the discovery of preference-specific solutions along the Pareto front. Moreover, inspired by how certain MOO methods (Désidéri, 2009; Liu et al., 2021; Xiao et al., 2023) deal with conflicting gradients and stochastic gradients, we propose a novel MORL-specific algorithm **PreCo**. We also conduct a comprehensive theoretical analysis, demonstrating that in MORL’s noisy gradient setting, PCRL with PreCo can achieve preference-controlled Pareto-stationary solutions.

¹Department of Mathematics and Computer Science, Eindhoven University of Technology, Eindhoven, The Netherlands
²Department of Computer Science, University of Maryland, College Park, Maryland, The United States ³Department of Computer Science, University of Liverpool, Liverpool, The United Kingdom.
Correspondence to: Yucheng Yang <y.yang@tue.nl>, Meng Fang <Meng.Fang@liverpool.ac.uk>.

Proceedings of the 42nd International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

We conducted experiments in environments with conflicting objectives (Felten et al., 2023) to empirically demonstrate that (1) our PCRL scheme is compatible with various MOO methods; and (2) PCRL with PreCo consistently achieves superior performance across multiple MORL environments. In particular, our method excels in cases with a large number of objectives or conflicting objectives.

2. Preliminaries

Multiple Objective Reinforcement Learning In the multi-objective RL (MORL) setting, agent needs to optimize possibly conflicting objectives with their separate reward function. MORL setting can be modeled as Multi-Objective Markov Decision Process (MOMDP). Unlike the scalar reward function in conventional MDP, the reward function in MOMDP is vector-valued. A MOMDP is defined as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, p_0, \gamma)$, with state space \mathcal{S} and action space \mathcal{A} , dynamics $P(s_{t+1}|s_t, a_t)$, initial state distribution $p_0(s_0)$, and discount factor $\gamma \in [0, 1]$. The vector-valued function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^m$ is a multi-objective reward function with m objectives. A policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ is a function mapping states to actions. The multi-objective value functions for a policy π are:

$$\mathbf{q}^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{i=0}^{\infty} \gamma^i \mathbf{r}(S_{t+i}, A_{t+i}) | S_t = s, A_t = a \right] \quad (1)$$

$$\mathbf{v}^\pi(s) = \mathbb{E}_\pi \left[\sum_{i=0}^{\infty} \gamma^i \mathbf{r}(S_{t+i}, A_{t+i}) | S_t = s \right] \quad (2)$$

Let $\mathbf{v}^\pi \in \mathbb{R}^m$ to be the multi-objective value vector of π under the initial state distribution p_0 :

$$\mathbf{v}^\pi = \mathbb{E}_{S_0 \sim p_0} [\mathbf{q}^\pi(S_0, \pi(S_0))] \quad (3)$$

Each entry of \mathbf{v}^π is a value for an objective. The Pareto Front is a set of nondominated multi-objective value functions $\mathcal{F} := \{\mathbf{v}^\pi \mid \nexists \pi' \text{ s.t. } \mathbf{v}^{\pi'} \succ \mathbf{v}^\pi\}$, where \succ is the relation of Pareto dominance such that $\mathbf{v}^{\pi'} \succ \mathbf{v}^\pi$ means $(\forall i, \mathbf{v}_i^{\pi'} \geq \mathbf{v}_i^\pi) \wedge (\exists j, \mathbf{v}_j^{\pi'} > \mathbf{v}_j^\pi)$. Intuitively, if \mathbf{v}^{π_1} is dominated by \mathbf{v}^{π_2} , then there is no objective where π_1 performs better so π_2 is always a better choice than π_1 . An optimal MORL agent should have its value vector on the Pareto front.

Preference Control Preference quantifies the trade-off among the multiple objectives. We define the set of preferences $\mathcal{P} := \{p \in \mathbb{R}^m : p^T \mathbf{1} = 1, p \succ 0\}$. The desired policy π for preference p should have the value \mathbf{v}^π optimizing a similarity metric $\Psi(p, \mathbf{v}^\pi)$, which can be cosine similarity or what we define in Definition 4.2. The optimal \mathbf{v}^π should be on the Pareto Front with a maximal similarity to p . In other words, the ideal \mathbf{v}^π for preference p should be

on the Pareto front and closest to the intersection of between the Pareto front and the ray from the origin to the direction of p .

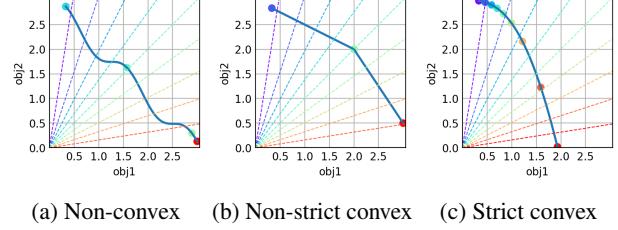


Figure 1: The plots show results and limitations of optimizing the two objectives using LS objective $\max_{\pi_p} \mathbf{p}^T \mathbf{v}^{\pi_p}$. The blue solid curve is the Pareto front, the colored dotted rays are the preference directions of p , and the same colored points are the resulted values \mathbf{v}^{π_p} . The Pareto front in the left is non-convex, which could happen with deterministic policies, the middle is convex but not strictly convex which often happens with discrete action space, and the right is strictly convex. We observe an obvious gap between the preferences and the achieved values in all situations. Linear scalarization can not always discover all Pareto optimal solutions and the discovered solutions are not in the intersection between the preference rays and the Pareto front. This explains why optimizing a similarity $\Psi(\mathbf{p}, \mathbf{v}^\pi)$ is necessary for preference control.

Previous works (Yang et al., 2019; Xu et al., 2020; Alegré et al., 2023) have focused on maximizing a linear scalarization of objectives $\mathbf{p}^T \mathbf{v}^\pi$. However, the solution to $\max_\pi \mathbf{p}^T \mathbf{v}^\pi$ or $\max_\theta \mathbf{p}^T \mathbf{v}^{\pi_\theta}$ is confined to the convex region of the Pareto front (Chapter 4.7, Boyd & Vandenberghe (2004)), excluding the non-convex regions. Lu et al. (2023) demonstrated that for stochastic policies, the Pareto front can be treated as convex. Even so, as shown in Fig. 1b, the solution of linear scalarization (LS) is restricted to a Convex Coverage Set (CCS) (Rojers et al., 2013), which is only a subset of the full Pareto front. Furthermore, even when the Pareto front is strictly convex, LS is not guaranteed to align closely with the direction of p , as illustrated in Fig. 1c. As a result, LS often fails to discover all optimal solutions and is not well-suited for preference alignment.

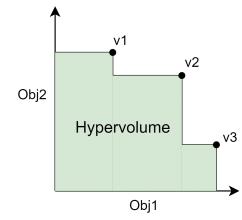


Figure 2: Illustration of hypervolume of three value vectors $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ for a two objective optimization. Their hypervolume is the volume of the union set of their dominated regions (the green shaded area), reflecting their diversity and coverage.

two requirements for the agent. One is to explore the Pareto front as much as possible, and the other is to have a performance trade-off close to the input preference. These two requirements can be evaluated for two metrics: **Hypervolume(HV)** for exploration of Pareto front and **Similarity** $\Psi(\mathbf{p}, \mathbf{v}^\pi)$ for controllability.

Multi-Objective Optimization Methods Previous Multi-Objective Optimization (MOO) methods deal with how to manipulate gradients from multiple objectives so that updating with the manipulated gradient can reach Pareto optimality. A typical method MGDA (Désidéri, 2009) can guarantee to update in a common ascending direction and stops when the Pareto stationary points are reached. Methods such as CAGrad (Liu et al., 2021) and SDMGrad (Xiao et al., 2023) can provide Pareto optimal solutions by linear scalarization with preference as weights. However, as mentioned above, optimizing linearly scalarized objective with weight \mathbf{p} can not guarantee a large similarity $\Psi(\mathbf{p}, \mathbf{v}^\pi)$.

Methods such as PMTL (Lin et al., 2019) and EPO (Mahapatra & Rajan, 2020) apply similarity constraints to reach the Pareto front with the desired preference, so they can be used for preference control purpose. In the next section, we show how these methods can be used for learning $\pi(a|s, \mathbf{p})$ and they will be used as baselines for our proposed new MOO algorithm.

3. Learning preference controllable agent

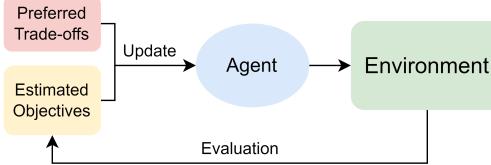


Figure 3: PCRL updates the agent based on its performance and the user preference of objectives.

We propose ‘‘Preference control (PC) RL’’ scheme to incorporate MOO algorithms to handle the trade-offs between multiple conflicting objectives. We train a single agent that can be conditioned on different performance preferences. Conditional preference \mathbf{p} controls the agent’s emphasis on different objectives and corresponds to a desired point on the Pareto front. We denote the policy conditioned on a preference $\pi(\cdot|\cdot, \mathbf{p})$ as $\pi_{\mathbf{p}}$. During training, $\mathbf{v}^{\pi_{\mathbf{p}}}$ is estimated for uniform sampled $\mathbf{p} \in \mathcal{P}$. Then we can evaluate similarity $\Psi(\mathbf{p}, \mathbf{v}^{\pi_{\mathbf{p}}})$ and obtain an update direction for $\pi_{\mathbf{p}}$. In PCRL scheme, the update direction can be obtained using any methods that can incorporate preference on the objectives, including LS (optimizing $\max_{\pi_{\mathbf{p}}} \mathbf{p}^T \mathbf{v}^{\pi_{\mathbf{p}}}$), or other MOO methods with extra optimization or regularization of the similarity (detailed implementations in Appendix B). We

propose a MOO approach specifically for PCRL while these existing MOO methods and LS will be tested as baselines. In the following section, we first introduce how to estimate the $\mathbf{v}^{\pi_{\mathbf{p}}}$ values then explain our proposed update method. We provide theoretical guarantee of our proposed update method in the next section.

3.1. Objective Estimation

Preference control aims to achieve the desired trade-off on conflicting objectives. In the previous RL experiments of MOO methods like (Yu et al., 2020; Liu et al., 2021; Xiao et al., 2023), the loss of the value function is used as the objective for MOO, and equal weight is given to all value losses to balance the multi-objectives. While this may be appropriate for RL tasks with minimal conflict, in our setting for preference control, it is essential to align the objective with the preference, so the objective to be aligned with the preference should be $\mathbf{v}^{\pi_{\mathbf{p}}}$ itself rather than its approximation loss. Here, we show how to estimate $\mathbf{v}^{\pi_{\mathbf{p}}}$ for mainstream RL algorithms.

When learning $\pi_{\mathbf{p}}$ with value-based methods like DDPG (Lillicrap et al., 2016), TD3 (Fujimoto et al., 2018) and SAC (Haarnoja et al., 2018), we can estimate $\mathbf{v}^{\pi_{\mathbf{p}}}$ by

$$\hat{\mathbf{v}}^{\pi_{\mathbf{p}}} = \mathbb{E}_{S_0 \sim p_0} [\mathbf{q}_\theta(S_0, \pi_{\mathbf{p}}(S_0), \mathbf{p})] \quad (4)$$

where \mathbf{q}_θ is multi-objective critic network that outputs a vector of Q-values, it is also conditioned on the preference \mathbf{p} , because \mathbf{p} controls the policy π thus controlling the value \mathbf{q}^π .

For policy-based methods such as A3C (Mnih et al., 2016), PPO (Schulman et al., 2017), they update with a whole episode so $\mathbf{v}^{\pi_{\mathbf{p}}}$ can be estimated by episodic returns.

$$\hat{\mathbf{v}}^{\pi_{\mathbf{p}}} = \mathbb{E}_{S_0 \sim p_0} \left[\sum_{t=0}^T \gamma^t \mathbf{r}(S_t, A_t) \right] \quad (5)$$

As a result, our scheme is applicable to both discrete action space and continuous action space. With the estimated value vector $\hat{\mathbf{v}}^{\pi_{\mathbf{p}}}$, we can evaluate the similarity $\Psi(\mathbf{p}, \hat{\mathbf{v}}^{\pi_{\mathbf{p}}})$.

3.2. Updating Procedure

After estimating objective values and similarity for preference control, we need to manipulate the gradients from different objectives and update the agent using the manipulated gradient. Our scheme has the following updating procedure:

1. Get the Jacobian matrix $\nabla_{\pi_{\mathbf{p}}} \hat{\mathbf{v}}^{\pi_{\mathbf{p}}}$:

Each row of the Jacobian matrix $\nabla_{\pi_{\mathbf{p}}} \hat{\mathbf{v}}^{\pi_{\mathbf{p}}}$ is a gradient for one objective. The gradient can be obtained by

conventional RL methods, such as the policy gradient and the deterministic policy gradient. An illustrative diagram (Fig. 9 in Appendix E) provides intuition and shows how to estimate $\nabla_{\pi_p} \hat{v}^{\pi_p}$ for different RL algorithms.

2. Get similarity gradient $\nabla_{\pi_p} \Psi(\mathbf{p}, \hat{v}^{\pi_p})$:

Ψ should effectively measure the similarity between the preference \mathbf{p} and the estimated value \hat{v}^{π_p} . While cosine similarity could be a reasonable choice for quantifying the closeness between the value vector and the preference, it may conflict with the underlying objectives. To address this, we propose a novel design for the similarity function, combined with a gradient manipulation approach, that is capable of both reaching the Pareto front and enabling effective preference control.

3. Manipulate the gradients and find the optimal update direction d^* by solving:

$$w^* \in \arg \min_{\mathbf{w}} \|\mathbf{d}\|, \quad (6)$$

$$\mathbf{d} \triangleq \nabla_{\pi_p}^T \hat{v}^{\pi_p} \mathbf{w} + \lambda \nabla_{\pi_p} \Psi(\mathbf{p}, \hat{v}^{\pi_p})$$

$$\mathbf{d}^* = \nabla_{\pi_p}^T \hat{v}^{\pi_p} \mathbf{w}^* + \lambda \nabla_{\pi_p} \Psi(\mathbf{p}, \hat{v}^{\pi_p}) \quad (7)$$

This is a min-norm problem similar to MGDA and SMGrad, but it adds a similarity gradient to every objective gradient, making the update not only ascent in a common improving direction but also closing the value v^{π_p} to the preference \mathbf{p} . We call this update *PREFerence COntrol(PreCo)* update. Intuitively, updating with the solution \mathbf{d}^* converges to where $\|\mathbf{d}^*\| = 0$, indicating no direction possible for common improvement of all objectives thus satisfying Pareto stationary. We will prove the convergence of this gradient under our proposed similarity function.

This is a general update procedure that can employ any RL algorithm for the calculation of the objective gradients $\nabla_{\pi_p} \hat{v}^{\pi_p}$. In the third step, the gradient manipulation can also be performed by not only PreCo but also existing MOO algorithms. In the experiment, we examine our scheme with PreCo against the baselines with existing MOOs such as EPO (Mahapatra & Rajan, 2020) CAGrad (Liu et al., 2021). Computationally, the min-norm problem in the third step is solved at the policy level with $\nabla_{\pi_p} \hat{v}^{\pi_p}$ instead of the parameter level with $\nabla_{\theta} \hat{v}^{\pi_p}$. The size for a sample of $\nabla_{\pi_p} \hat{v}^{\pi_p}$ is only $m \times B$ for a batch of B transitions, while $\nabla_{\theta} \hat{v}^{\pi_p}$ of size $m \times M$ could have a parameter size $M \gg m$. M can even be billions for large models. For those cases, solving the min-norm problem at the parameter level could be memory-inefficient and computationally intractable. A pseudo-code for the PCRL scheme with PreCo update and more details on the definitions of policy level gradients and parameter level gradients can be found in Appendix E.

Algorithm 1 PreCo in the theoretical analysis setting

```

1: Initialize: Preference  $\mathbf{p}$ , preference-conditioned policy  $\pi_p$ , and weights  $\mathbf{w}_0$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Rollout and estimate the value to get data  $\xi, \xi', \zeta$ 
4:    $\mathbf{w}_t = \Pi_{\mathcal{W}}(\mathbf{w}_{t-1} - \beta_t [G(\pi_{\mathbf{p}, t}; \xi)^T (G(\pi_{\mathbf{p}, t}; \xi') \mathbf{w}_{t-1} + \lambda_t g_s(\pi_{\mathbf{p}}; \xi'))])$ 
5:    $\pi_{\mathbf{p}, t+1} = \pi_{\mathbf{p}, t} + \alpha_t (G(\pi_{\mathbf{p}, t}; \zeta) \mathbf{w}_{t-1} + \lambda_t g_s(\pi_{\mathbf{p}, t}; \zeta))$ 
6: end for

```

4. Theoretical Analysis

In this section, we provide the formal definition of our proposed similarity function $\Psi(\cdot, \cdot)$ and the theoretical analysis for the PreCo update. We will prove that it converges to Pareto stationary points, and the resulting similarity $\Psi(\mathbf{p}, v^{\pi_p})$ will also converge to stationary points.

Definition 4.1. We define our similarity function as

$$\Psi(\mathbf{p}, \mathbf{v}) = -\frac{1}{2} \left\| \max_i \frac{\mathbf{v}_i}{\mathbf{p}_i} \mathbf{p} - \mathbf{v} \right\|^2. \quad (8)$$

Intuitively, the similarity gradient $\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})$ encourages to focus on the less optimal objectives to reach the preference \mathbf{p} . A visualization for $\Phi(\mathbf{p}, \cdot)$ can be found in Appendix F.

Deep reinforcement learning is inherently stochastic and sensitive to sample complexity. Therefore, we analyze the convergence rate of the proposed PreCo update in the stochastic gradient setting. The PreCo algorithm that we analyze in this case is Algorithm 1, where \mathbf{w} is the coefficient defined in Equation (6) and $\Pi_{\mathcal{W}}$ means the projection to the set $\mathcal{W} := \{\mathbf{w} \in \mathbb{R}^m : \mathbf{w}^T \mathbf{1} = 1, \mathbf{w} \succ 0\}$. Data ξ, ξ', ζ are different noise samples when estimating $G(\pi_p)$ and $g_s(\pi_p)$ gradients, which are defined as:

$$G(\pi_p) = \mathbb{E}[G(\pi_p; \xi)] = \nabla_{\pi_p}^T \hat{v}^{\pi_p} = \mathbb{E}[\nabla_{\pi_p}^T \hat{v}^{\pi_p}], \quad (9)$$

$$\begin{aligned} g_s(\pi_p) &= \mathbb{E}[g_s(\pi_p, \xi)] \\ &= G(\pi_p) \nabla_{\mathbf{v}} \Psi(\mathbf{p}, \hat{v}^{\pi_p}) \\ &= \mathbb{E}[G(\pi_p; \xi) \nabla_{\mathbf{v}} \Psi(\mathbf{p}, \hat{v}^{\pi_p})]. \end{aligned} \quad (10)$$

where the expectation is taken w.r.t. the noise ξ , the i th column of $G(\pi_p; \xi)$ is the gradient of i th objective and $g_s(\pi_p)$ is the similarity gradient.

Algorithm 1 is only for theoretical analysis; In practice, the weight \mathbf{w} does not need to be updated only once every iteration but can be fully optimized for the min-norm problem (6) and a more practical Algorithm 2 is provided in Appendix E.

4.1. Convergence Analysis

First, we define what Pareto stationary is:

Definition 4.2. We define π is an ϵ -accurate Pareto stationary policy if $\mathbb{E}[\min_{\mathbf{w}} \|G(\pi_{\mathbf{p}})\mathbf{w}\|] \leq \epsilon$, where \mathbf{w} is a convex coefficient.

We assume the continuity and smoothness of the objectives.

Assumption 4.1. For every objective $i \in [m]$, $v_i(\pi_{\mathbf{p}})$ is l_i -Lipschitz continuous and $\nabla v_i(\pi_{\mathbf{p}})$ is $l_{i,1}$ -Lipschitz continuous for any preference conditioned policy $\pi_{\mathbf{p}}$.

This assumption is quite common in RL setting. By the “branched returns bound” in (Janner et al., 2019),

$$|v_i(\pi_1) - v_i(\pi_2)| \leq 2r_{\max,i} \left(\frac{\gamma\epsilon_{\pi}}{(1-\gamma)^2} + \frac{\epsilon_{\pi}}{1-\lambda} \right), \quad (11)$$

where $r_{\max,i} = \max_{s,a} r_i(s, a)$ and ϵ_{π} can be any scalar satisfying $\epsilon_{\pi} \geq \max_s D_{TV}(\pi_1(\cdot|s), \pi_2(\cdot|s))$. Because

$$\max_s D_{TV}(\pi_1(\cdot|s), \pi_2(\cdot|s)) \leq D_{TV}(\pi_1, \pi_2) = \frac{1}{2}|\pi_1 - \pi_2|, \quad (12)$$

we can derive

$$|v_i(\pi_1) - v_i(\pi_2)| \leq r_{\max,i} \left(\frac{\gamma}{(1-\gamma)^2} + \frac{1}{1-\lambda} \right) |\pi_1 - \pi_2|, \quad (13)$$

and L_i can be $r_{\max,i} \left(\frac{\gamma}{(1-\gamma)^2} + \frac{1}{1-\lambda} \right)$. Therefore, the Lipschitz continuity of objectives is naturally satisfied for conventional RL settings, and we only need to assume the gradients are also Lipschitz continuous.

Next, we make an assumption on the bias and variance of the stochastic gradient $g_i(\pi; \xi)$.

Assumption 4.2. For every objective $i \in [m]$, the gradients $g_i(\pi_{\mathbf{p}}; \xi)$ is unbiased estimate of $g_i(\pi_{\mathbf{p}})$, and the variances is bounded by $\mathbb{E}_{\xi}[\|g_i(\pi_{\mathbf{p}}; \xi) - g_i(\pi_{\mathbf{p}})\|^2] \leq \sigma^2$.

We also assume bounded gradient.

Assumption 4.3. There exists a constant C_g such that $\|G(\pi_{\mathbf{p}})\| \leq C_g$.

Lemma 4.1. The similarity function $\Psi(\mathbf{p}, \cdot)$ is $(1 + \max_i \frac{\|\mathbf{p}\|}{\|P_i\|})$ -Lipschitz smooth and $g_s(\cdot)$ is Lipschitz continuous under Assumption 4.1 and Assumption 4.3.

This lemma shows that our proposed similarity function is Lipschitz smooth. The detailed proof is in Appendix I.1. PreCo and SDMgrad (Xiao et al., 2023) both belong to MGDA-variant methods that solve a min-norm problem for gradient manipulation. Leveraging the fact that $g_s(\pi_{\mathbf{p}})$ is a positive linear combination of $G(\pi_{\mathbf{p}})$ and the Lipschitz smoothness property, we can therefore build upon their results to prove that PreCo converges to Pareto stationary points.

Theorem 4.1. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function $\Psi(\mathbf{p}, \cdot)$,

we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|] = \mathcal{O}(mT^{-\frac{1}{2}})$. To achieve an ϵ -accurate Pareto stationary point, it requires $T = \mathcal{O}(m^2\epsilon^{-2})$ updates.

Theorem 4.1 shows PreCo converges to Pareto stationary points when λ is a constant. This theorem applies to our proposed similarity function $\Psi(\mathbf{p}, \cdot)$.

Theorem 4.2. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a Lipschitz smooth similarity function with $g'_s(\pi_{\mathbf{p},t})$ being convex combination of $g_i(\pi_{\mathbf{p},t})$ for all t , there can be an increasing $\lambda = \Theta(\log T)$ and we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|] = \mathcal{O}(mT^{-\frac{1}{2}} \log T)$.

Theorem 4.2 consider a case requiring similarity gradient to be a convex combination of objective gradients, of which its design is discussed in Appendix F.2. In this case λ can increase without an upper limit and eventually g_s will dominate the min-norm solution of (6). Proofs are in Appendix I.2.

Remark 4.1. In practice, **Theorem 4.1** still applies to cases where λ increases but with an upper limit. Because after λ gets close to the limit, it can be considered constant. This offers theoretical justification for implementing PreCo with $\Psi(\mathbf{p}, \cdot)$ and an increasing λ .

Remark 4.2. The Pareto front is convex (but not necessarily strictly convex) for MORL (Lu et al., 2023), so local Pareto stationarity (**Definition 4.2**) is equivalent to Pareto optimality.

4.2. Controllability Analysis

Controllability in our setting is the similarity between the desired preference \mathbf{p} and the value $\mathbf{v}^{\pi_{\mathbf{p}}}$ of the preference-conditioned policy $\pi_{\mathbf{p}}$. It is measured by $\Psi(\mathbf{p}, \mathbf{v}^{\pi_{\mathbf{p}}})$. We provide the following results to show how $\mathbf{v}^{\pi_{\mathbf{p}}}$ will converge to the point close to the \mathbf{p} direction.

Theorem 4.3. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function like $\Psi(\mathbf{p}, \cdot)$, we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p}})\|] - \frac{2C_g^2}{\lambda^2} = \mathcal{O}(mT^{-\frac{1}{2}})$.

Theorem 4.3 provides an intuitive result, that with constant λ , the norms of the similarity gradient $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p}})\|]$ will converge and be bounded. The larger λ , the lower the bound $\frac{2C_g^2}{\lambda^2}$, and the closer the solution will reach the stationary points for maximizing similarity.

Theorem 4.4. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function like $\Psi(\mathbf{p}, \cdot)$, there can be an increasing $\lambda = \Theta(T^{\frac{1}{2}})$ and we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p}})\|] = \mathcal{O}(mT^{-\frac{1}{2}} \log T)$.

Theorem 4.4 shows PreCO with increasing λ will converge to the stationary points for the similarity objective. The proofs of the theorems are in Appendix I.3.

Remark 4.3. *Similar to Theorem 4.1, Theorem 4.3 applies to practical implementations where λ increases but with an upper limit.*

Remark 4.4. *The converged stationary points do not guarantee to have always high similarity metrics. For example, when using $\Psi(\mathbf{p}, \cdot)$, our results show $g_s(\pi_p) = G(\pi_p)\nabla_v\Psi(\mathbf{p}, v^{\pi_p})$ converges to 0. However, the value v^{π_p} coincide with the preference \mathbf{p} only when $\|\nabla_v\Psi(\mathbf{p}, v^{\pi_p})\| = 0$. $\|g_s(\pi_p)\|$ can also be 0 when $\|\nabla_v\Psi(\mathbf{p}, v^{\pi_p})\| > 0$, with $G(\pi_p)$ and $\nabla_v\Psi(\mathbf{p}, v^{\pi_p})$ being orthogonal or $G(\pi_p) = 0$. These situations means the points desired by the preference might not exist on the Pareto front. We discuss in practice how to deal with unreachable regions of Pareto front in Appendix H.*

The theoretical results show that PreCo can discover not only Pareto stationary solutions but also preference-specific solutions. Using a 2-objective MOO example with noisy gradients in Appendix G, we highlight PreCo’s ability to find preference-specific solutions in the context of stochastic and conflicting multi-objective problems.

5. Experiments

Beyond the toy example MOO problem in Appendix G, which illustrates the theoretical advantages of PreCo, we conduct experiments in MORL environments to empirically demonstrate the practical effectiveness of PCRL with PreCo.

Benchmarking Environments Common continuous control environments like MO-Hopper and MO-Ant (Felten et al., 2023) feature higher-dimensional spaces but symmetric objectives (e.g., moving north and east), resulting in strictly convex Pareto fronts. In such cases, low-level RL implementations (e.g., HER (Andrychowicz et al., 2017), curriculum design (Alegre et al., 2023)) that improve sample efficiency are more critical than algorithm design. To isolate the effects of different multi-objective methods from implementation details, we focus on empirical analyses in environments with discrete action spaces, more objectives, and non-strictly convex Pareto fronts, such as Fruit-Tree and MO-Reacher. Results for MO-Hopper and MO-Ant are in Appendix C, where we also demonstrate that PreCo is compatible with these low-level implementation improvements.

Fruit-Tree: A discrete environment with up to a 6-dimensional reward, presenting significant challenges due to its six objectives and a non-strictly convex Pareto front. LS-based methods can only find limited Pareto solutions.

MO-Reacher: A robotic control environment with a continuous state space and a discrete action space. The four

objectives are highly conflicting.

Evaluation metrics We evaluated the results using two metrics: **hyperVolume(HV)** for Pareto front exploration and **Cosine Similarity(CS)** for controllability evaluation. They are measured in test time with preference samples unseen in training (Appendix D). We report the mean and standard deviation results of 5 seeds.

Baselines We compare PreCo with existing MOO gradient manipulation methods in the PCRL scheme and existing MORL algorithms.

Linear Scalarization (LS): It optimizes $\max_{\pi_p} \mathbf{p}^T \hat{\mathbf{v}}^{\pi_p}$. Existing MORL methods such as Yang et al. (2019); Xu et al. (2020); Alegre et al. (2023) all optimize the LS objective with implementation-level modifications.

MOO algorithms: Such as EPO (Mahapatra & Rajan, 2020), CAGrad (Liu et al., 2021), SDMgrad (Xiao et al., 2023). They can find a common ascending direction to handle conflicting gradients.

SOTA MORL methods: State-Of-The-Art (SOTA) MORL methods such as CAPQL (CAP) (Lu et al., 2023) and PDMORL (PDM) (Basaklar et al., 2023). They try to address the limitations of the LS-based MORL.

More details about the baselines and their implementations can be found in Appendix B and Table 3.

5.1. Fruit Tree

We evaluate our method against MOO baselines in settings with reward dimensions ranging from 3 to 6. The HV and CS results, presented in Table 1, demonstrate our proposed PreCo outperforms the baselines, particularly in higher-dimensional reward scenarios. The results in Table 2 highlight the significant advantage of our PCRL scheme with EPO/PreCo (ours) over SOTA MORL methods.

This environment exemplifies a scenario where LS methods learn only limited solutions for a **non-strictly convex** Pareto front. As shown in Fig. 4, in the 3-D reward setting, PCRL with PreCo successfully discovers the blue points representing all Pareto solutions. In contrast, the LS agent learns only a single constant v^{π_p} at the red point (a CCS solution), regardless of the preference input. This occurs because the red point optimizes the LS objective $\mathbf{p}^T \hat{\mathbf{v}}^{\pi_p}$ for most \mathbf{p} (shown as red rays). The singular learned value reveals that the LS agent is uncontrollable by \mathbf{p} , despite being conditioned on it.

In the 6-D setting, our implemented LS baseline demonstrates performance identical to another SOTA MORL algorithm, GPI-LS/PD (Alegre et al., 2023). This observation highlights that LS-based methods inherently face an

Method	3D	4D	5D	6D	
LS	0.12 ± 0.01	0.78 ± 0.03	0.33 ± 0.13	0.76 ± 0.05	1.59 ± 0.29
SDMgrad	0.14 ± 0.01	0.78 ± 0.03	0.66 ± 0.02	0.72 ± 0.00	0.74 ± 0.01
EPO	0.15 ± 0.01	0.84 ± 0.02	1.04 ± 0.05	0.89 ± 0.02	3.98 ± 0.48
CAGrad	0.14 ± 0.02	0.78 ± 0.02	0.30 ± 0.06	0.87 ± 0.01	1.23 ± 0.14
PreCo(Ours)	0.15 ± 0.01	0.84 ± 0.02	1.09 ± 0.02	0.91 ± 0.01	4.33 ± 0.21
					0.87 ± 0.01
					15.61 ± 0.75
					0.78 ± 0.03

Table 1: “HV|CS” (higher is better for both) in fruit-tree environment with HV in the scale of 10^3 . Our method consistently achieves the best optimality (HV) and controllability (CS) from 3-6 objectives.

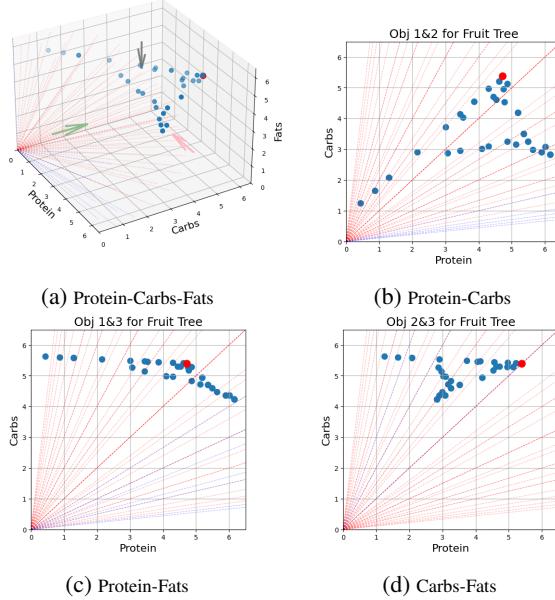


Figure 4: (a) shows the 3-D values v^{π_p} achieved under difference preference input p . Blue points are v^{π_p} of PreCo while LS only learns the red point for some runs. The red point is the optimal LS solution for the preference directions represented by the red rays. (b) shows the Protein-Carbs view, (c) shows the Protein-Fats view, and (d) shows the Protein-fats view.

upper limit when applied to non-strictly convex problems, regardless of lower-level modifications or enhancements. PDMORL explicitly optimizes for cosine similarity, which may introduce conflicts with other objectives. This could explain why PDMORL achieves high CS but fails to deliver competitive HV results..

5.2. MO-Reacher

The MO-Reacher presents a challenging environment with 4-dimensional rewards, requiring the agent to reach four targets illustrated in Fig. 5a. Fig. 5 shows the quantitative results of HV and CS and Fig. 6 shows the state coverage of robotic arm tip positions.

The primary difficulty of this environment lies in the fact that the four objectives are often highly conflicting—moving to

Method	HV	CS
LS/GPI-LS	5.74 ± 0.88	0.72 ± 0.01
CAPQL (entropy coef = 0.01)	5.95 ± 1.12	0.72 ± 0.02
PDMORL (with HER)	9.30 ± 0.08	0.89 ± 0.05
EPO	14.97 ± 2.29	0.77 ± 0.03
PreCo (Ours)	15.61 ± 0.75	0.78 ± 0.03

Table 2: PCRL with advanced MOO algorithms such as EPO and our proposed PreCo outperforms SOTA MORL methods in Fruit-tree environment.

ward one target typically increases the distance from others. Therefore, the LS gradient, $\nabla_{\pi_p}^T \hat{v}^{\pi_p} p$, associated with different preferences p , can exhibit significant conflicts. Moreover, a local optimum for the LS objective $\max_{\pi_p} p \hat{v}^{\pi_p}$ may simply involve the agent staying near the origin—a position that is “equidistant” from all four goals but fails to fully achieve any of them. As shown in Fig. 6a, unsurprisingly, LS and SDMgrad that optimize $\max_{\pi_p} p^T \hat{v}^{\pi_p}$ learned only the ‘equidistant’ local optimum. Moreover, when preferences p are sampled uniformly, the gradients can exhibit significant variance. This explains why CAGrad also fails to learn meaningful results, as demonstrated in Appendix G and discussed in (Xiao et al., 2023), where it is shown that CAGrad struggles with stochastic gradients. Consequently, these methods learn only limited solutions that cannot be effectively controlled by the preference parameter p .

As for existing MORL methods, while both CAPQL and PDMORL address some limitations of LS, they have notable shortcomings. CAPQL relies on reward augmentation, which may lead to suboptimal behavior due to the information loss of the original objectives. Like PCRL with EPO or PreCo, PDMORL optimizes for similarity but directly adds cosine similarity to the original objectives. However, the gradients of cosine similarity and the original objectives can conflict, as it does not leverage conflict-avoidance techniques from MOO algorithms, thereby failing to deliver competitive results in this environment.

Fig. 6b shows the state coverage of PreCo and EPO controlled by 4 different p . From left to right they are

$$[0, 1, 0, 0], [0, 0.67, 0.33, 0], [0, 0.33, 0.67, 0], [0, 0, 1, 0].$$

The preference $[0, 1, 0, 0]$ means full focus on closing to the top target, while $[0, 0, 1, 0]$ full focus on the left target. The first row is EPO and the second is PreCo. Their state

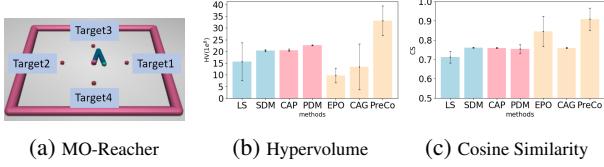


Figure 5: HV and CS of MO-Reacher. The dotted line is the performance of a randomly initialized agent as a reference.

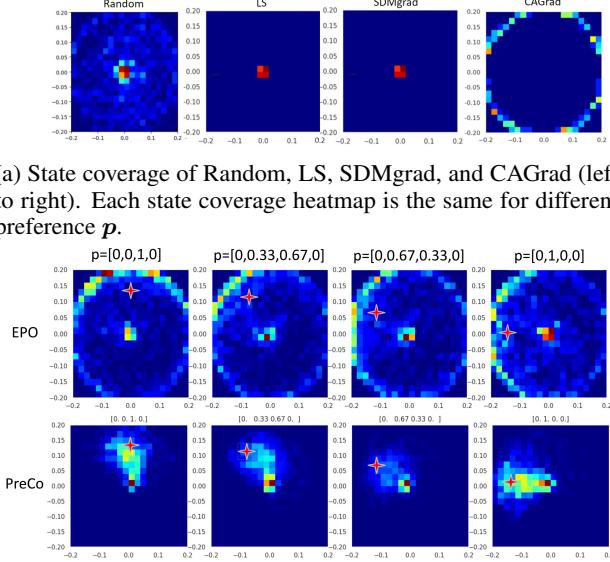


Figure 6: State coverage heatmaps for the positional states of the tip of the robotic arm. Red for higher density of coverage and blue for lower. EPO and PreCo exhibit different state coverage controlled by different p , while random, LS, SDMgrad, and CAGrad show the same state coverage for different preferences.

coverage can be smoothly controlled from more density to the top to more density to the left. This accords with their higher results in CS. More experiments in higher dimensional continuous control environments can be found in Appendix C. Our PCRL scheme with advanced MOO algorithms is shown to have competitive performance against existing MORL methods. Especially when incorporated with our proposed PreCo algorithm, it consistently outperforms existing methods across all environments.

6. Related Works

Existing MORL methods that learn a similar preference-conditioned policy include (Abels et al., 2019; Chen et al., 2019; Lu et al., 2023; Basaklar et al., 2023), of which (Abels et al., 2019; Lu et al., 2023) are LS methods and they care more about discovering all Pareto optimal policies rather than the similarity between the weight input and the resulted value. Chen et al. (2019) employs a setting most similar to our PCRL since they optimize a Tchebycheff Scalarized (TS) (Ehrgott, 2005) objective for solutions aligned

with the preference directions. However, Xu et al. (2020) reported the TS aproach has suboptimal performance in practice, due to the oscillation and stagnation issues, as noted by Mahapatra & Rajan (2021). In particular, for MORL, which is sensitive to stochasticity and conflicting gradients, the convergence of TS can be problematic. Our baseline implementation of EPO with a small constraint threshold can be considered as a version of TS tailored to the MORL setting, designed to mitigate oscillation when near preference direction. PDMORL (Basaklar et al., 2023) incorporates cosine similarity for preference alignment but it is designed specifically for off-policy value-based RL and its performance relies heavily on the HER (Andrychowicz et al., 2017) technique since it often fails without HER as explained by "underrepresentation" in (Basaklar et al., 2023). Another potential reason could be gradient conflicts between the cosine similarity optimization and the original objectives. CAPQL (Lu et al., 2023) addressed LS limitations by adding a concave augmentation term to the reward, transforming the original Pareto front into a strictly convex one. However, this introduces information loss, making the approach sensitive to the augmentation term's magnitude.

To our best knowledge, PCRL is the first preference-conditioned MORL framework to integrate recent advances in Multi-Objective Optimization (MOO). Among the MOO algorithms, PMTL and EPO (Lin et al., 2019; Mahapatra & Rajan, 2020) can be viewed as enhanced versions of TS for preference-specific optimization. CAGrad and SDMGrad (Liu et al., 2021; Xiao et al., 2023) originally optimize only for the average objective, but they can handle conflicting gradients. MGDA and SDMgrad (Désidéri, 2009; Xiao et al., 2023) also address a min-norm problem like (6) for gradient manipulation. However, a significant advantage of our PreCo is that it not only identifies a common ascending direction but also discovers preference-specific solutions. Our theoretical analysis of PreCo builds on some lemmas from Xiao et al. (2023), but incorporating the similarity gradient from Ψ is a novel contribution, which complicates the proof. A more detailed discussion of a broader range of related works can be found in Appendix A.

7. Conclusion

We propose PCRL for preference control in multi-objective trade-offs, integrating recent MOO algorithms into MORL. We also introduce PreCo, a novel MOO approach, with a convergence analysis supporting its ability to learn preference-specific Pareto-optimal solutions and handle stochastic gradients. Experiments across multiple RL environments show that PCRL with PreCo consistently outperforms baselines, demonstrating its effectiveness in learning preference-controllable agents. Future work could enhance the efficiency through curriculum learning.

Acknowledgements

This work was supported by the TKI Smart2 project, a collaboration between KPN and Eindhoven University of Technology (TU/e).

Impact Statement

This work addresses key limitations of existing multi-objective reinforcement learning methods by leveraging recent advances in multi-objective optimization. We propose a general and flexible approach applicable across diverse domains, including gaming agents, robotics, and fine-tuning large language models. While the method itself poses no direct risks, its impact depends on the context and manner in which it is applied.

References

- Abels, A., Roijers, D. M., Lenaerts, T., Nowé, A., and Steckelmacher, D. Dynamic weights in multi-objective deep reinforcement learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 11–20. PMLR, 2019. URL <http://proceedings.mlr.press/v97/abels19a.html>.
- Alegre, L. N., Bazzan, A. L. C., Roijers, D. M., Nowé, A., and da Silva, B. C. Sample-efficient multi-objective learning via generalized policy improvement prioritization. In *Proceedings of the 22nd International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pp. 2003–2012. International Foundation for Autonomous Agents and Multiagent Systems (IFAAMAS), 2023. URL <https://arxiv.org/abs/2301.07784>.
- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, P., and Zaremba, W. Hindsight experience replay. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, 2017. URL <https://papers.nips.cc/paper/7090-hindsight-experience-replay>.
- Basaklar, T., Gumussoy, S., and Ogras, U. Y. Pd-morl: Preference-driven multi-objective reinforcement learning algorithm. In *Proceedings of the 11th International Conference on Learning Representations (ICLR)*, 2023. URL <https://arxiv.org/abs/2208.07914>.
- Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Brahmanage, J. C., Ling, J., and Kumar, A. Flowpg: Action-constrained policy gradient with normalizing flows. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Chen, X., Ghadirzadeh, A., Björkman, M., and Jensfelt, P. Meta-learning for multi-objective reinforcement learning. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2019, Macau, SAR, China, November 3-8, 2019*, pp. 977–983. IEEE, 2019. doi: 10.1109/IROS40897.2019.8968092. URL <https://doi.org/10.1109/IROS40897.2019.8968092>.
- Désidéri, J.-A. Multiple-Gradient Descent Algorithm (MGDA). Research Report RR-6953, June 2009. URL <https://inria.hal.science/inria-00389811>. In this report, the problem of minimizing simultaneously n smooth and unconstrained criteria is considered. A descent direction common to all the criteria is identified, knowing all the gradients. An algorithm is defined in which the optimization process is carried out in two phases : one that is cooperative yielding to the Pareto front, and the other optional and competitive.
- Ehrgott, M. *Multicriteria optimization*. 2005.
- Eysenbach, B., Gupta, A., Ibarz, J., and Levine, S. Diversity is all you need: Learning skills without a reward function. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, 2019. URL <https://arxiv.org/abs/1802.06070>.
- Eysenbach, B., Salakhutdinov, R., and Levine, S. The information geometry of unsupervised reinforcement learning. In *Proceedings of the 10th International Conference on Learning Representations (ICLR)*, 2022. URL <https://arxiv.org/abs/2110.02719>.
- Felten, F., Alegre, L. N., Nowé, A., Bazzan, A. L. C., Talbi, E. G., Danoy, G., and Silva, B. C. da. A toolkit for reliable benchmarking and research in multi-objective reinforcement learning. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS 2023)*, 2023.
- Fujimoto, S., van Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In Dy, J. G. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1582–1591. PMLR, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>.

- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Dy, J. G. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1856–1865. PMLR, 2018. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- Hansen, S., Dabney, W., Barreto, A., Warde-Farley, D., de Wiele, T. V., and Mnih, V. Fast task inference with variational intrinsic successor features. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=BJeAHkrYDS>.
- Hung, W., Huang, B.-K., Hsieh, P.-C., and Liu, X. Q-pensieve: Boosting sample efficiency of multi-objective RL through memory sharing of q-snapshots. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=AwWaBXLIJE>.
- Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust your model: Model-based policy optimization. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 12498–12509, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.
- Lee, Y., Sun, S., Somasundaram, S., Hu, E. S., and Lim, J. J. Composing complex skills by learning transition policies. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=rygrBhC5tQ>.
- Li, Z., Li, T., Smith, V., Bilmes, J., and Zhou, T. Many-objective multi-solution transport. In *The Thirteenth International Conference on Learning Representations (ICLR)*, Boston, MA, USA, 2025.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In Bengio, Y. and Le-Cun, Y. (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1509.02971>.
- Lin, X., Zhen, H., Li, Z., Zhang, Q., and Kwong, S. Pareto multi-task learning. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 12037–12047, 2019.
- Lin, X., Zhang, X., Yang, Z., Liu, F., Wang, Z., and Zhang, Q. Smooth tchebycheff scalarization for multi-objective optimization. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, 2024. URL <https://arxiv.org/abs/2402.19078>.
- Liu, B., Liu, X., Jin, X., Stone, P., and Liu, Q. Conflict-averse gradient descent for multi-task learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Liu, M., Zhu, M., and Zhang, W. Goal-conditioned reinforcement learning: Problems and solutions. In Raedt, L. D. (ed.), *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pp. 5502–5511. ijcai.org, 2022. doi: 10.24963/IJCAI.2022/770. URL <https://doi.org/10.24963/ijcai.2022/770>.
- Lu, H., Herman, D., and Yu, Y. Multi-objective reinforcement learning: Convexity, stationarity and pareto optimality. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=TjEzIsyEsQ6>.
- Mahapatra, D. and Rajan, V. Multi-task learning with user preferences: Gradient descent with controlled ascent in pareto optimization. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 6597–6607. PMLR, 2020. URL <http://proceedings.mlr.press/v119/mahapatra20a.html>.
- Mahapatra, D. and Rajan, V. Exact pareto optimal search for multi-task learning: Touring the pareto front. *CoRR*, abs/2108.00597, 2021. URL <https://arxiv.org/abs/2108.00597>.
- Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X., and Gao, J. Large language models: A survey. *arXiv preprint arXiv:2402.06196*, 2024.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T. P., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In Balcan, M. and Weinberger, K. Q. (eds.), *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA*, 2016.

- June 19-24, 2016, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 1928–1937. JMLR.org, 2016. URL <http://proceedings.mlr.press/v48/mnihal6.html>.
- Nam, T., Sun, S., Pertsch, K., Hwang, S. J., and Lim, J. J. Skill-based meta-reinforcement learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=jeLW-Fh9bV>.
- Narvekar, S., Peng, B., Leonetti, M., Sinapov, J., Taylor, M. E., and Stone, P. Curriculum learning for reinforcement learning domains: A framework and survey. *J. Mach. Learn. Res.*, 21:181:1–181:50, 2020. URL <https://jmlr.org/papers/v21/20-212.html>.
- Nguyen, D., Chen, J., and Zhou, T. Multi-objective linguistic control of large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 4336–4347, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.257. URL <https://aclanthology.org/2024.findings-acl.257>.
- Portelas, R., Colas, C., Weng, L., Hofmann, K., and Oudeyer, P. Automatic curriculum learning for deep RL: A short survey. In Bessiere, C. (ed.), *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pp. 4819–4825. ijcai.org, 2020. doi: 10.24963/IJCAI.2020/671. URL <https://doi.org/10.24963/ijcai.2020/671>.
- Rojers, D. M., Vamplew, P., Whiteson, S., and Dazeley, R. A survey of multi-objective sequential decision-making. *J. Artif. Intell. Res.*, 48:67–113, 2013. doi: 10.1613/JAIR.3987. URL <https://doi.org/10.1613/jair.3987>.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Sekar, R., Rybkin, O., Daniilidis, K., Abbeel, P., Hafner, D., and Pathak, D. Planning to explore via self-supervised world models. In *International Conference on Machine Learning*, pp. 8583–8592. PMLR, 2020.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575 (7782):350–354, 2019.
- Wang, Z., Hunt, J. J., and Zhou, M. Diffusion policies as an expressive policy class for offline reinforcement learning. In *Proceedings of the 11th International Conference on Learning Representations (ICLR)*, 2023. URL <https://arxiv.org/abs/2208.06193>.
- Xiao, P., Ban, H., and Ji, K. Direction-oriented multi-objective learning: Simple and provable stochastic algorithms. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Xu, J., Tian, Y., Ma, P., Rus, D., Sueda, S., and Matusik, W. Prediction-guided multi-objective reinforcement learning for continuous robot control. In *International conference on machine learning*, pp. 10607–10616. PMLR, 2020.
- Yang, R., Sun, X., and Narasimhan, K. A generalized algorithm for multi-objective reinforcement learning and policy adaptation. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 14610–14621, 2019.
- Yang, R., Lu, Y., Li, W., Sun, H., Fang, M., Du, Y., Li, X., Han, L., and Zhang, C. Rethinking goal-conditioned supervised learning and its connection to offline rl. In *Proceedings of the 10th International Conference on Learning Representations (ICLR)*, 2022. URL <https://arxiv.org/abs/2202.04478>.
- Yang, R., Pan, X., Luo, F., Qiu, S., Zhong, H., Yu, D., and Chen, J. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, 2024a. URL <https://arxiv.org/abs/2402.10207>.
- Yang, Y., Zhou, T., He, Q., Han, L., Pechenizkiy, M., and Fang, M. Task adaptation from skills: Information geometry, disentanglement, and new objectives for unsupervised reinforcement learning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, spotlight, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024b. URL <https://openreview.net/forum?id=zSxpKh1yS>.

Yu, T., Kumar, S., Gupta, A., Levine, S., Hausman, K., and Finn, C. Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*, volume 33, pp. 5824–5836, 2020. URL <https://arxiv.org/abs/2001.06782>.

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

A. Related Works

A.1. Meta-policies

Existing MORL methods that learn a similar meta-policy conditioned on a weight or preference include (Abels et al., 2019; Chen et al., 2019; Lu et al., 2023), of which (Abels et al., 2019; Lu et al., 2023) are LS methods and they care more about discovering all Pareto optimal policies rather than the similarity between the weight input and the resulted value. (Chen et al., 2019) employs a setting most similar to our PCRL since they optimize a Tchebycheff Scalarized (TS) (Ehrgott, 2005) objective for solutions aligned with the preference directions. However, (Xu et al., 2020) reported that (Chen et al., 2019) has suboptimal performance in practice. This might be due to the oscillation and stagnation issue inherent in the TS approach, as noted by (Mahapatra & Rajan, 2021). In particular, for MORL, which is sensitive to stochasticity and conflicting gradients, the convergence of TS can be problematic. Our baseline implementation of EPO with a small constraint threshold can be considered as a version of TS tailored to the MORL setting, designed to mitigate oscillation when near preference direction. More recent MORL algorithm (Basaklar et al., 2023) incorporates cosine similarity for preference alignment but it is designed specifically for off-policy value-based RL and its performance relies heavily on the HER (Andrychowicz et al., 2017) technique. In contrast, our PCRL with PreCo is a broader MORL framework compatible with both on-policy and off-policy RL, capable of learning quality policies without HER. Nonetheless, HER can still be integrated into PCRL with PreCo in off-policy settings to enhance sample efficiency. While not targeting exact preference alignment, (Lu et al., 2023) addressed LS limitations by adding a concave augmentation term to the reward, transforming the original Pareto front into a strictly convex one. However, this introduces information loss, making the approach sensitive to the augmentation term's magnitude. Their implementation is limited to SAC, using policy entropy as the augmentation term. Compared to (Lu et al., 2023), our method possesses the theoretical advantages of no information loss and exact preference alignment, which can be empirically demonstrated by our additional experiment in Appendix B.3. Because of the memory-efficient design of our PCRL with PreCo, they are also suitable for multi-objective fintuning of Large language models, like Nguyen et al. (2024); Yang et al. (2024a).

Other RL paradigms employing meta-policies include Goal-Conditioned RL (GCRL) (Sekar et al., 2020; Yang et al., 2022; Liu et al., 2022) and skill-based RL (SBRL) (Nam et al., 2022; Lee et al., 2019). GCRL is controlled by an additional input of a target state that it aims to reach. SBRL is conditioned by a skill latent z that often has a lower dimension than the state for a specific primitive skill. Similar to SBRL, the skill learning methods of Unsupervised Reinforcement Learning (Eysenbach et al., 2019; Hansen et al., 2020) learns skills without external task rewards by optimizing a Mutual Information Skill Learning (MISL) (Eysenbach et al., 2022; Yang et al., 2024b) objective $I(s; z) = H(s) - H(s|z)$. Maximizing $I(s; z)$ encourages the state space coverage to be high and the state distribution to be certain when controlled by a skill z . The concept of Preference Control (PC) has a resemblance to optimizing $I(s; z) = H(s) - H(s|z)$. The purpose of PCRL can also be interpreted as optimizing $I(v; p) = H(v) - H(v|p)$ to encourage diverse values on the Pareto front and the distribution of the values needs to be controlled by preference p .

A.2. Multi-objective optimization

We have already introduced PMTL and EPO (Lin et al., 2019; Mahapatra & Rajan, 2020) that could find preference-specific solutions and CAGrad, SDMGrad (Liu et al., 2021; Xiao et al., 2023) that optimize for the average objective but can deal with conflicting gradients. (Désidéri, 2009; Xiao et al., 2023) has the most similarity to our proposed PreCo because they all solve a min norm problem like 6 for gradient manipulation. The advantage of our PreCo is not only like SDMGrad, which can provably deal with stochastic gradients, but also can follow a preference like EPO. Our theoretical analysis of PreCo is based on some results from (Xiao et al., 2023), but incorporating the similarity gradient from Ψ makes it more complicated and novel. Li et al. (2025) is a method large objective numbers, which could also inspire future MORL design for simultaneously achieving a large number of tasks.

A.3. Training schemes

We uniformly sample p for every episode during training. Techniques from curriculum reinforcement learning (Narvekar et al., 2020; Portelas et al., 2020) can also potentially improve the training of PCRL by using a progressing p preference distribution instead of uniform $p \in \mathcal{P}$.

Baselines	LS obj	Similarity obj	Conflict-avoidance	Characteristic
LS	✓			Yang et al. (2019); Xu et al. (2020); Alegre et al. (2023) are LS with implementation-level modifications. It is used as an ablation for PreCo with LS regularization instead of similarity regularization.
SDMgrad	✓		✓	Its reward augmentation can cause information loss of the original problem.
CAPQL				
PDMORL		✓		It directly adds cosine similarity in the objectives without considering potential gradient conflict.
EPO		✓	✓	Similar to Tchebycheff Scalarization (TS), but it addresses the oscillation issue of TS.
CAGrad		✓	✓	Only convergence guarantee with average gradient regularization, but we implemented it with similarity.
PreCo	✓	✓	✓	Our proposed algorithm specifically for preference controllable MORL with similarity.

Table 3: PCRL with advanced MOO algorithms outperforms SOTA MORL methods in Fruit-tree environment. PDMORL results are means reported in Basaklar et al. (2023).

B. Implementation of baseline methods

We modify existing MOO algorithms EPO (Mahapatra & Rajan, 2020), CAGrad (Liu et al., 2021), and SDMGrad (Xiao et al., 2023) for our proposed PCRL scheme and use them as baselines for our proposed PreCo algorithm. In this section, we explain how these existing baseline MOO algorithms work. Appendix E offers more details for practical implementation.

B.1. Linear scalarization preference control

Linear Scalarization (LS) For a preference p conditioned policy π_p , it is updated by $d = \nabla \hat{v}^{\pi_p}$. Equivalently, it can be $d = p^T \nabla \hat{v}^{\pi_p}$, which means to linearly combine the objective gradients with a coefficient equal to the preference p .

SDMGrad Similar to PreCo that needs to solve the min-norm problem (6) for update direction, Our implementation of SDMGrad solves:

$$\min_{\mathbf{w}} \|\nabla_{\pi_p}^T \hat{v}^{\pi_p} \mathbf{w} + \lambda p^T \nabla_{\pi_p} \hat{v}^{\pi_p}\| \quad (14)$$

$$d = \nabla_{\pi_p}^T \hat{v}^{\pi_p} \mathbf{w}^* + \lambda \nabla_{\pi_p} \hat{v}^{\pi_p}, \quad (15)$$

where w^* is the solution for problem (14). The update direction is for π_p is d . We can see that this SDMGrad implementation and our proposed PreCo differ only in the term multiplied by λ . SDMGrad uses LS gradient for preference alignment while PreCo uses gradient of our proposed similarity function Ψ . This is why SDMGrad can be used as a case for ablation study of our method.

B.2. Similarity preference control

Exact Pareto Optimal (EPO) MOO methods such as PMTL (Lin et al., 2019) and EPO (Mahapatra & Rajan, 2020) apply similarity constraints and have two modes for situations of low and high similarity. Based on this idea, we implement the EPO baseline as: When similarity is low, only similarity gradients $\nabla_{\pi_p} \Psi(\mathbf{w}, \hat{v}^{\pi_p})$ will be used for update. When similarity is high enough, a common ascent direction calculated by MGDA (Désidéri, 2009) is used for update.

$$d = \nabla_{\pi_p} \Psi(\mathbf{w}, \hat{v}^{\pi_p}), \quad \text{if } \Psi'(\mathbf{w}, \hat{v}^{\pi_p}) > \epsilon, \quad (16)$$

$$d = \nabla_{\pi_p}^T \hat{v}^{\pi_p} \arg \min_{\mathbf{w}} \|\nabla_{\pi_p}^T \hat{v}^{\pi_p} \mathbf{w}\|, \quad \text{if } \Psi'(\mathbf{w}, \hat{v}^{\pi_p}) \leq \epsilon, \quad (17)$$

where ϵ is a threshold of similarity and Ψ' can be cosine similarity or our proposed Ψ . Equation (17) is the min-norm update from (Désidéri, 2009), which is equivalent to finding a common ascent direction that maximizes the least improvement among the objectives:

$$d = \arg \max_d \min_i \nabla_{\pi_p}^T \hat{v}_i^{\pi_p} d \quad (18)$$

Because most of the time during training, similarity is not high enough and only similarity gradients are applied in updates, this implementation of EPO can also be seen as an implementation of a relaxed Tchebycheff Scalarization, which avoids gradient oscillation as claimed in (Mahapatra & Rajan, 2020).

Conflict-Averse Gradient (CAGrad) CAGrad (Liu et al., 2021) tries to find a common ascent direction that is not too far from the average gradient. In our setting, we modify it to be a common ascent direction not is not too far from the similarity

gradient.

$$\mathbf{d} = \arg \max_d \min_i \nabla_{\pi_p}^T \hat{\mathbf{v}}_i^{\pi_p} d \quad \text{s.t. } \|\mathbf{d} - \nabla_{\pi_p} \Psi(\mathbf{w}, \hat{\mathbf{v}}^{\pi_p})\| \leq c \nabla_{\pi_p} \Psi(\mathbf{w}, \hat{\mathbf{v}}^{\pi_p}), \quad (19)$$

where $c \in \{r \in \mathcal{R} \mid 0 < r < 1\}$ is a constraint constant to keep d close to the similarity gradient $\nabla_{\pi_p} \Psi(\mathbf{w}, \hat{\mathbf{v}}^{\pi_p})$. This implementation might not apply to the convergence analysis in (Liu et al., 2021). However, as shown by the empirical results, it works in practice for our PCRL scheme.

B.3. Implementation of SOTA MORL methods

Implementation of CAPQL (Lu et al., 2023) has tried to address LS limitations by adding a concave augmentation term to the reward, transforming the original Pareto front into a strictly convex one. Then for this "more convex" new problem, LS can find more optimal solutions. Their implementation only included SAC (Haarnoja et al., 2018), as the entropy maximization in SAC serves as the reward augmentation. To ensure a fair comparison independent of settings, code-level implementations, and algorithmic techniques (such as HER (Andrychowicz et al., 2017)), we modified our original LS with PPO into a "maximum entropy PPO". The modified multi-objective advantages are:

$$\hat{\mathbf{A}}(s, a) = \mathbf{R} + E - \hat{\mathbf{v}}(s, p) \quad (20)$$

where \mathbf{R} represents the vector of multi-objective episodic returns, E denotes the sum of future policy entropies in the sampled episode, and $\hat{\mathbf{v}}(s, p)$ is the multi-objective vector value conditioned on preference p , approximating both expected returns and entropies:

$$\min_{\hat{\mathbf{v}}} \mathbb{E} [||\mathbf{R} + E - \hat{\mathbf{v}}(s, p)||^2] \quad (21)$$

With these modifications, our modified "maximum entropy multi-objective PPO" with Linear Scalarization(LS) is optimizing the concave augmented objective in Eq.(10) from Lu et al. (2023).

To showcase the advantage of our PCRL (Ours) framework with similarity-based methods EPO and PreCo (ours), we test in the 'simple but hard' fruit-tree environment. It is simple for RL due to small discrete state and action spaces but challenging for MORL with 6 objectives and a non-strictly convex Pareto front. This comparison isolates MORL performance from lower-level RL factors, directly highlighting our method's strengths.

α/Method	Hyper volume(1e3)	Cosine Similarity
0 (LS)	5.74 ± 0.88	0.718 ± 0.040
0.01	5.95 ± 1.12	0.722 ± 0.006
0.05	5.18 ± 0.36	0.718 ± 0.040
0.10	1.75 ± 1.35	0.633 ± 0.141
EPO	14.97 ± 2.29	0.77 ± 0.03
PreCo(ours)	15.61 ± 0.75	0.78 ± 0.03

Table 4: HV and CS performance in 6D Fruit-tree environment, the HV value has a unit of $1e^3$. The comparison is between PCRL (Ours) framework with similarity-based methods such as EPO and PreCo (Ours) and LS with different strength of concave augmentation from (Lu et al., 2023)

The results in Table 4 show when the augmentation strength $\alpha = 0.01$, the performance of CAPQL-modified PPO is marginally better than the original LS ($\alpha = 0$), but still significantly worse than similarity-based methods (EPO, PreCo(ours)). Larger α values lead to performance drops. This result aligns with Remark 5 and Figure 9 in (Lu et al., 2023), which highlights that such augmentation can cause information loss in the original problem, and excessive augmentation results in performance degradation. In contrast, our method has the theoretical advantage of overcoming the LS limitation without any reward augmentation, thus avoiding information loss from the original problem.

Implementation of PDMORL (Basaklar et al., 2023) We implement PDMORL using its official GitHub codebase.

Unfortunately, we encountered several issues with their implementation. It is computationally inefficient, requiring 14 hours to complete training in the Fruit-Tree environment and over 70 hours for MO-Hopper, even when using an NVIDIA A100

GPU and 72 CPU cores. In contrast, our PCRL implementation completes training in just 20 minutes for Fruit-Tree and 4 hours for MO-Hopper.

Additionally, while we successfully reproduced the Fruit-Tree results reported in their paper, we were unable to reproduce the MO-Hopper results. The training appears to fail due to the "underrepresentation" problem mentioned in their paper. This aligns with the intuition that directly adding cosine similarity to the original objectives, as done in PDMORL, can introduce gradient conflicts, especially since it lacks conflict-avoidance mechanisms commonly used in MOO algorithms.

Due to these computational inefficiencies and reproducibility issues, we are unable to compare PDMORL for MO-Hopper and MO-Ant.

C. Additional Experimental results

In this section, we present experimental results in high-dimensional continuous environments and demonstrate the compatibility of our method with low-level sample efficiency modifications.

C.1. Continuous Environments

As discussed in [Section 5](#), typical robotic control MORL environments, such as MO-Hopper and MO-Ant, present challenges for sample-efficient RL but often feature simple and strictly convex Pareto fronts. As a result, these environments fail to fully reveal the limitations of common LS-based MORL methods. Consequently, the advantages of our PCRL and PreCo are less pronounced in these settings.

To provide a more comprehensive evaluation, we introduce an additional metric—Overall Non-dominated Ratio (ONR)—to demonstrate that, while our method may not show a clear advantage in terms of HV and CS, it still outperforms existing approaches in ONR.

In our setting, the ONR of one method is defined as:

$$ONR = \frac{\text{the number of this method's samples that are non-dominated by all samples of all methods}}{\text{total sample number of this method}} \quad (22)$$

We attempted to implement PDMORL using its official GitHub codebase but encountered reproducibility issues and computational inefficiencies mentioned in [Appendix B.3](#). As a result, we were unable to compare PDMORL in the continuous control environments.

C.1.1. MO-HOPPER

The MO-Hopper is a classic continuous robotic control environment, with one objective rewards for going forward in the x-axis, and the other rewards for jumping high in the z-axis as shown in [Fig. 7a](#). The two objectives are less symmetric than MO-Ant and there is a clear trade-off in directions. The HV and CS results are shown in [Fig. 7](#). LS and CAPQL have

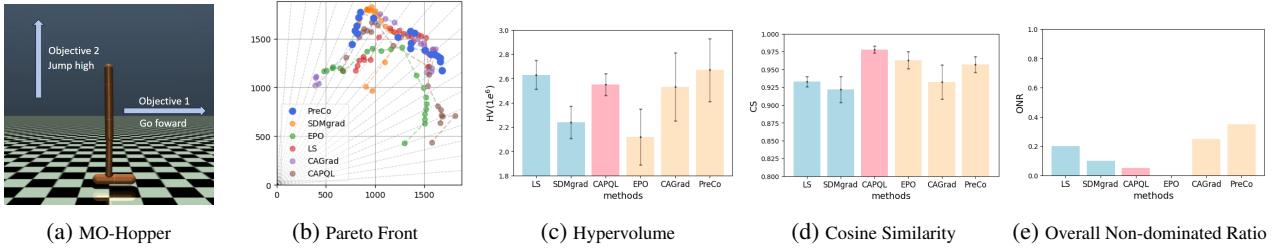


Figure 7: Optimality (HV) and Controllability (CS) of MO-Hopper.

relatively much better performance for MO-Hopper than Fruit-Tree and MO-Reacher. This is because the Pareto front of MO-Hopper is a simple 2-D front that is already strictly convex, adding a concave term will not cause an unacceptable information loss of the original problem. Our method demonstrated superior performance in HV, while its CS was only slightly lower than that of CAPQL. CAPQL adds a concave term to the reward to make the augmented Pareto more strictly

convex, resulting in better CS but the augmentation caused information loss of the original problem, thus causing suboptimal HV and ONR. Unlike EPO, which employs hard constraints on similarity, our proposed PreCo utilizes soft constraints. This could be the reason why PreCo can sacrifice a small degree of controllability for a significant enhancement in optimality. The hopper has to be able to jump before jumping forward, this is why objective 2 is higher than objective 1 for most methods. As a result, the asymmetric objectives make the discovered Pareto front not as symmetric as that of MO-Ant. We have a calibration approach to further improve controllability in Appendix H.

C.1.2. MO-ANT

The MO-Ant is a higher dimensional continuous control environment challenging for RL but not necessarily hard for MORL. The reward is 2-dimensional, with one for x-velocity and one for y-velocity. Although the robotic agent has more complex dynamics, the objectives appear to be very similar since both involve the movement of the agent, making it relatively easy for preference control. As shown in Fig. 8, the Pareto front has an intuitive convex shape. The similarity approaches such as PreCo and EPO have high CS metrics over 0.98. This indicates that our proposed PCRL scheme with similarity optimization is scalable for higher dimensional environments and has better preference controllability than LS methods such as LS and SDMgrad. They also have much better ONRs.

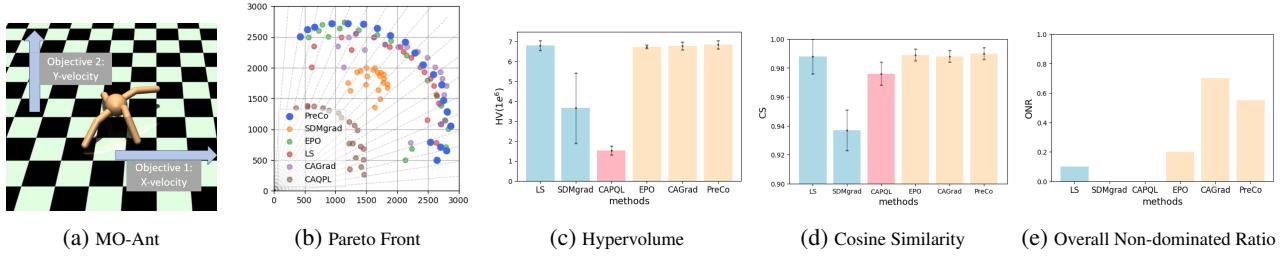


Figure 8: Optimality (HV) and Controllability (CS) of MO-Ant. PreCo (ours) achieves the best CS. Though being the second best on HV, it achieves the widest spread on the Pareto front in (d). In (b)-(c), methods of blue bars are based on linear scalarization, while methods of orange bars optimize the similarity.

The relative performance of CAPQL drops from lower dimensional Hopper to higher dimensional Ant. The potential reason could be that it requires learning an actor outputting not only the means but also the log standard deviations of the action distributions, which is equivalent to learning with an action space of double the original size. As a result, it could be less sample efficient than other methods that do not need to learn variance or standard deviation, especially for the high dimensional Ant environment. In contrast, LS has competitive HV and CS in this environment. One reason is that the symmetric and strictly convex Pareto front does not expose the limitations of LS, and another reason could be that its simplicity does not complicate the learning of lower-level robotic control.

C.2. Compatibility with lower-level modifications

Mainstream MORL methods have focused on optimizing the LS objective, often incorporating different low-level modifications to improve sample efficiency. As summarized in Table 5, Q-envelope (Yang et al., 2019) uses HER (Andrychowicz et al., 2017), Q-pensive (Hung et al., 2023) combines Q-buffer and HER.

Table 5: Summary of methods and their low-level modifications

Method	Modifications
Q-envelope	HER
Q-pensive	Q-buffer, HER

Intuitively, HER allows reusing experiences from one preference for others, while Q-buffer maintains Q-value snapshots that optimistically estimate values to leverage this stored information. Since Q-pensive essentially combines Q-envelope with Q-buffer, we compare the performance of our PreCo with HER and Q-buffer against Q-pensive. The results are shown in Table 6.

Table 6: Comparison of Q-Pensieve and PreCo (ours) on Fruit-Tree and MO-Reacher environments.

Environment	Metric	Q-Pensieve	PreCo (Ours)
Fruit-Tree	HV	6.86 ± 0.01	15.61 ± 0.75
	CS	0.72 ± 0.01	0.78 ± 0.03
MO-Reacher	HV	47.36 ± 0.73	49.55 ± 0.74
	CS	0.89 ± 0.02	0.91 ± 0.02

Our PreCo with HER and Q-buffer consistently outperforms Q-Pensieve, the state-of-the-art LS method, across environments and metrics, achieving higher HV and CS scores.

In non-strictly convex environments such as Fruit-Tree, the advantage is particularly pronounced, as Q-Pensieve’s reliance on the LS objective limits its ability to discover a broad set of optimal policies. In contrast, in strictly convex settings like Reacher, the performance gap narrows, as expected.

These results further underscore PreCo’s superiority over LS-based methods and demonstrate its compatibility with low-level sample efficiency enhancements.

D. Experimental details

The test preferences are $p \in \mathcal{P}$ with a resolution of 0.1 for each dimension.

For instance, in 3-D cases, these preferences include

$$[0, 0, 1], [0, 0.1, 0.9], \dots, [0, 1, 0], [0.1, 0, 0.9], \dots, [0.9, 0.1, 0], [1, 0, 0],$$

with a quantity of 66. There are 286 test preferences for 4-D, 1001 for 5-D, and 3003 for 6-D.

During training, the preferences were sampled uniformly from the convex coefficient set \mathcal{P} , making the probability of sampling an exact test preference nearly zero. Therefore, high CS metric in test time means the ability to generalize to unseen preferences.

We run 5 seeds for each environment setting, and for each run, we select the best-performing agent as a candidate for testing. The results are presented as the mean and standard deviation of the 5 candidates.

D.1. MO-Ant

The exact data for bar charts in Fig. 8 is shown in Table 7.

	LS	SDMgrad	CAPQL	EPO	CAGrad	PreCO
HV(*1e ⁶)	6.81 ± 0.24	3.67 ± 1.76	1.53 ± 0.22	6.75 ± 0.08	6.79 ± 0.20	6.85 ± 0.21
CS	0.988 ± 0.012	0.937 ± 0.014	0.976 ± 0.008	0.989 ± 0.004	0.988 ± 0.0004	0.990 ± 0.004

 Table 7: HV and CS performance in MO-Ant environment, the HV value has a unit of 1e⁶.

For MO-Ant, both SDMgrad and PreCo have λ that increase linearly with each update from 1 to 5. EPO has a constraint threshold of $\epsilon = 3e - 4$ for cosine similarity Ψ' , which is very small, making it comparable to Tchebycheff Scalarization and also similar to PreCo with a large constant $\lambda_t = \lambda >> 1$. CAGrad has constraint $c = 0.2$. The definitions of c and ϵ can be found in Appendix B.

D.2. MO-Hopper

The exact data for bar charts in Fig. 7 is shown in Table 9.

For MO-Hopper, both SDMgrad and PreCo have λ increasing linearly with every update from 3 to 11. EPO has a constraint threshold of $\epsilon = 3e - 4$ for cosine similarity Ψ' . CAGrad has constraint $c = 0.1$.

Table 8: Hyper-parameters settings MO-Ant.

Hyper-parameter	Value
Discount (γ)	0.99
Optimizer	Adam (Kingma & Ba, 2015)
Learning rate for networks	3×10^{-4}
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Activation function	ReLU
Batch size	256
Buffer Size	1×10^6
Starting timesteps	2.5×10^3
Gradient clipping	False
Exploration method	Noise
Noise distribution	$\mathcal{N}(0, 0.1^2)$
Noise clipping limit	0.5
Policy frequency (delay)	2
Target network update rate (τ)	5×10^{-3}
Maximum episode timesteps	500
Preference sampling	every new episode until max total steps is reached
Evaluation episodes for each test preference	10

	LS	SDMgrad	CAPQL	EPO	CAGrad	PreCO
HV(* $1e^6$)	2.63 ± 0.12	2.24 ± 0.13	2.55 ± 0.09	2.53 ± 0.28	0.94	2.67 ± 0.26
CS	0.933 ± 0.007	0.922 ± 0.018	0.978 ± 0.005	0.963 ± 0.023	0.932 ± 0.024	0.957 ± 0.011

Table 9: HV and CS performance in MO-hopper environment, the HV value has a unit of $1e^6$.

D.3. MO-Reacher

The exact data for Fig. 5 is shown in Table 11.

For MO-Reacher, both SDMgrad and PreCo have λ_t increasing linearly with every update from 10 to 20. EPO has a constraint threshold of $\epsilon = 3e - 4$ for cosine similarity Ψ' , which is very small, making it comparable to Tchebycheff Scalarization and also similar to PreCo with a large constant $\lambda_t = \lambda \gg 1$. CAGrad employs a constraint constant of $c = 0.1$.

D.4. Fruit Tree

For 3-D reward, most runs of LS only learn a very limited number of values like [4.71, 5.39, 5.40] and the values SDMGrad for most test preferences lie at [4.01, 7.17, 1.47]. They, as the LS approach, discover much fewer Pareto optimal policies than methods of the similarity approach, which have one value for each test preference. This shows the limitation of linear scalarization methods. In theory, LS methods have the potential to discover all Pareto optimal policies for MORL (Lu et al., 2023). However, in practice, this is often not the case. Possible reasons could be the numerical instability inherent in deep RL, limitations of model capacity, and the fact that the value space is usually **not strictly convex**.

E. Practical implementation of the update procedure

E.1. Algorithm framework

Our goal is to train a single agent that can be conditioned on different performance preferences and 0-shot adapt to user preference at test time. During training, we need to uniformly sample preferences from \mathcal{P} and let the agent learn to find Pareto optimal policies with values aligned to p . The procedure is shown in Algorithm 2.

Table 10: Hyper-parameters settings MO-Hopper.

Hyper-parameter	Value
Discount (γ)	0.99
Optimizer	Adam (Kingma & Ba, 2015)
Learning rate for networks	3×10^{-4}
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Activation function	ReLU
Batch size	256
Buffer Size	1×10^6
Starting timesteps	2.5×10^3
Gradient clipping	False
Exploration method	Noise
Noise distribution	$\mathcal{N}(0, 0.1^2)$
Noise clipping limit	0.5
Policy frequency (delay)	2
Target network update rate (τ)	5×10^{-3}
Maximum episode timesteps	500
Preference sampling	every new episode until max total steps is reached
Evaluation episodes for each test preference	10

Algorithm 2 PCRL with PreCo update

- 1: **Initialize:**
 \mathcal{B} : Buffer.
 N : Number of training samples for p ,
 E : Number of training episodes for every p sample,
 π_θ : Preference-conditioned actor model,
 \mathbf{Q}_ϕ for DDPG/TD3 or \mathbf{V}_ϕ for A3C/PPO: Preference-conditioned critic model with m -dimensional output, where m is the objective number.
 - 2: **for** $n = 0, 1, \dots, N - 1$ **do**
 - 3: Sample preference $p \in \mathcal{P}$
 - 4: **for** $e = 0, 1, \dots, E - 1$ **do**
 - 5: Rollout with policy $\pi_\theta(\cdot | \cdot, p)$
 - 6: Store transitions (s, a, r, p) in \mathcal{B}
 - 7: **end for**
 - 8: Update \mathbf{Q}_ϕ or \mathbf{V}_ϕ by minimizing TD error for every objective.
 - 9: Estimate policy-level gradient $\nabla_{\pi_p}^T \hat{v}^{\pi_p}$ by Eq.24/26 for TD3 or Eq.30/37 for PPO.
 - 10: Estimate similarity gradient $\nabla_{\pi_p} \Psi(p, \hat{v}^{\pi_p}) = \nabla_{\pi_p}^T \hat{v}^{\pi_p} \nabla_v \Psi(p, \hat{v}^{\pi_p})$
 - 11: Get policy-level update direction d^* by solving Eq.6 with $\nabla_{\pi_p}^T \hat{v}^{\pi_p}$ and $\nabla_{\pi_p} \Psi(p, \hat{v}^{\pi_p})$
 - 12: Update θ by solving Eq.25 for TD3 or Eq.32/34 for PPO with d^*
 - 13: **end for**
-

	random	LS	SDMgrad	CAPQL	PDMORL	EPO	CAGrad	PreCO
HV($\times 10^8$)	13.47	15.66 ± 8.03	20.37 ± 0.37	20.56 ± 0.47	22.69 ± 0.15	9.87 ± 3.11	13.46 ± 9.71	33.11 ± 6.29
CS	0.758	0.652 ± 0.030	0.761 ± 0.001	0.760 ± 0.001	0.754 ± 0.002	0.845 ± 0.078	0.760 ± 0.002	0.906 ± 0.002

 Table 11: HV and CS performance in MO-reacher environment, the HV value has a unit of 10^6 .

Table 12: Hyper-parameters settings MO-Reacher.

Hyper-parameter	Value
Discount (γ)	0.99
Optimizer	Adam (Kingma & Ba, 2015)
Learning rate for networks	3×10^{-4}
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Activation function	ReLU
Batch size	250
Gradient clipping	False
Exploration method	Policy Entropy
Entropy Coefficient	0.001
epsilon-clip for PPO	0.001
Epochs per PPO update	3
Timesteps every update	100
Maximum episode timesteps	250
Number of episodes per preference sample	40
Number of preference samples (for 4D reward)	600
Evaluation episode for each test preference	10

E.2. Policy-level gradient

Solving the min-norm problem (6) with parameter-level gradients $\nabla_{\theta} \hat{v}^{\pi_p}$ at every gradient update can be memory and computationally expensive when $|\theta|$ is large. Video game playing agents like AlphaZero (Silver et al., 2018) and AlphaStar (Vinyals et al., 2019) can have millions of model parameters. Besides, Large models with billions of model parameters have become very common with recent developments in Large Language Models (LLMs) (Zhao et al., 2023; Minaee et al., 2024). To circumvent this issue, we suggest solving the min-norm problem (6) before the gradient propagates to the model parameter θ . Therefore, ideally, we want to solve the min-norm problem with gradients at the policy-level $\nabla_{\pi_p} \hat{v}^{\pi_p}$, which has only a size of batch size B of hundreds at each update; In practice for deep RL implementations, as shown by Fig. 9, $\nabla_{\pi_p} \hat{v}^{\pi_p}$ can also be replaced by the gradients of the value \hat{v}^{π_p} with respect to the policy model outputs, such as $\nabla_{l_p} \hat{v}^{\pi_p}$ for logits of categorical distribution policies and $\nabla_{\mu_p, \sigma_p} \hat{v}^{\pi_p}$ for means and standard deviations of diagonal Gaussian distribution policies.

E.2.1. CONTINUOUS ACTION SPACE

Value-based methods like TD3 (Fujimoto et al., 2018) and SAC (Haarnoja et al., 2018) are often used for continuous action spaces. To avoid computing min-norm with parameter gradient

$$\nabla_{\theta} \hat{v}^{\pi_p} = \mathbb{E}[\nabla_{\theta} \hat{Q}(s, a, p) |_{a \sim \pi_{\theta}(s, p)}] \quad (23)$$

We look at their policy formulations. Their policy $\pi(a|s, p)$ is often a Gaussian or squashed Gaussian distribution with parameters mean $\mu_p(s)$ and log standard deviation $\log \sigma_p(s)$. We denote a distribution parameter vector ρ_p with $\rho_p(s) = [\mu_p(s), \log \sigma_p(s)]^T$ and we can get

$$\nabla_{\rho_p} \hat{v}^{\pi_p} = \mathbb{E}[\nabla_{\rho_p} \hat{Q}(s, a, p)] \quad (24)$$

For each update, the size of each objective gradient $\nabla_{\rho_p} \hat{v}_i^{\pi_p}$ is $2|\mathcal{A}| \times B$, where B is the batch size and ρ has a size of

Table 13: Hyper-parameters settings Fruit-tree.

Hyper-parameter	Value
Discount (γ)	0.99
Optimizer	Adam (Kingma & Ba, 2015)
Learning rate for networks	3×10^{-4}
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Activation function	ReLU
Batch size	100
Gradient clipping	False
Exploration method	Policy Entropy
Entropy Coefficient	0.001
epsilon-clip for PPO	0.001
Epochs per PPO update	3
Timesteps every update	100
Maximum episode timesteps	100
Number of episodes per preference sample	20
Number of preference samples (for 4D reward)	3000
Evaluation episode for each test preference	10

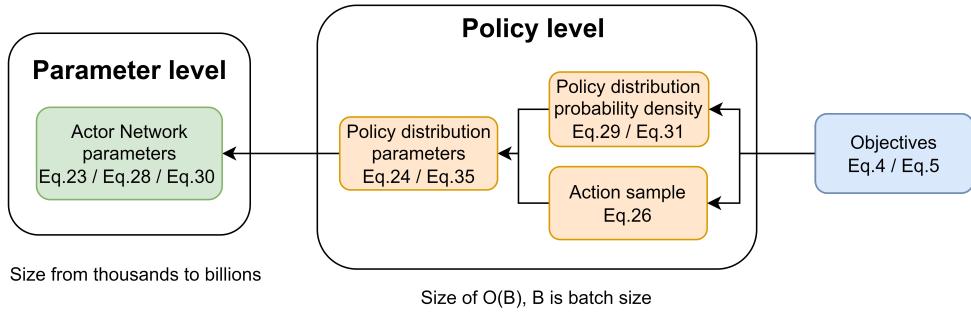


Figure 9: Backward path of policy update. we can see that the gradient from objectives to the parameters of the actor network first backpropagate through the probability density of action distribution (for policy-based methods such as A3C/PPO) or action sample (for value-based methods such as DDPG/TD3), then propagate through the distribution parameters of policies such as the logits for categorical distribution or μ, Σ for Gaussian distributions. We consider these the policy-level gradients. They often have a size of $\mathcal{O}(B)$, where B is the batch size. Since B is often limited to a few hundreds, the size of a policy-level gradient would be much smaller than the size of the neural network parameter.

$2|\mathcal{A}|$. This means that $\nabla_{\rho_p} \hat{v}_i(s, p)$ could have a much lower dimension than $\nabla_\theta \hat{v}_i(s, p)$, thus increasing the memory and computational efficiency.

After getting the update direction d for ρ_p by solving the min norm problem (6) with $\nabla_{\rho_p} \hat{v}^{\pi_p}$, we update model parameter θ by solving

$$\max_{\theta} \left\{ d^\top \nabla_{\theta} \rho_p \text{ s.t. } \|\rho - \rho_{\text{old}}\|_2 < \delta \right\},$$

which is a trust region formulation that updates ρ in the direction of d while keeping in a local region where d is valid. A simple and practical implementation for parameter update can be as follows:

$$\max_{\theta} \mathcal{J}(\theta) = \mathbb{E}_{s,a} [\text{clip}(\rho_p(s, a), \rho_p(s, a) - \epsilon, \rho_p(s, a) + \epsilon) d(s, a)]. \quad (25)$$

The update of every entry of π is clipped to ϵ , so $\|\pi_\theta - \pi_\theta\|_2 \leq \sqrt{B * \epsilon^2} = \delta$, where B is the batch size.

For more expressive models such as diffusion models (Wang et al., 2023) or normalizing flows (Brahmanage et al., 2023), the mean and covariance gradients would not be adequate. We can instead use the gradient of action samples as policy-level gradients and we get an update direction

$$d(s, a) = \nabla_a \hat{Q}(s, a, p) \quad (26)$$

for every (s, a) sample in the batch. Then we can perform min-norm with d and update the more expressive policy networks by reparameterization techniques.

E.2.2. DISCRETE ACTION SPACE

Policy-based methods like A3C (Mnih et al., 2016) and PPO (Schulman et al., 2017) are often used for discrete action spaces. We can approximate the multi-objective value function $\hat{v}^{\pi_p}(s)$ by a function $\hat{v}(s, p)$ that takes s and p as inputs, sample the episodic returns as vector \mathbf{R} , and calculate the multi-objective advantage function as

$$\hat{\mathbf{A}}(s, a) = \mathbf{R} - \hat{v}(s, p) \quad (27)$$

Then, the policy gradient in the model parameter space is

$$\nabla_{\theta} \hat{v}^{\pi_p} = \mathbb{E} \left[\frac{\nabla_{\theta} \pi(a|s, p)}{\pi(a|s, p)} \hat{\mathbf{A}}(s, a, p) \right], \quad (28)$$

And for gradient at the policy space $d = \nabla_{\pi_p} \hat{v}(s, p)$, we have

$$d(s, a) = \frac{1}{\pi(a|s, p)} \hat{\mathbf{A}}(s, a, p) \quad (29)$$

for every (s, a) sample. When using policy optimization methods like PPO/TRPO they are

$$\nabla_{\theta} \hat{v}^{\pi_p} = \mathbb{E} \left[\frac{\nabla_{\theta} \pi(a|s, p)}{\pi_{\text{old}}(a|s, p)} \hat{\mathbf{A}}(s, a, p) \right], \quad (30)$$

$$d(s, a) = \frac{1}{\pi_{\text{old}}(a|s, p)} \hat{\mathbf{A}}(s, a, p) \quad (31)$$

In practical situations, at every update, the size of each objective gradient $\nabla_{\pi_p} \hat{v}_i(s, p)$ is the batch size B , and the min norm problem (6) can be performed with gradients of batch size B , which could be much smaller than the parameter size of deep neural networks, especially when implemented for large language models.

After getting the update direction $d = \nabla_{\pi_p} \hat{v}(s, p)$ for π_p , we optimize model parameters by

$$\max_{\theta} \left\{ d^\top \nabla_{\theta} \pi_p \text{ s.t. } \|\pi_p - \pi_{p,\text{old}}\|_2 < \delta \right\}, \quad (32)$$

which is a trust region formulation that updates π_p in the direction of d while keeping in a local region where d is valid. This can be practically implemented by an objective as follows:

$$\mathcal{J}(\theta) = \mathbb{E}_{s,a} [\text{clip}(\pi_\theta(a|s, p), \pi_{\theta_{\text{old}}}(a|s, p) - \epsilon, \pi_{\theta_{\text{old}}}(a|s, p) + \epsilon) d(s, a)]. \quad (33)$$

The update of every entry of π is clipped to ϵ , so $\|\pi_\theta - \pi_{\theta,\text{old}}\|_2 \leq \sqrt{B * \epsilon^2} = \delta$, where B is the batch size.

Trust region formulation with KL-divergence could be more suitable for categorical distribution π_p , so another formulation of parameter update could be

$$\max_{\theta} \left\{ d^\top \nabla_{\theta} \pi_p \quad s.t. \quad \pi_p \pi_{p,\text{old}} < \delta \right\}. \quad (34)$$

Whether a KL divergence trust region is theoretically compatible with the solution d of the min norm problem (6) will be further researched in our future work.

One potential issue of $\nabla_{\pi_p} \hat{v}(s, p)$ is that $(\pi_p + \alpha \nabla_{\pi_p} \hat{v}(s, p))$ may not be in the probability simplex. As a result, projecting it back onto the probability simplex could cause it to deviate from the intended update direction. Since policies for discrete action spaces are often categorical distributions, one way to avoid this issue is to consider the gradient of l_p , which denotes the logits for policy $\pi(\cdot | \cdot, p)$ conditioned on preference p , and $l_p(s)$ are the logits for $\pi(\cdot | s, p)$. The logits do not have the constraint to be in the probability simplex.

$$\begin{aligned} & \nabla_{l_p} \pi_p(s) \\ &= -\pi(a|s, p) [\pi(a_1|s, p), \pi(a_2|s, p), \dots, \pi(a|s, p) - 1, \dots, \pi(a_{|\mathcal{A}|}|s, p)]^T, \end{aligned} \quad (35)$$

where $\nabla_{l_p} \pi_p(s)$ is the s -th entry of the jacobian $\nabla_{l_p} \pi_p$, and

$$[\pi(a_1|s, p), \pi(a_2|s, p), \dots, \pi(a|s, p) - 1, \dots]^T$$

is a vector of action space size $|\mathcal{A}|$.

Then, we can get

$$\begin{aligned} & \nabla_{l_p} \hat{v}^{\pi_p}(s) \\ &= \mathbb{E} \left[-\frac{\pi(a|s, p)}{\pi_{\text{old}}(a|s, p)} \hat{\mathbf{A}}(s, a, p) [\pi(a_1|s, p), \pi(a_2|s, p), \dots, \pi(a|s, p) - 1, \dots, \pi(a_{|\mathcal{A}|}|s, p)] \right], \end{aligned} \quad (36)$$

where $\nabla_{l_p} \hat{v}^{\pi_p}(s)$, of size $m \times |\mathcal{A}|$, is the index $[:, s, :]$ for $\nabla_{l_p} \hat{v}^{\pi_p}$ tensor, of size $m \times |\mathcal{S}| \times |\mathcal{A}|$.

In every update, the size of the objective gradient $\nabla_{l_p} \hat{v}^{\pi_p}(i, s)$ for the objective i has a size of $B \times |\mathcal{A}|$. For large language models, the action space could be the vocabulary size of tens of thousands, so $B \times |\mathcal{A}|$ could be in millions, but is still much smaller than the parameter size that is often in billions. Moreover, for large action spaces, $\pi(a|s, p) \cdot \pi(a'|s, p)$ could be much smaller than $\pi(a|s, p)$, so $\nabla_{l_p} \hat{v}^{\pi_p}(s, a)$ can be approximated by

$$\nabla_{l_p} \hat{v}^{\pi_p}(s, a) \approx \frac{\pi(a|s, p)}{\pi(a|s, p)} \hat{\mathbf{A}}(s, a, p) = \hat{\mathbf{A}}(s, a, p). \quad (37)$$

This approximation of $\nabla_{l_p} \hat{v}^{\pi_p}$ is what we implemented to replace $\nabla_{\pi_p} \hat{v}^{\pi_p}$ in the min norm problem (6) to avoid solving min norm with large parameter gradient $\nabla_{\theta} \hat{v}^{\pi_p}$. The model parameters are updated by solving (32) using $\nabla_{l_p} \hat{v}^{\pi_p}$ as d .

F. More details about the similarity objective

F.1. About proposed similarity function Ψ

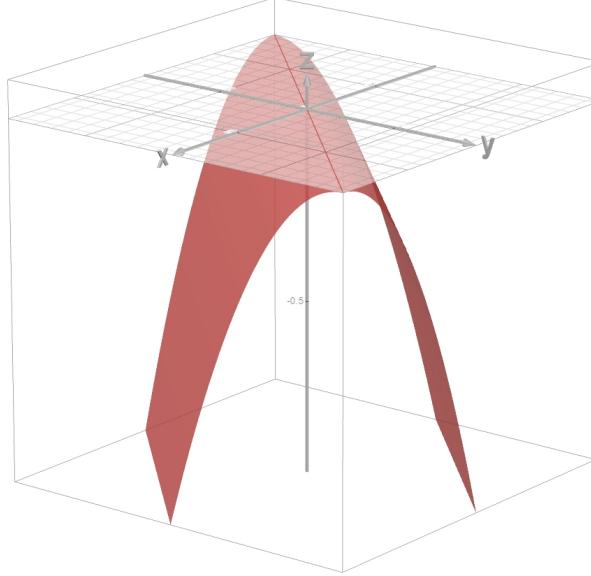
For two objective cases, when $p = [0.5, 0.5]$, $\Psi(p, v)$ is shown in Fig. 10. The x-axis is the first element of v and y-axis is the second element of v . The z-axis is the value of $\Psi(p, v)$.

We can see that the similarity is maximized to 0 only when $\frac{v_0}{p_0} = \frac{v_1}{p_1}$. It is also smooth as proved by Lemma 4.1.

F.2. Similarity objective design for better theoretical properties

Theorem 4.2 requires the similarity gradient to be both Lipschitz continuous and convex combinations of the objective gradients. Which formally means that for an similarity objective $\Psi'(p, \cdot)$,

$$\nabla_v \Psi'(p, v) \in \mathcal{W},$$


 Figure 10: $\Psi(\mathbf{p}, \cdot)$ when $\mathbf{p} = [0.5, 0.5]^T$

which can not be satisfied by $\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}) = \max_i \frac{\mathbf{v}_i}{p_i} \mathbf{p} - \mathbf{v}$, because $\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})$ will be $\mathbf{0}$ when \mathbf{p} and \mathbf{v} are perfectly aligned. Moreover, we can not directly normalize $\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})$ by dividing $\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})\|_1$, because this will make the normalized gradient not Lipschitz continuous (the gradient changes drastically when \mathbf{v} passes the direction of \mathbf{p}).

Our hints to design such a similarity function $\Psi'(\mathbf{p}, \cdot)$ are as follows:

- Its similarity gradient $\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v})$ could get close to \mathbf{p} , when \mathbf{v} has a high similarity to the preference \mathbf{p} ;
- $\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v})$ should be close to the normalized gradient $\frac{\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})}{\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})\|_1}$ when the similarity is low.

One possible design is:

$$\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v}) = \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \mathbf{p} + (1 - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}}) \frac{\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})}{\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})\|_1} \quad (38)$$

where CS is the cosine similarity. In this design, $\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v})$ is always a convex combination between \mathbf{p} and $\frac{\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})}{\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})\|_1}$, so itself is always a convex coefficient.

We have for $\mathbf{v}, \mathbf{v}' \neq \mathbf{0}$:

$$\begin{aligned} & \frac{\|\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v}) - \nabla_{\mathbf{v}'} \Psi'(\mathbf{p}, \mathbf{v}')\|}{\|\mathbf{v}' - \mathbf{v}\|} \\ &= \left\| \left(\exp^{\frac{1-CS(\mathbf{p}, \mathbf{v}')}{\eta}} - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \right) \mathbf{p} + \left(1 - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v}')}{\eta}} \right) \frac{\nabla_{\mathbf{v}'} \Psi(\mathbf{p}, \mathbf{v}')}{\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}')\|_1} - \left(1 - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \right) \frac{\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})}{\|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v})\|_1} \right\| \\ &\leq \left\| \left(\exp^{\frac{1-CS(\mathbf{p}, \mathbf{v}')}{\eta}} - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \right) \mathbf{p} \right\| + \left\| \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v}')}{\eta}} - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \right\| \\ &\leq (1 + \|\mathbf{p}\|) \left\| \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v}')}{\eta}} - \exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}} \right\| \end{aligned} \quad (39)$$

which is equivalent to proving $\exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}}$ is Lipschitz continuous.

Let $p \neq 0$ be fixed. Consider:

$$CS(\mathbf{v}, \mathbf{p}) = \frac{\langle \mathbf{v}, \mathbf{p} \rangle}{\|\mathbf{v}\| \|\mathbf{p}\|}, \quad \mathbf{v} \neq 0.$$

Compute the gradient:

$$\nabla_{\mathbf{v}} CS(\mathbf{v}, \mathbf{p}) = \frac{1}{\|\mathbf{p}\| \|\mathbf{v}\|^2} \left(\|\mathbf{v}\| \mathbf{p} - \frac{\langle \mathbf{v}, \mathbf{p} \rangle}{\|\mathbf{v}\|} \mathbf{v} \right).$$

Bounding its norm:

$$\|\nabla_{\mathbf{v}} CS(\mathbf{v}, \mathbf{p})\| \leq \frac{1}{\|\mathbf{p}\| \|\mathbf{v}\|^2} (\|\mathbf{v}\| \|\mathbf{p}\| + \|\mathbf{v}\| \|\mathbf{p}\|) = \frac{2}{\|\mathbf{v}\|}. \quad (40)$$

Thus, for $\|\mathbf{v}\| \geq \delta > 0$, $\|\nabla_{\mathbf{v}} CS(\mathbf{v}, \mathbf{p})\| \leq \frac{2}{\delta}$, showing CS is Lipschitz continuous with constant $\frac{2}{\delta}$.

Then for $\exp^{\frac{1-CS(\mathbf{p}, \mathbf{v})}{\eta}}$, let

$$f(\mathbf{v}) = \exp\left(\frac{1-CS(\mathbf{v}, \mathbf{p})}{\eta}\right)$$

and we have

$$\nabla_{\mathbf{v}} f = -\frac{1}{\eta} \exp\left(\frac{1-CS(\mathbf{v}, \mathbf{p})}{\eta}\right) \nabla_{\mathbf{v}} CS(\mathbf{v}, \mathbf{p}).$$

By Equation (40), we get

$$\|\nabla_{\mathbf{v}} f\| \leq \frac{2}{\eta \|\mathbf{v}\|} \exp\left(\frac{2}{\eta}\right) \quad (41)$$

Then for $\|\mathbf{v}\| \geq \delta$:

$$\|\nabla_{\mathbf{v}} f\| \leq \frac{2 \exp(2/\eta)}{\eta \delta}.$$

Hence, $f(\mathbf{v})$ is Lipschitz continuous on $\|\mathbf{v}\| \geq \delta$ with constant $\frac{2 \exp(2/\eta)}{\eta \delta}$.

Therefore, $\nabla_{\mathbf{v}} \Psi'(\mathbf{p}, \mathbf{v})$ defined in Equation (38) is both Lipschitz continuous and is a convex coefficient when $\|\mathbf{v}\| \geq \delta$, which is very common in practice with non-zero values.

G. MOO Toy Example

This is an toy example used in SDMGrad (Xiao et al., 2023) to show that in MOO, our proposed PreCo can achieve better or comparable performance under stochastic settings. Besides, PreCo can find the Pareto optimal point optimizing the similarity function $\Psi(\mathbf{p}, \cdot)$.

The two objectives $L_1(x)$ and $L_2(x)$ shown in Fig. 11 are defined on $x = (x_1, x_2)^T \in \mathbb{R}^2$,

$$L_1(x) = f_1(x)g_1(x) + f_2(x)h_1(x) \text{ and } L_2(x) = f_1(x)g_2(x) + f_2(x)h_2(x),$$

where the functions are given by

$$\begin{aligned} f_1(x) &= \max(\tanh(0.5x_2), 0) \\ f_2(x) &= \max(\tanh(-0.5x_2), 0) \\ g_1(x) &= \log\left(\max(|0.5(-x_1 - 7) - \tanh(-x_2)|, 0.000005)\right) + 6 \\ g_2(x) &= \log\left(\max(|0.5(-x_1 + 3) - \tanh(-x_2) + 2|, 0.000005)\right) + 6 \\ h_1(x) &= ((-x_1 + 7)^2 + 0.1(-x_1 - 8)^2)/10 - 20 \\ h_2(x) &= ((-x_1 - 7)^2 + 0.1(-x_1 - 8)^2)/10 - 20. \end{aligned}$$

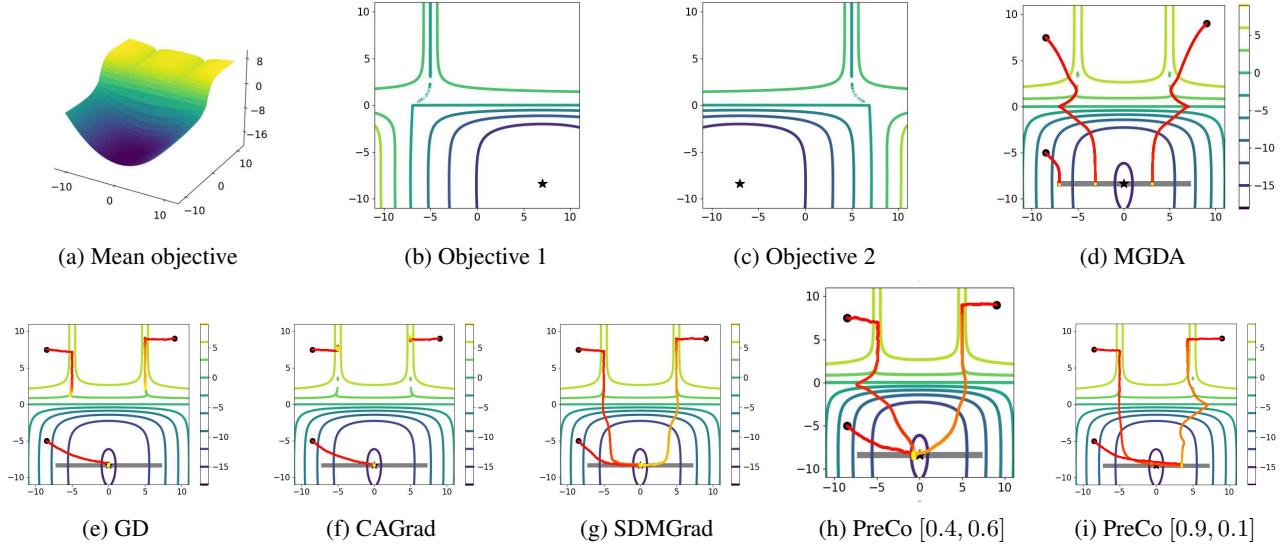


Figure 11: A two-objective toy example.

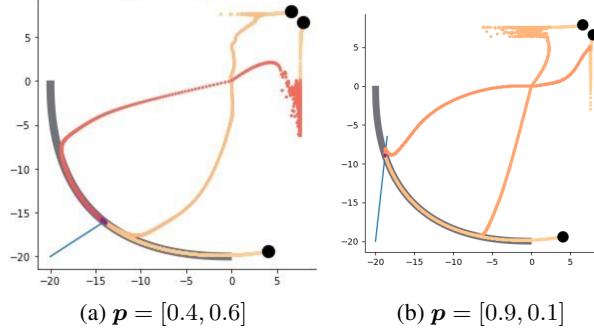


Figure 12: Plots showing the Pareto front: The x-axis is L_1 and the y-axis is L_2 . The blue line is the direction of the preference, for (a) it is $\mathbf{p} = [0.4, 0.6]$, and for (b) it is $\mathbf{p} = [0.9, 0.1]$. All three initial points converged to the Pareto optimal point that intersects with the line of \mathbf{p} direction.

Initializations points are from $\{(-8.5, 7.5), (-8.5, 5), (9, 9)\}$. The optimization trajectories are visualized in Fig. 11. The starting point of every trajectory in Fig. 11d-Fig. 11g is given by the \bullet symbol, and the color of every trajectory changes gradually from red to yellow. The gray horizontal line illustrates the Pareto front, and the \star symbol denotes the global optimum for the mean objective $L_0 = 0.5L_1 + 0.5L_2$. The setting is the same as in (Xiao et al., 2023) and all other methods except PreCo optimize for L_0 . Zero-mean Gaussian noise is added to the gradient of each objective for all the methods except MGDA. Adam optimizer is adopted with learning rate of 0.002 and 70000 iterations for each run. We can see that GD and CAGrad can fail to converge to the Pareto front in certain circumstances. Only SDMGrad and our proposed PCGrad converge to the Pareto front in all cases. Notice that, the preference PreCo in Fig. 11h is $\mathbf{p} = [0.4, 0.6]$, as shown in Fig. 12, we can see that it converges to a point optimizing $\Psi(\mathbf{p}, [L_1, L_2]^T)$. In addition, Figs. 11i and 12b show for the case where $\mathbf{p} = [0.9, 0.1]$, PreCo also updates to the preference specific Pareto optimal point.

Below is a figure comparing PreCo with existing preference following MOO algorithms such as Tchebycheff scalarization (Lin et al., 2024)(TS) (Ehrhart, 2005) and smooth Tchebycheff scalarization

In Fig. 13, PreCo first converges to the Pareto front, then as λ goes up, it converges to the preference desired solution.

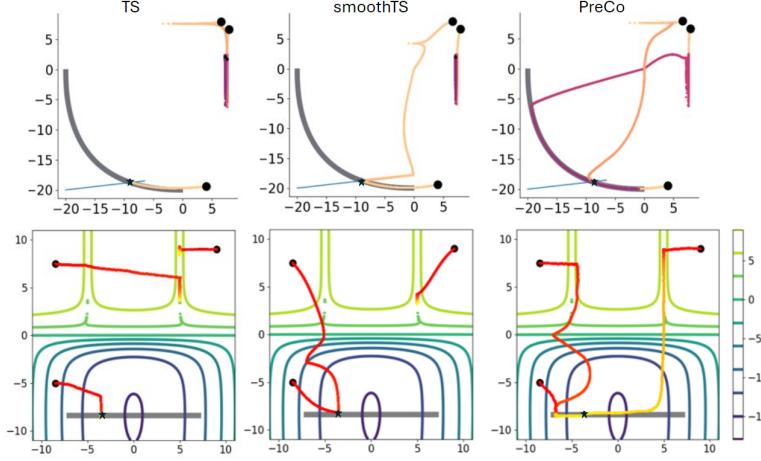


Figure 13: Plots showing the results of TS, SmoothTS and PreCo for preference $p = [0.1, 0.9]$. Our proposed PreCo converges to the preference-aligned Pareto front for all initialization.

H. Calibration

The reachable Pareto front for PCRL is often not the entire \mathbb{R}^m value space, and there are often gaps between the desired preferences and the values reached. To calibrate the possible misalignment between the input preference and the reached value in a sample-efficient way, we employ a Gaussian process (GP) based method to model the relationship between the input of the desired preference and the actual values reached by the agent.

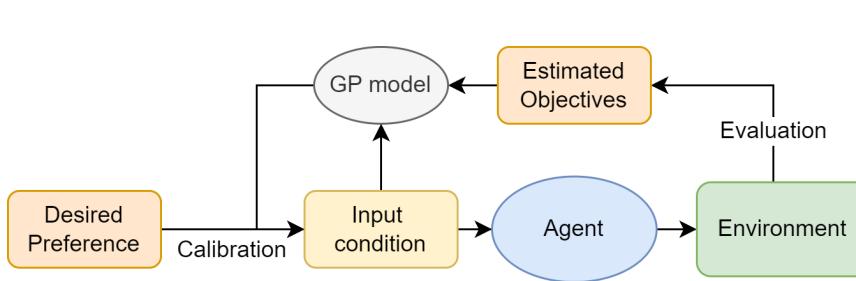


Figure 14: After training, due to general errors in deep learning or unreachable regions of the Pareto front, there could still be gaps between the actually reached objectives ratios and the desired preference. The calibration procedure obtains the reachable Pareto front and modifies the desired preference into an input that results in performance more aligned with the desired preference.

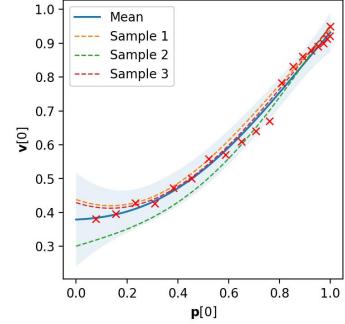


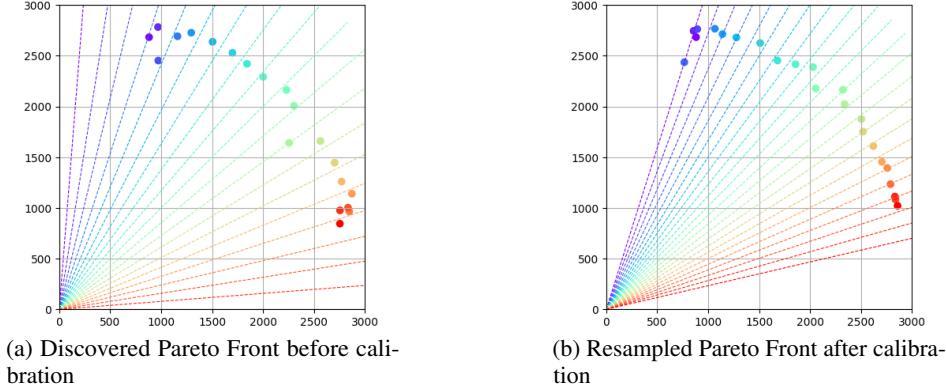
Figure 15: Example of GP regression for the (p, v^{π_p}) samples, this case has two objectives and shows the relation between the first element of the value and the first element of the preference. Each red cross point in this plot is a $(p_0, v_0^{\pi_p})$, where $p_0, v_0^{\pi_p}$ are the first elements of p, v^{π_p} .

After training, there might still be an input p' with better $\Psi(p, v^{\pi_p})$ than v^{π_p} . We want to find p' that solves $\max_{p'} E[\Psi(p, v^{\pi_{p'}})]$ for any p . we first uniformly sample the values v^{π_p} reached by giving the agent preference input from $\{p \in \mathbb{R}^m : p^T 1 = 1\}$, then perform a GP regression for the (p, v^{π_p}) samples. As shown in Fig. 15, some samples can provide a Gaussian distribution of the mapping $\phi(p)$ from p to v^{π_p} . Based on the distribution of ϕ , for a desired preference p , we can find a good input p' to solve

$$\max_{p'} \mathbb{E}[\Psi(p, \phi(p'))] \quad (42)$$

This procedure learns the reachable regions of the agent and calibrates the desired preference into the best input for reaching the preference. Also, it is general and can be applied to any preference control approach. Here is an empirical example for calibration:

Figure 16: Pareto front before and after calibration.



The colored rays are preferences \mathbf{p} . The points with the same color as the preference vector \mathbf{p} are value $\hat{\mathbf{v}}^{\pi_p}$ of preference conditioned policy π_p . The left plot is before calibration, \mathbf{p} is directly used as the input for π_p . The right plot is after calibration, we know which regions can be reached, so preferences \mathbf{p} are not all directions but reachable directions. Also, the input for π_p is \mathbf{p}' by solving (42), resulting in higher similarity and the CS metric improved from 0.991 to 0.997.

I. Theoretical Proofs

I.1. Proof for Lemma 4.1

Lemma 4.1. *The similarity function $\Psi(\mathbf{p}, \cdot)$ is $(1 + \max_i \frac{|\mathbf{p}|}{|\mathbf{p}_i|})$ -Lipschitz smooth and $g_s(\cdot)$ is Lipschitz continuous under Assumption 4.1 and Assumption 4.3.*

Proof. By definition:

$$\Psi(\mathbf{p}, \mathbf{v}) = -\frac{1}{2} \left\| \max_i \frac{\mathbf{v}_i}{\mathbf{p}_i} \mathbf{p} - \mathbf{v} \right\|^2 \quad (43)$$

and

$$\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}) = \max_i \frac{\mathbf{v}_i}{\mathbf{p}_i} \mathbf{p} - \mathbf{v} = d(\mathbf{v}, \mathbf{p}) \mathbf{p} - \mathbf{v} \quad (44)$$

. where $d(\mathbf{v}, \mathbf{p})$ denotes $\max_i \frac{\mathbf{v}_i}{\mathbf{p}_i}$, we have

$$\begin{aligned} \|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}) - \nabla_{\mathbf{v}'} \Psi(\mathbf{p}, \mathbf{v}')\| &= \|d(\mathbf{v}, \mathbf{p}) \mathbf{p} - \mathbf{v} - d(\mathbf{v}', \mathbf{p}) \mathbf{p} + \mathbf{v}'\| \\ &= \|d(\mathbf{v}, \mathbf{p}) \mathbf{p} - d(\mathbf{v}', \mathbf{p}) \mathbf{p} - (\mathbf{v} - \mathbf{v}')\| \\ &\leq |d(\mathbf{v}, \mathbf{p}) \mathbf{p} - d(\mathbf{v}', \mathbf{p})| \|\mathbf{p}\| + \|\mathbf{v} - \mathbf{v}'\| \end{aligned} \quad (45)$$

Without loss of generality, we first consider the case where $d(\mathbf{v}, \mathbf{p}) - d(\mathbf{v}', \mathbf{p}) \geq 0$. We denote $i_{\mathbf{v}} = \arg \max_j \frac{\mathbf{v}_j}{\mathbf{p}_j}$

$$\begin{aligned} \|d(\mathbf{v}, \mathbf{p}) \mathbf{p} - d(\mathbf{v}', \mathbf{p}) \mathbf{p}\| + \|\mathbf{v} - \mathbf{v}'\| &= (d(\mathbf{v}, \mathbf{p}) - d(\mathbf{v}', \mathbf{p})) \|\mathbf{p}\| + \|\mathbf{v} - \mathbf{v}'\| \\ &\leq \left(\frac{\mathbf{v}_{i_{\mathbf{v}}}}{|\mathbf{p}_{i_{\mathbf{v}}}|} - \frac{\mathbf{v}'_{i_{\mathbf{v}}}}{|\mathbf{p}_{i_{\mathbf{v}}}|} \right) \|\mathbf{p}\| + \|\mathbf{v} - \mathbf{v}'\| \\ &\leq \frac{\|\mathbf{p}\|}{|\mathbf{p}_{i_{\mathbf{v}}}|} \|\mathbf{v} - \mathbf{v}'\| + \|\mathbf{v} - \mathbf{v}'\| \\ &\leq \max_i \frac{\|\mathbf{p}\|}{|\mathbf{p}_i|} \|\mathbf{v} - \mathbf{v}'\| + \|\mathbf{v} - \mathbf{v}'\| \\ &\leq (1 + \max_i \frac{\|\mathbf{p}\|}{|\mathbf{p}_i|}) \|\mathbf{v} - \mathbf{v}'\| \end{aligned} \quad (46)$$

The first inequality is because i_v is optimal for v but not necessarily for v' . The case where $d(v, p) - d(v', p) < 0$ can be proved by the same procedure by denote $i'_v = \arg \max_j \frac{v'_i}{p_i}$, then

$$\begin{aligned}
 \|d(v, p)p - d(v', p)p\| + \|v - v'\| &= (d(v', p) - d(v, p))\|p\| + \|v - v'\| \\
 &\leq (\frac{v'_{i'_v}}{\|p_{i'_v}\|} - \frac{v_{i'_v}}{\|p_{i'_v}\|})\|p\| + \|v - v'\| \\
 &\leq \frac{\|p\|}{\|p_{i'_v}\|}\|v - v'\| + \|v - v'\| \\
 &\leq \max_i \frac{\|p\|}{\|p_i\|}\|v - v'\| + \|v - v'\| \\
 &\leq (1 + \max_i \frac{\|p\|}{\|p_i\|})\|v - v'\|.
 \end{aligned} \tag{47}$$

Therefore, we have proven

$$\|\nabla_v \Psi(p, v) - \nabla_{v'} \Psi(p, v')\| \leq (1 + \max_i \frac{\|p\|}{\|p_i\|})\|v - v'\| \tag{48}$$

and the similarity function for a preference p is $1 + \max_i \frac{\|p\|}{\|p_i\|}$ Lipschitz smooth.

Next, we prove that $g_s(\pi_p) = G(\pi_p)\nabla_v \Psi(p, v^{\pi_p})$ is Lipschitz continuous.

We have:

$$\begin{aligned}
 \|g_s(x) - g_s(y)\| &= \|G(x)\nabla_v \Psi(p, v^x) - G(y)\nabla_v \Psi(p, v^y)\| \\
 &= \|(G(x) - G(y))\nabla_v \Psi(p, v^x) + G(y)(\nabla_v \Psi(p, v^x) - \nabla_v \Psi(p, v^y))\| \\
 &\leq \sum_{i=1}^m \|g_i(x) - g_i(y)\|\|\nabla_v \Psi(p, v^x)\| + \|G(y)\|\|\nabla_v \Psi(p, v^x) - \nabla_v \Psi(p, v^y)\|
 \end{aligned} \tag{49}$$

where the inequality is by Cauchy-Schwartz. Since under [Assumption 4.1](#) and [Assumption 4.3](#), $\|G(y)\| \leq C_g$ and $\|g_i(x) - g_i(y)\| \leq l_{i,1}\|x - y\|$, and

$$\|\nabla_v \Psi(p, v^x) - \nabla_v \Psi(p, v^y)\| \leq (1 + \max_i \frac{\|p\|}{\|p_i\|})\|v^x - v^y\| \leq (1 + \max_i \frac{\|p\|}{\|p_i\|})\|l\|\|x - y\|, \tag{50}$$

where $l = [l_1, l_2, \dots, l_m]^T$ is the vector of Lipschitz constants of all objectives. Denoting $L_m = (1 + \max_i \frac{\|p\|}{\|p_i\|})\|l\|$, we have

$$\|g_s(x) - g_s(y)\| \leq \left(\|\nabla_v \Psi(p, v^x)\| \sum_{i=1}^m l_{i,1} + C_g L_m \right) \|x - y\| \tag{51}$$

By definition in [Equation \(44\)](#):

$$\|\nabla_v \Psi(p, v^x)\| = \left\| \max_i \frac{v_i^x}{p_i} p - v^x \right\| \leq \left\| \max_v \max_i \frac{v_i}{p_i} p \right\|, \tag{52}$$

where the inequality is because the values of x should be no larger than the maximum values for the objectives. Denoting

$$L_p = \left\| \max_v \max_i \frac{v_i}{p_i} p \right\|, \tag{53}$$

we have

$$\|g_s(x) - g_s(y)\| \leq \left(L_p \sum_{i=1}^m l_{i,1} + C_g L_m \right) \|x - y\|, \tag{54}$$

and we define $L_s = (L_p \sum_{i=1}^m l_{i,1} + C_g L_m)$. We have proven $g_s(\cdot)$ is to be L_s -Lipshitz continuous. Therefore, both claims of this lemma have been proven. \square

I.2. Proof for Theorem 4.1 and 4.2

Before proving [Theorem 4.2](#), we prove [Lemma I.1](#) for the requirements to use the proof idea in ([Xiao et al., 2023](#)) for their theorem.3 and obtain [Lemma I.2](#). To be consistent with their proof, we consider minimizing the negative value and similarity with gradient descent.

Lemma I.1. *Under the Assumptions 4.1-4.3, we have*

$$\|g_s(\pi_{\mathbf{p},t})\| \leq L_{\mathbf{p}} C_g, \quad \mathbb{E} [\|g_s(\pi_{\mathbf{p},t}; \xi) - g_s(\pi_{\mathbf{p},t})\|^2] \leq L_{\mathbf{p}}^2 m \sigma^2 \quad (55)$$

$$\mathbb{E} [\|G(\pi_{\mathbf{p},t}; \xi)^T (G(\pi_{\mathbf{p},t}; \xi') w_t + \lambda g_s(\pi_{\mathbf{p},t}; \xi'))\|^2] \leq \underbrace{8(m\sigma^2 + C_g^2)^2 + 8L_{\mathbf{p}}^2 \lambda^2 (m\sigma^2 + C_g^2)^2}_{C_1} \quad (56)$$

$$\mathbb{E} [\|G(\pi_{\mathbf{p},t}; \zeta) w_t + \lambda g_s(\pi_{\mathbf{p},t}; \zeta)\|^2] \leq \underbrace{4m\sigma^2 + 4C_g^2 + 4\lambda^2 L_{\mathbf{p}} m \sigma^2 + 4\lambda^2 L_{\mathbf{p}}^2 C_g^2}_{C_2} \quad (57)$$

$$\mathbb{E} [\|(G(\pi_{\mathbf{p},t}) w + \lambda \pi_{\mathbf{p},t}(\pi_{\mathbf{p},t}))^T G(\pi_{\mathbf{p},t})\| \|w_t - w_{t+1}\|] \leq \underbrace{2(1 + L_{\mathbf{p}} \lambda)^2 C_g^2 (m\sigma + C_g)^2}_{C_2} \quad (58)$$

Proof. Under the Assumptions 4.1-4.3, by (52) and (53), we have

$$\|g_s(\pi_{\mathbf{p},t})\| \leq \|G(\pi_{\mathbf{p},t}) \nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}^{\pi_{\mathbf{p},t}})\| \leq L_{\mathbf{p}} C_g, \quad (59)$$

and

$$\mathbb{E} [\|g_s(\pi_{\mathbf{p},t}; \xi) - g_s(\pi_{\mathbf{p},t})\|^2] \leq \mathbb{E} [\|G(\pi_{\mathbf{p},t}; \xi) - G(\pi_{\mathbf{p},t})\|^2 \|\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}^{\pi_{\mathbf{p},t}})\|^2] \leq L_{\mathbf{p}}^2 m \delta^2 \quad (60)$$

The first two claims in (55) are proven.

Next, we have

$$\begin{aligned} & \mathbb{E} [\|G(\pi_{\mathbf{p},t}; \xi)^T (G(\pi_{\mathbf{p},t}; \xi') w_t + \lambda g_s(\pi_{\mathbf{p},t}; \xi'))\|^2] \\ & \stackrel{(i)}{\leq} 2 \mathbb{E} [\underbrace{\|G(\pi_{\mathbf{p},t}; \xi)^T G(\pi_{\mathbf{p},t}; \xi') w_t\|^2}_{N_1} + 2\lambda^2 \underbrace{\|G(\pi_{\mathbf{p},t}; \xi)^T g_s(\pi_{\mathbf{p},t}; \xi')\|^2}_{N_2}], \end{aligned} \quad (61)$$

where (i) is by the Young's inequality. Next, we provide bounds for $\mathbb{E}[N_1]$ and $\mathbb{E}[N_2]$, separately:

$$\begin{aligned} \mathbb{E}[N_1] & \stackrel{(i)}{\leq} \mathbb{E} [\|(G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T + G(\pi_{\mathbf{p},t})^T)(G(\pi_{\mathbf{p},t}; \xi') - G(\pi_{\mathbf{p},t}) + G(\pi_{\mathbf{p},t}))\|^2] \\ & = \mathbb{E} [\|(G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T)(G(\pi_{\mathbf{p},t}; \xi') - G(\pi_{\mathbf{p},t})) + (G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T)G(\pi_{\mathbf{p},t}) \\ & \quad + G(\pi_{\mathbf{p},t})^T(G(\pi_{\mathbf{p},t}; \xi') - G(\pi_{\mathbf{p},t})) + G(\pi_{\mathbf{p},t})^T G(\pi_{\mathbf{p},t})\|^2] \\ & \stackrel{(ii)}{\leq} 4 \mathbb{E} [\|G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T\|^2 \|G(\pi_{\mathbf{p},t}; \xi') - G(\pi_{\mathbf{p},t})\|^2 + \|G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T\|^2 \|G(\pi_{\mathbf{p},t})\|^2 \\ & \quad + \|G(\pi_{\mathbf{p},t})^T\|^2 \|(G(\pi_{\mathbf{p},t}; \xi') - G(\pi_{\mathbf{p},t}))\|^2 + \|G(\pi_{\mathbf{p},t})^T G(\pi_{\mathbf{p},t})\|^2] \\ & \stackrel{(iii)}{\leq} 4m^2 \sigma^4 + 8m\sigma^2 C_g^2 + 4C_g^4 = 4(m\sigma^2 + C_g^2)^2, \end{aligned} \quad (62)$$

where (i) follows from Cauchy–Schwarz inequality and w_t is a convex coefficient, (ii) follows from Young's inequality and (iii) follows from [Assumption 4.2](#) and [Assumption 4.3](#). For another term,

$$\begin{aligned} \mathbb{E}[N_2] & = \mathbb{E} [\|(G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T + G(\pi_{\mathbf{p},t})^T)(g_s(\pi_{\mathbf{p},t}; \xi') - g_s(\pi_{\mathbf{p},t}) + g_s(\pi_{\mathbf{p},t}))\|^2] \\ & \stackrel{(i)}{\leq} 4 \mathbb{E} [\|(G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T)(g_s(\pi_{\mathbf{p},t}; \xi') - g_s(\pi_{\mathbf{p},t}))\|^2 + \|(G(\pi_{\mathbf{p},t}; \xi)^T - G(\pi_{\mathbf{p},t})^T)g_s(\pi_{\mathbf{p},t})\|^2 \\ & \quad + \|G(\pi_{\mathbf{p},t})^T(g_s(\pi_{\mathbf{p},t}; \xi') - g_s(\pi_{\mathbf{p},t}))\|^2 + \|G(\pi_{\mathbf{p},t})^T g_s(\pi_{\mathbf{p},t})\|^2] \\ & \stackrel{(ii)}{\leq} 4L_{\mathbf{p}}^2 m^2 \sigma^4 + 8L_{\mathbf{p}}^2 m \sigma^2 C_g^2 + 4L_{\mathbf{p}}^2 C_g^4 = 4L_{\mathbf{p}}^2 (m\sigma^2 + C_g^2)^2, \end{aligned} \quad (63)$$

where (i) follows from Young's inequality, (ii) follows from (59) and (60). Then substituting (62) and (63) into (61), we can obtain,

$$\mathbb{E}[\|G(\pi_{\mathbf{p},t}; \xi)^T (G(\pi_{\mathbf{p},t}; \xi') \mathbf{w}_t + \lambda g_s(\pi_{\mathbf{p},t}; \xi'))\|^2] \leq 8(m\sigma^2 + C_g^2)^2 + 8L_p^2\lambda^2(m\sigma^2 + C_g^2)^2 = C_1.$$

We have proved (56). Then, we look at (57) :

$$\begin{aligned} & \mathbb{E}[\|G(\pi_{\mathbf{p},t}; \zeta) \mathbf{w}_t + \lambda g_s(\pi_{\mathbf{p},t}; \zeta)\|^2] \\ &= \mathbb{E}[\|G(\pi_{\mathbf{p},t}; \zeta) \mathbf{w}_t - G(\pi_{\mathbf{p},t}) \mathbf{w}_t + G(\pi_{\mathbf{p},t}) \mathbf{w}_t + \lambda g_s(\pi_{\mathbf{p},t}; \zeta) - \lambda g_s(\pi_{\mathbf{p},t}) + \lambda g_s(\pi_{\mathbf{p},t})\|^2] \\ &\stackrel{(i)}{\leq} 4\mathbb{E}[\|G(\pi_{\mathbf{p},t}; \zeta) - G(\pi_{\mathbf{p},t})\|^2] + 4\mathbb{E}[\|G(\pi_{\mathbf{p},t})\|^2] + 4\lambda^2\mathbb{E}[\|g_s(\pi_{\mathbf{p},t}; \zeta) - g_s(\pi_{\mathbf{p},t})\|^2] \\ &\quad + 4\lambda^2\mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] \\ &\stackrel{(ii)}{\leq} \underbrace{4m\sigma^2 + 4C_g^2 + 4\lambda^2 L_p^2 m\sigma^2 + 4\lambda^2 L_p^2 C_g^2}_{C_2} \end{aligned} \tag{64}$$

where (i) follows from Young's inequality, and (ii) follows from 59 and 60.

Finally,

$$\begin{aligned} & \mathbb{E}[\|(G(\pi_{\mathbf{p},t}) \mathbf{w} + \lambda \pi_{\mathbf{p},t}(\pi_{\mathbf{p},t}))^T G(\pi_{\mathbf{p},t})\| \|\mathbf{w}_t - \mathbf{w}_{t+1}\|] \\ &= \beta_t \mathbb{E}[\|(G(\pi_{\mathbf{p},t}) \mathbf{w} + \lambda \pi_{\mathbf{p},t}(\pi_{\mathbf{p},t}))^T G(\pi_{\mathbf{p},t})\| \|G(\pi_{\mathbf{p},t}; \xi)^T (G(\pi_{\mathbf{p},t}; \xi') \mathbf{w}_t + \lambda \pi_{\mathbf{p},t}(\theta; \xi'))\|] \\ &\leq \beta_t \mathbb{E}[\|(G(\pi_{\mathbf{p},t}) \mathbf{w} + \lambda \pi_{\mathbf{p},t}(\pi_{\mathbf{p},t}))^T G(\pi_{\mathbf{p},t})\| (\|G(\pi_{\mathbf{p},t}; \xi)^T (G(\pi_{\mathbf{p},t}; \xi') \mathbf{w}_t\| + \lambda \|G(\pi_{\mathbf{p},t}; \xi) \pi_{\mathbf{p},t}(\theta; \xi')\|))] \\ &= \beta_t \mathbb{E}[\|(G(\pi_{\mathbf{p},t}) \mathbf{w} + \lambda \pi_{\mathbf{p},t}(\pi_{\mathbf{p},t}))^T G(\pi_{\mathbf{p},t})\| (\sqrt{N_1} + \sqrt{N_2})] \\ &\leq \beta_t 2(1 + L_p \lambda)^2 C_g^2 (m\sigma + C_g)^2 = \beta_t C_3, \end{aligned} \tag{65}$$

□

Lemma I.2. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}} T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1} T^{-\frac{1}{2}})$, the updates by our method satisfy

$$\begin{aligned} E[\|G(\pi_{\mathbf{p},t}) \mathbf{w}_{t,\lambda} + \lambda_t g_s(\pi_{\mathbf{p},t})\|^2] &\leq \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\ &\quad + \frac{\beta_t}{2} C_1(\lambda_t) + \frac{l'_1 \alpha_t}{2} C_2(\lambda_t) + \beta_t C_3(\lambda_t) \end{aligned} \tag{66}$$

where w is a fixed convex coefficient, and

$$l'(\pi_{\mathbf{p},t}) = -\mathbf{w}^T \mathbf{v}(\pi_{\mathbf{p},t}) - \lambda_t \Psi(\mathbf{p}, \pi_{\mathbf{p},t}), \tag{67}$$

$$l'_1 = \max_i l_{i,1} + \lambda L_s \tag{68}$$

$$C_1 = 8(m\sigma^2 + C_g^2)^2 + 8L_p^2\lambda^2(m\sigma^2 + C_g^2)^2, \tag{69}$$

$$C_2 = 4m\sigma^2 + 4C_g^2 + 4\lambda^2 L_p^2 m\sigma^2 + 4\lambda^2 L_p^2 C_g^2, \tag{70}$$

$$C_3 = 2(1 + L_p \lambda)^2 C_g^2 (m\sigma + C_g)^2, \tag{71}$$

where L_s is the Lipschitz constant for $g_s(\cdot)$, defined in (54).

Under Assumptions.(4.1-4.3), previous results show Lemma 4.1 and Lemma I.1 hold. Therefore, we can replace g_0 in their analysis with g_s and apply their (33) in our case and it becomes Equation (66). Stochastic gradient samples like $G(\pi_{\mathbf{p},t}, \xi)$ have been taken expectations and become $G(\pi_{\mathbf{p},t})$ or σ .

This is an intuitive result of convergence analysis for smooth non-convex objective functions using conventional techniques. Next we prove our main theoretical contributions based on Lemma I.2.

Theorem 4.1. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function $\Psi(\mathbf{p}, \cdot)$, we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|] = \mathcal{O}(mT^{-\frac{1}{2}})$. To achieve an ϵ -accurate Pareto stationary point, it requires $T = \mathcal{O}(m^2\epsilon^{-2})$ updates.

Proof. By definition,

$$\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}) = \max_i \frac{\mathbf{v}_i}{\mathbf{p}_i} \mathbf{p} - \mathbf{v} > 0. \quad (72)$$

So $g_s(\pi_{\mathbf{p}}) = G(\pi_{\mathbf{p}})\nabla_{\mathbf{v}} \Psi(\mathbf{p}, \mathbf{v}^{\pi_{\mathbf{p}}})$ can be considered as a positive linear combination of objective gradients. We have

$$\begin{aligned} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda g_s(\pi_{\mathbf{p},t})\|^2] &= E[\|(G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda G(\pi_{\mathbf{p},t})\tilde{\mathbf{w}}_t)\|^2] \\ &\geq E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2]. \end{aligned} \quad (73)$$

For every time step t , by Equation (66) from Lemma I.2 and constant λ ,

$$\begin{aligned} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] &\leq \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\ &\quad + \frac{\beta_t}{2} C_1(\lambda) + \frac{l'_1 \alpha_t}{2} C_2(\lambda) + \beta_t C_3(\lambda) \end{aligned} \quad (74)$$

We take $\alpha_t = \alpha$ and $\beta_t = \beta$ as constants and telescope (74),

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] &\leq \frac{1}{\alpha T} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \frac{1}{2\beta T} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\ &\quad + \frac{1}{T} \sum_{t=0}^{T-1} \frac{\beta}{2} C_1(\lambda) + \frac{1}{T} \sum_{t=0}^{T-1} \frac{l'_1 \alpha}{2} C_2(\lambda) + \frac{1}{T} \sum_{t=0}^{T-1} \beta C_3(\lambda) \\ &\leq \mathcal{O}\left(\frac{1}{\alpha T} + \frac{1}{\beta T} + \beta m^2 + \alpha m\right) \end{aligned} \quad (75)$$

By setting $\alpha = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta = \Theta(m^{-1}T^{-\frac{1}{2}})$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] = \mathcal{O}(mT^{-\frac{1}{2}}).$$

To achieve an ϵ -accurate Pareto stationary point, it requires $T = \mathcal{O}(m^2\epsilon^{-2})$ updates. \square

After proving for cases with constant λ , we need to prove further for cases with increasing $\lambda = \Theta(T^{\frac{1}{2}})$.

Theorem 4.2. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a Lipschitz smooth similarity function with $g'_s(\pi_{\mathbf{p},t})$ being convex combination of $g_i(\pi_{\mathbf{p},t})$ for all t , there can be an increasing $\lambda = \Theta(\log T)$ and we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|] = \mathcal{O}(mT^{-\frac{1}{2}} \log T)$.

Proof. Because the similarity gradients $g'_s(\pi_{\mathbf{p},t})$ are convex combinations of $G(\pi_{\mathbf{p},t})$, let $g'_s(\pi_{\mathbf{p},t}) = G(\pi_{\mathbf{p},t})\tilde{\mathbf{w}}_t$ where $\tilde{\mathbf{w}}_t$ is a convex coefficient, then

$$\begin{aligned} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda_t g'_s(\pi_{\mathbf{p},t})\|^2] &= E[\|(G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda_t G(\pi_{\mathbf{p},t})\tilde{\mathbf{w}}_t)\|^2] \\ &\geq E[(1 + \lambda_t)^2 \min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] \end{aligned} \quad (76)$$

holds because $(\mathbf{w}_{t,\lambda} + \lambda_t \tilde{\mathbf{w}}_t)$ is also a convex coefficient which can not be more optimal than $\arg \min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2$. For

every time step t , by Equation (66) from Lemma I.2,

$$\begin{aligned}
 E[(1 + \lambda_t)^2 \min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] &\leq \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\
 &\quad + \frac{\beta_t}{2} C_1(\lambda_t) + \frac{l'_1 \alpha_t}{2} C_2(\lambda_t) + \beta_t C_3(\lambda_t) \\
 E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] &\leq \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\
 &\quad + \frac{\beta_t}{2(1 + \lambda_t)^2} C_1(\lambda_t) + \frac{l'_1 \alpha_t}{2(1 + \lambda_t)^2} C_2(\lambda_t) + \frac{\beta_t}{(1 + \lambda_t)^2} C_3(\lambda_t) \\
 &\leq \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\
 &\quad + \frac{\beta_t}{2\lambda_t^2} C_1(\lambda_t) + \frac{l'_1 \alpha_t}{2\lambda_t^2} C_2(\lambda_t) + \frac{\beta_t}{(1 + \lambda_t)^2} C_3(\lambda_t) \\
 &\stackrel{(i)}{\leq} \frac{1}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{2\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\
 &\quad + \frac{\beta_t}{2\lambda_t^2} C_1(\lambda_t) + \frac{(\max_i l_{i,1} + \lambda_T L_s)\alpha_t}{2\lambda_t^2} C_2(\lambda_t) + \frac{\beta_t}{(1 + \lambda_t)^2} C_3(\lambda_t)
 \end{aligned} \tag{77}$$

where (i) is by the definition of l'_1 in (68). In the proofs of Theroem 1 and 3 of (Xiao et al., 2023), l'_1 was considered as constant. However, for more rigor (it increases with λ), we upper bound it here with $\mathcal{O}(\log T)$. We take $\alpha_t = \alpha$ and $\beta_t = \beta$ as constants and telescope (77), and by $\lambda_t = \Theta(\log t)$, we have

$$\begin{aligned}
 \frac{1}{T} \sum_{t=0}^{T-1} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] &\leq \frac{1}{\alpha T} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \frac{1}{2\beta T} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\
 &\quad + \frac{1}{T} \sum_{t=0}^{T-1} \frac{\beta}{2\lambda_t^2} C_1(\lambda_t) + \frac{1}{T} \sum_{t=0}^{T-1} \frac{(\max_i l_{i,1} + \lambda_T L_s)\alpha}{2\lambda_t^2} C_2(\lambda_t) + \frac{1}{T} \sum_{t=0}^{T-1} \frac{\beta}{(1 + \lambda_t)^2} C_3(\lambda_t) \\
 &= \mathcal{O}\left(\frac{1}{\alpha T} + \frac{1}{\beta T} + \frac{\beta m^2}{\log T} + \beta m^2 + \alpha m \log T + \alpha m\right)
 \end{aligned} \tag{78}$$

By setting $\alpha = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta = \Theta(m^{-1}T^{-\frac{1}{2}})$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] = \mathcal{O}(mT^{-\frac{1}{2}} \log T),$$

and proof is done. \square

Remark I.1. The convergence rate for Theorem 4.2 seems slower than results from (Xiao et al., 2023) because we rigorously considered the changes in the Lipschitz constant of the $(G(\pi_{\mathbf{p},t})\mathbf{w}_t + \lambda g_s(\pi_{\mathbf{p},t}))$ caused by increasing λ .

I.3. Proof for Theorem 4.3 and 4.4

To be consistent with previous results in MOO literature, we consider minimizing the negative objectives and similarity with gradient descent.

Theorem 4.3. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}}T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1}T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function like $\Psi(\mathbf{p}, \cdot)$, we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|] - \frac{2C_g^2}{\lambda^2} = \mathcal{O}(mT^{-\frac{1}{2}})$.

Proof. By Equation (66) from Lemma I.2 and constant λ , we have

$$\begin{aligned}\mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] &\leq \frac{2}{\lambda^2} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda g_s(\pi_{\mathbf{p},t})\|^2] + \frac{2}{\lambda^2} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda}\|^2] \\ &\leq \frac{2}{\lambda^2 \alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{\lambda^2 \beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\ &\quad + \frac{\beta_t}{\lambda^2} C_1(\lambda) + \frac{l'_1 \alpha}{\lambda^2} C_2(\lambda) + \frac{2\beta_t}{\lambda^2} C_3(\lambda) + \frac{2C_g^2}{\lambda^2}.\end{aligned}\tag{79}$$

By the definition of l'_1 in (68), it is a constant when λ is constant. Take $\alpha_t = \alpha$ and $\beta_t = \beta$ as constants telescope (79), we get

$$\begin{aligned}\frac{1}{T} \sum_{t=t_0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] - \frac{2C_g^2}{\lambda^2} &\leq \frac{2}{\lambda^2 \alpha T} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \frac{1}{\lambda^2 \beta T} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\ &\quad + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{\beta}{\lambda^2} C_1(\lambda) + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{l'_1 \alpha}{\lambda^2} C_2(\lambda) + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{2\beta}{\lambda^2} C_3(\lambda) + \frac{2C_g^2}{\lambda^2} \\ &\leq \mathcal{O}\left(\frac{1}{\alpha T} + \frac{1}{\beta T} + \beta m^2 + \alpha m\right)\end{aligned}\tag{80}$$

By setting $\alpha = \Theta(m^{-\frac{1}{2}} T^{-\frac{1}{2}})$, $\beta = \Theta(m^{-1} T^{-\frac{1}{2}})$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} E[\min_{\mathbf{w}_t} \|G(\pi_{\mathbf{p},t})\mathbf{w}_t\|^2] - \frac{2C_g^2}{\lambda^2} = \mathcal{O}(m T^{-\frac{1}{2}}).$$

To achieve an ϵ -accurate stationary point, it requires $T = \mathcal{O}(m^2 \epsilon^{-2})$ updates. \square

Theorem 4.4. Under the Assumptions 4.1-4.3, setting $\alpha_t = \Theta(m^{-\frac{1}{2}} T^{-\frac{1}{2}})$, $\beta_t = \Theta(m^{-1} T^{-\frac{1}{2}})$, with a constant λ and Lipschitz smooth similarity function like $\Psi(\mathbf{p}, \cdot)$, there can be an increasing $\lambda = \Theta(T^{\frac{1}{2}})$ and we have $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p}})\|] = \mathcal{O}(m T^{-\frac{1}{2}} \log T)$.

Proof. Suppose for all time steps $t > t_0$, $\lambda_t \geq 1$, by Equation (79) we have

$$\begin{aligned}\mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] &\leq \frac{2}{\lambda_t^2} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda} + \lambda_t g_s(\pi_{\mathbf{p},t})\|^2] + \frac{2}{\lambda_t} \mathbb{E}[\|G(\pi_{\mathbf{p},t})\mathbf{w}_{t,\lambda}\|^2] \\ &\leq \frac{2}{\alpha_t} \mathbb{E}[l'(\pi_{\mathbf{p},t}) - l'(\pi_{\mathbf{p},t+1})] + \frac{1}{\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2] \\ &\quad + \frac{\beta_t}{\lambda_t^2} C_1(\lambda_t) + \frac{l'_1 \alpha}{\lambda_t^2} C_2(\lambda_t) + \frac{2\beta_t}{\lambda_t^2} C_3(\lambda_t) + \frac{2C_g^2}{\lambda_t^2}.\end{aligned}\tag{81}$$

Take $\alpha_t = \alpha$ and $\beta_t = \beta$ as constants telescope (81) and by the definition of l'_1 in (68) we get

$$\begin{aligned}\sum_{t=t_0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] &\leq \sum_{t=t_0}^{T-1} \frac{2}{\alpha} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \sum_{t=t_0}^{T-1} \frac{1}{\beta} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\ &\quad + \sum_{t=t_0}^{T-1} \frac{\beta}{\lambda_t^2} C_1(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{l'_1 \alpha}{\lambda_t^2} C_2(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{2\beta}{\lambda_t^2} C_3(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{2C_g^2}{\lambda_t^2} \\ &\leq \sum_{t=t_0}^{T-1} \frac{2}{\alpha} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \sum_{t=t_0}^{T-1} \frac{1}{\beta} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\ &\quad + \sum_{t=t_0}^{T-1} \frac{\beta}{\lambda_t^2} C_1(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{(\max_i l_{i,1} + \lambda_T L_s) \alpha}{\lambda_t^2} C_2(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{2\beta}{\lambda_t^2} C_3(\lambda_t) + \sum_{t=t_0}^{T-1} \frac{2C_g^2}{\lambda_t^2},\end{aligned}\tag{82}$$

then by $\lambda_t = \Theta(\log t)$ we have:

$$\begin{aligned} \frac{1}{T} \sum_{t=t_0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] &\leq \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{2}{\alpha} \mathbb{E}[l'(\pi_{\mathbf{p},0}) - l'(\pi_{\mathbf{p},T})] + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{1}{\beta} \mathbb{E}[\|\mathbf{w}_0 - \mathbf{w}\|^2 - \|\mathbf{w}_T - \mathbf{w}\|^2] \\ &+ \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{\beta}{\lambda_t^2} C_1(\lambda_t) + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{(\max_i l_{i,1} + \lambda_T L_s)\alpha}{\lambda_t^2} C_2(\lambda_t) + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{2\beta}{\lambda_t^2} C_3(\lambda_t) + \frac{1}{T} \sum_{t=t_0}^{T-1} \frac{2C_g^2}{\lambda_t^2} \quad (83) \\ &= \mathcal{O}\left(\frac{1}{\alpha T} + \frac{1}{\beta T} + \frac{\beta m^2}{\log T} + \beta m^2 + \alpha m \log T + \alpha m + \frac{1}{(\log T)^2}\right) \end{aligned}$$

Adding the terms before t_0 , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] &= \frac{1}{T} \sum_{t=0}^{t_0} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] + \frac{1}{T} \sum_{t=t_0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] \quad (84) \\ &\leq \mathcal{O}\left(\frac{1}{T} + \frac{1}{\alpha T} + \frac{1}{\beta T} + \frac{\beta m^2}{\log T} + \beta m^2 + \alpha m \log T + \alpha m + \frac{1}{(\log T)^2}\right). \end{aligned}$$

By setting $\alpha = \Theta(m^{-\frac{1}{2}} T^{-\frac{1}{2}})$, $\beta = \Theta(m^{-1} T^{-\frac{1}{2}})$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|g_s(\pi_{\mathbf{p},t})\|^2] = \mathcal{O}(m T^{-\frac{1}{2}} \log T), \quad (85)$$

and the proof is done. \square