# One Model for All: Multi-Objective Controllable Language Models

**Anonymous authors**
**Paper under double-blind review**

### Abstract

Aligning large language models (LLMs) with human preferences is critical to enhancing LLMs' safety, helpfulness, humor, faithfulness, etc. Current reinforcement learning from human feedback (RLHF) mainly focuses on a fixed reward learned from average human ratings, which may weaken the adaptability and controllability of varying preferences. However, creating personalized LLMs requires aligning LLMs with individual human preferences, which is non-trivial due to the scarce data per user and the diversity of user preferences in multi-objective trade-offs, such as prioritizing humor and empathy in one context, while seeking efficiency and precision in another. *Can we train one LLM to produce personalized outputs for different user preferences on the Pareto front?* In this paper, we introduce <u>M</u>ulti-<u>O</u>bjective <u>C</u>ontrol (MOC), which trains a single LLM to directly generate responses in the preference-defined regions of the Pareto front. Our approach introduces multi-objective optimization (MOO) principles into RLHF to train an LLM as a preference-conditioned policy network. We improve the computational efficiency of MOC by applying MOO at the policy level, enabling us to fine-tune a 7B-parameter model on a single A6000 GPU. Extensive experiments demonstrate the advantages of MOC over baselines in three aspects: (i) controllability of LLM outputs w.r.t. user preferences on the trade-off among multiple rewards; (ii) quality and diversity of LLM outputs, measured by the hyper-volume of multiple solutions achieved; and (iii) generalization to unseen preferences. These results highlight MOC's potential for real-world applications requiring scalable and customizable LLMs.

## 1 Introduction

Large language models (LLMs) have gained significant attention for their impressive performance across a wide range of tasks, including machine translation (Vaswani et al., 2017; Radford & Narasimhan, 2018; Devlin et al., 2019), text generation (Touvron et al., 2023; OpenAI, 2023), and conversational agents (Ouyang et al., 2022; Bai et al., 2022). However, these models are generally aligned with fixed preferences predetermined by developers (Ouyang et al., 2022; Touvron et al., 2023; Bai et al., 2023; Dubey et al., 2024), limiting the degree of personalization available to users. In real-world scenarios, users often have diverse preferences for LLMs behaviors. For instance, one user might prefer a humorous and empathetic response for emotional support, while another might prioritize a more efficient, task-oriented assistant. Despite this variability, the inherent flexibility of current LLMs (Dubey et al., 2024; OpenAI, 2023) is limited in providing fully personalized interactions.

The ability of LLMs to dynamically balance the trade-offs between different objectives according to diverse user-specified preferences is called *multi-objective controllability*, a crucial feature for enhancing user satisfaction. Training separate models for each preference order, however, is neither practical nor scalable due to the high computational costs. This highlights the need to enable one-time LLMs training while accommodating a broad range of preferences.

Can we **control** the trade-offs in a single, once-trained LLM to meet diverse human preferences? Our answer is **yes**. This paper aims to (i) enable LLMs to generate customized responses for diverse user preferences and (ii) achieve this goal with a once-trained model. To this end, we introduce a novel algorithm, <u>M</u>ulti-<u>O</u>bjective

Table 1: Comparison with the state-of-the-art MOO methods. MOC addresses MOO more principally and efficiently. $M-$ the number of preferences, $N-$ the number of reward models (objectives).

| Algorithms | Explicit policy improvement | Num of trained LLMs | Inference adaptation | Preference data | Loss |
|---|---|---|---|---|---|
| MORLHF | ✓ | $M$ | × | No | PPO |
| Rewarded Soups (Ramé et al., 2023) | × | $N$ | ✓ | No | PPO |
| MODPO (Zhou et al., 2024) | ✓ | $M$ | × | Yes | DPO |
| RiC (Yang et al., 2024b) | × | 1 | ✓ | No | SFT |
| MOC (Ours) | ✓ | 1 | ✓ | No | PPO |

Control (MOC), which leverages a carefully designed multi-objective optimization (MOO) algorithm. MOC requires only one training, incorporates explicit policy improvement, and does not rely on human preference data. Moreover, its training cost is comparable to single-objective RLHF (Ouyang et al., 2022), and we made it feasible to fine-tune a 7-billion model on a single A6000 GPU with LoRA (Hu et al., 2022).

We first formulate multi-objective controllability as an MOO problem with preference vector constraints, inspired by recent advancements in MOO (Désidéri, 2009; Sener & Koltun, 2018; Xiao et al., 2023). This formulation presents two primary challenges. The first one is identifying the target to be controlled. Existing MOO methods typically focus on optimizing different loss functions(Liu et al., 2021; 2023) or linearized utility functions (Yang et al., 2019), which do not effectively capture the quality or behavior of LLMs. In contrast, MOC selects the reward signal as the control target, enabling direct manipulation of the model's behavior. The second challenge is to solve this optimization problem within feasible computational limits. Our formulated optimization problem involves complex trade-offs among multiple objectives under different preference constraints. To address this, we relax the problem into a new form of MOO, where the preference constraint is treated as an additional objective. MOC scalarizes the objectives with dynamic weighting in different steps, ensuring that the computational cost is comparable to the widely used single-objective RLHF (Ouyang et al., 2022). Table 1 presents a detailed comparison of MOC and baseline methods.

In extensive experiments, MOC consistently outperforms baseline methods (Ouyang et al., 2022; Ramé et al., 2023; Yang et al., 2024b) across multiple tasks. It demonstrates strong performance in three key areas: (i) controllability, as it effectively aligns model behavior with diverse preference vectors and ensures a clear monotonic relationship between input preferences and outcomes; (ii) quality of solution set, measured by the hyper-volume metric, where MOC achieves a superior Pareto front while maintaining a diverse set of solutions; and (iii) generalization, as it robustly handles unseen preferences while preserving both the alignment quality and diversity. Compared to baseline methods, MOC offers a more efficient and flexible approach to personalizing LLMs, managing different trade-offs among multiple objectives within a single model and seamlessly adapting to new preferences. These findings highlight MOC's potential for real-world applications requiring scalable and customizable personalization.

Our contributions are as follows: (i) We introduce the MOC algorithm, which takes comparable computation as single-objective RLHF and finetunes LLMs only once to accommodate diverse user preferences; (ii) We empirically demonstrate MOC's superior performance in terms of controllability, solution quality, and generalization, including its ability to generalize to unseen user preferences.

## 2 Multi-Objective Controllable Language Models

Modern LLM alignment is commonly implemented via RLHF (Ouyang et al., 2022): a reward model is first fit from preference pairs by maximizing

$$\mathcal{L}_{\mathrm{RM}} = \mathbb{E}_{(x,y^w,y^l)\sim\mathcal{D}}\Big[\log\sigma\big(r(x,y^w) - r(x,y^l)\big)\Big],$$

and the policy is then optimized against this learned reward, typically with PPO (Schulman et al., 2017),

$$\arg\max_{\pi(\cdot|x;\theta)} \mathbb{E}_{x\sim\mathcal{D},\, y\sim\pi(\cdot|x)}\left[r(x,y) - \beta\log\frac{\pi(y\mid x;\theta)}{\pi_{\mathrm{old}}(y\mid x)}\right].$$

While effective, this pipeline implicitly commits to a *single*, developer-chosen average preference. The resulting model is well-aligned to that preference but not easily *controllable*: at inference time it does not natively adapt when different users prioritize different objectives (e.g., helpfulness vs. conciseness, safety vs. creativity).

We therefore ask whether a model trained *once* can later be steered to different trade-offs on demand. Our answer is **yes**: we aim (i) to support a broad spectrum of user preferences, and (ii) to achieve this within one training run.

**Controllability vs. Alignment.** We use the term *controllability* to mean that the model systematically varies its behavior in response to user-specified preferences, producing outputs consistent with those preferences. By contrast, *alignment* refers to optimizing toward a fixed, global preference that remains unchanged at inference.

## 2.1 Problem Formulation

To represent user preferences, we define a preference vector $\mathbf{p} = [p_1, p_2, \cdots, p_N], \sum_{i=1}^N p_i = 1, p_i \geq 0$, where each element in $\mathbf{p}$ reflects the importance of a specific objective. Inspired by recent work on multi-objective learning (Xu et al., 2020; Ma et al., 2020; Yang et al., 2022), we use this preference vector to regulate the model's output in the objective space. Given a set of preference vectors $\{\mathbf{p}^i\}_{i=1}^M$, training a controllable LLM by the preference vectors is formulated as the following optimization problem:

$$
\begin{aligned}
\max_\theta \mathbf{J}(\pi(\cdot; \theta, \mathbf{p})) &\stackrel{\text{def}}{=} \max_\theta \big(J^1(\pi(\cdot; \theta, \mathbf{p})), J^2(\pi(\cdot; \theta, \mathbf{p})), \cdots, J^N(\pi(\cdot; \theta, \mathbf{p}))\big)^\top, \\
\text{s.t. } \Phi\Big(\pi(\cdot; \theta, \mathbf{p}) \,\Big\|\, \mathbf{p}\Big) &\leq \phi, \ \forall \mathbf{p} \in \{\mathbf{p}^1, \mathbf{p}^2, \cdots, \mathbf{p}^M\},
\end{aligned}
\tag{1}
$$

where $J^i$ denotes the RLHF objective associated with reward $R^i$. Equation (1) ensures LLMs align with a set of preferences. The LLM is a policy $\pi$ parameterized by $\theta$ and takes a preference vector $\mathbf{p}$ as an input condition. In addition, $\Phi$ is a distance or divergence metric between the policy $\pi$ and the preference vector $\mathbf{p}$, and the controllability requires a distance upper bounded by $\phi$. The objective $J^i$ is typically selected as a PPO loss (Schulman et al., 2017; Ouyang et al., 2022), unless stated otherwise.

Traditional methods for solving constrained optimization problems, such as the Lagrangian method, are inefficient for handling the complexity of Equation (1) due to the multiple constraints, diverse preferences, and the high dimensionality of language model parameters. This limitation renders developing new solutions imperative.

## 2.2 What Should the Preference Vector Align With?

Existing multi-objective learning methods (Yang et al., 2019; Liu et al., 2023; 2021) typically focus on balancing multiple loss functions. However, RL loss does not reflect the agent's true performance and is therefore not suitable as the target for controllability. In contrast, the value function or episodic return provides a reliable performance measure. In the context of RLHF for LLMs, the reward is evaluated by a reward model, which serves as the episodic return. Therefore, we choose a multi-dimensional reward signal as the primary target for control. To maintain simplicity, we select mean squared error (MSE) as the similarity metric between the reward signal and the preference vector. Formally, the constraint in Equation (1) is specified as

$$
\Phi\Big(\pi(\cdot; \theta, \mathbf{p}) \| \mathbf{p}\Big) \stackrel{\text{def}}{=} MSE\Big(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x, y), \mathbf{p}\Big) \leq \phi,
\tag{2}
$$

where $x$ represents the prompt/query, $y \sim \pi(x; \theta, \mathbf{p})$ is LLM-generated response, and $\mathcal{D}$ is the prompt dataset. The reward vector $\mathbf{R}(x, y) = (R^1(x, y), R^2(x, y), \cdots, R^N(x, y))$ is associated with the $N$ optimization objectives $\{J^i\}_{i=1}^N$. The sampled response $y$ depends on the policy parameters $\theta$, enabling the optimization of $\mathbf{J}(\pi)$ with respect to $\theta$ via standard RLHF. Equation (2) enforces that the reward vector $\mathbf{R}(x, y)$ aligns closely with the preference vector $\mathbf{p}$. In other contexts such as typical RL settings, the value function can be the target of control. Further details are provided in Appendix E.

**Re-labeling the prompt.** In MOC, the policy $\pi$ takes an additional condition: the user's preference vector $\mathbf{p} = [p_1, p_2, \cdots, p_N]$. To accommodate this, we modify the original prompt by appending the preference

vector to it, as follows:

$$\text{Re-labeled prompt} = \langle R1 \rangle p_1 \langle R2 \rangle p_2 \cdots \langle RN \rangle p_N \{\text{prompt}\}. \tag{3}$$

## 2.3 Multi-Objective Controllability

To solve the multi-objective learning problem with inequality constraints, we introduce our MOC algorithm, which builds on recent advances of multi-objective learning (Désidéri, 2009; Sener & Koltun, 2018). MOC simultaneously optimizes all the objectives while maximizing the alignment between the objective value vector and the preference vector. For simplicity, we optimize the following similarity objective:

$$\max_{\theta} J^{\Phi} \stackrel{\text{def}}{=} \max_{\theta} -ReLU(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x, y), \mathbf{p}) - \phi), \tag{4}$$

where $ReLU(x) = \max(x, 0)$ penalizes constraint violations when the error exceeds the threshold $\phi$. This ensures that optimization respects the trade-offs between rewards and preferences. The gradient of Equation (4) can be approximated as

$$\nabla_{\theta} \text{ReLU}(MSE(\mathbb{E}\mathbf{R}(x, y), \mathbf{p}) - \phi) = \mathbf{1}_{MSE(\mathbb{E}\mathbf{R}(x,y),\mathbf{p})-\phi>0} \sum_{k=1}^{N} [(\mathbb{E}R^k - p_k) \nabla_{\theta} \mathbb{E}R^k(x, y)] \tag{5}$$

where $\mathbf{1}_{(\cdot)}$ is the indicator function, $R^k$ represents the $k$-th entry of $\mathbf{R}$, $p_k^j$ means the $k$-th entry of preference vector $p^j$. $\mathbb{E}$ is short for $\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(x; \theta, p)}$ in Equation (5). The term $\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(x; \theta, p)} R^k(x, y)$ aims at maximizing the corresponding reward, which is also the goal of the gradient of the PPO loss. Thus, one could use the PPO objective $\nabla_{\theta} J^k(\pi(\cdot; \theta, \mathbf{p})$ to compute $\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(x; \theta, p)} R^k(x, y)$.

Solving the original optimization problem in Equation (2) is computationally challenging due to the involvement of $N$ objectives and $M$ preferences. Thus, we reformulate it as

$$\max_{\theta} \widehat{\mathbf{J}}(\pi(\cdot; \theta, \mathbf{p})) \stackrel{\text{def}}{=} \max_{\theta} \left( \mathbf{p}^{\top} \mathbf{J}(\pi(\cdot; \theta, \mathbf{p})), -\text{ReLU}\left(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x, y), \mathbf{p}) - \phi\right) \right)^{\top}, \tag{6}$$

where $\mathbf{J}(\pi(\cdot; \theta, \mathbf{p}))$ is defined in Equation (2). This reformulation offers two significant advantages: (i) It significantly reduces optimization complexity by transforming the original $N$-objective optimization into a bi-objective optimization; (ii) It retains control over the preference vectors in the newly formulated optimization problem. Scalarization simplifies the problem even further:

$$\max_{\theta} \left\{ c^{(1)} \mathbf{p}^{\top} \mathbf{J}(\pi(\cdot; \theta, \mathbf{p})) - c^{(2)} \text{ReLU}\left(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x, y), \mathbf{p}) - \phi\right) \Big| \sum_{i=1}^{2} c^{(i)} = 1, \ c^{(i)} \geq 0 \right\}, \tag{7}$$

where $c^{(i)}$ is an $i$-objective related co-efficient, determined by solving a min-norm problem

$$\min_{c^{(1)}, c^{(2)}} \left\{ \left\| c^{(1)} \mathbf{p}^{\top} \nabla_{\theta} \mathbf{J}(\pi(\cdot; \theta, \mathbf{p})) - c^{(2)} \nabla_{\theta} \text{ReLU}\left(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x, y), \mathbf{p}) - \phi\right) \right\|_2^2 \Big| \sum_{i=1}^{2} c^{(i)} = 1, \ c^{(i)} \geq 0 \right\}. \tag{8}$$

As demonstrated by Désidéri (2009), either: (i) The solution to this min-norm problem is zero, in which case the resulting point satisfies the KKT conditions; or (ii) The solution yields a gradient direction that improves all objectives.

## 2.4 Scalable Multi-Objective Control of LLMs with Surrogate Objectives

However, in the context of LLMs, directly solving this optimization remains computationally intractable because: (i) the need to backpropagate $N + 1$ times to compute the gradient for each objective; and (ii) the expense of solving the min-norm problem in the gradient space for LLM parameters is prohibitively expensive in computation. To overcome this computational burden, we introduce a more computationally efficient surrogate, which is an upper bound to the original objective, circumventing the need for costly backpropagation operations.

**Theorem 1.** *The upper bound of Equation* (8) *is*

$$\left\| c^{(1)} \sum_{j=1}^{N} p_j I(\hat{A}_j) - c^{(2)} \mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) I(\hat{A}_j) \right\|_2^2 \times \left\| \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2, \qquad (9)$$

*where*

$$I(A) = \begin{cases} 0, & \text{if } (A > 0 \text{ and } z > 1 + \epsilon) \\ & \text{or } (A < 0 \text{ and } z < 1 - \epsilon), \\ A, & \text{if } (A > 0 \text{ and } z \le 1 + \epsilon) \\ & \text{or } (A < 0 \text{ and } z \ge 1 - \epsilon), \end{cases} \qquad (10)$$

$$\sum_{i=1}^{2} c^{(i)} = 1, \quad c^{(i)} \ge 0 \quad \forall i; \qquad (11)$$

*where the advantage function $A$, the clip hyper-parameter $\epsilon$, and the ratio $z = \frac{\pi}{\pi_{old}}$ are introduced by the PPO loss (Schulman et al., 2017).*

The proof is deferred to Appendix A. Theorem 1 provides an upper bound on Equation (8), which yields two key advantages: (i) Both $I(\hat{A}_i)$ and $\mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) I(\hat{A}_j)$ can be efficiently computed without any additional expensive back-propagation; (ii) $\nabla_\theta \pi(\cdot; \theta, \mathbf{p})$ is no longer required by the min-norm problem as it does not depend on $c^{(i)}$. Therefore, we achieve the following computationally efficient surrogate problem of optimizing $c^{(1)}$ and $c^{(2)}$:

$$\min_{c^{(1)}, c^{(2)}} \left\{ \left\| c^{(1)} \sum_{j=1}^{N} p_j I(\hat{A}_j) - c^{(2)} \mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) I(\hat{A}_j) \right\|_2^2 \middle| \sum_{i=1}^{2} c^{(i)} = 1, \ c^{(i)} \ge 0, \ \forall i \right\}. \qquad (12)$$

Compared to the intractable original optimization in Equation (8), the surrogate optimization problem in Equation (12) offers the following advantages: (i) Computational efficiency: The term $I(\hat{A}_i)$ can be computed through a simple forward pass in a language model without requiring gradient calculations; (ii) Solution efficiency: Note that the objective function is a quadratic function of the variables $c^{(i)}$. The general min-norm problem is solvable by the existing Frank-Wolfe algorithm (Jaggi, 2013), a well-established convex optimization method. Equation (12) has a closed-form solution (Sener & Koltun, 2018) because Equation (12) only involves two terms.

As a result, the multi-objective learning problem in Equation (7) can be solved by iterating two steps: (i) Solving the min-norm problem in Equation (12) to achieve the dynamic weights $\{c^{(i)}\}_{i=1}^{2}$, and (ii) Optimizing scalarized objective in Equation (7) with the $\{c^{(i)}\}_{i=1}^{2}$. Finally, by integrating PPO's advantage function $A$ into Equation (12), our MOC algorithm can train a policy taking any preference vector to control the multi-objective alignment. This algorithm is summarized in Appendix D.

**Advantages** of MOC. i) Flexible preference handling: MOC can accommodate a wide range of preferences through a single training process. This is achieved through its design, which enables adaptation to diverse preferences without requiring repeated training. ii) Computational Efficiency: By introducing the surrogate objective in Equation (12), MOC significantly reduces computational costs, making its computational cost comparable to that of widely used single-objective RLHF.

## 2.5 An Illustrative Example

To demonstrate the capability of our proposed MOC algorithm, we perform an illustrative experiment on the fishwood task (Felten et al., 2023), where the agent controls a fisherman who can fish or gather wood, receiving corresponding rewards upon completion of each task. The rewards have two dimensions: one for gathering wood and one for fishing. Each collected wood or fished item increases the respective reward by 1. Detailed experimental setup can be found in Appendix G. The results are reported in Figure 1. MOC aims for (i) multi-objective optimization: The solutions should reach the Pareto front, meaning the points should
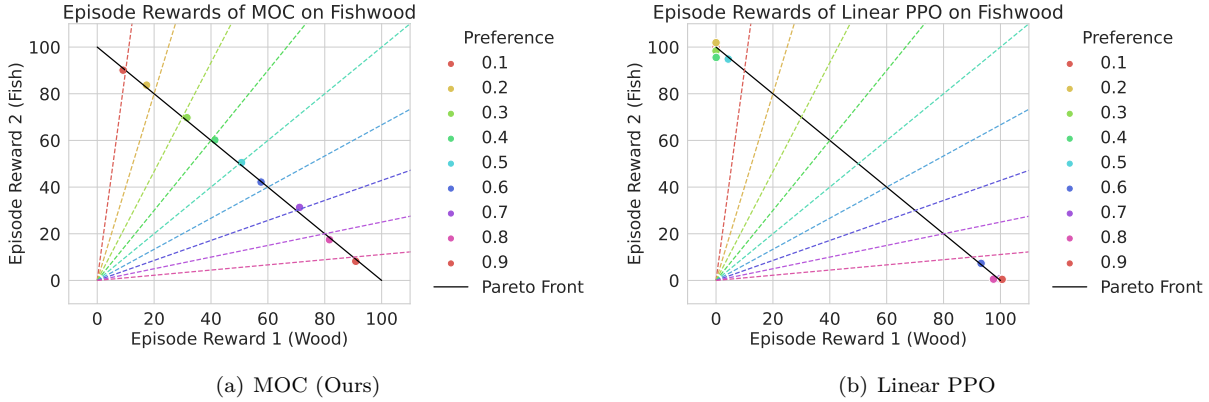
Figure 1: Solutions of MOC and Linear PPO on fishwood task and the Pareto front (line in black). MOC demonstrates advantages in both multi-objective optimization (solutions lie on the Pareto front) and multi-objective control (its solutions align closely with their corresponding preference vectors, shown as the colored dashed rays). The single model trained by MOC can handle diverse preference vectors. In contrast, Linear PPO optimizes a linear scalarization of the objectives and fails to follow the preference vectors, with solutions dominated by one objective. The examined preference weights of "episode reward 1 (wood)" are listed below "Preference".

be close to the black solid line. (ii) Multi-objective control: The solutions should align closely with their respective preference vectors, indicated by the dashed lines. The results demonstrate that MOC achieves both goals: (i) Its solutions lie on the Pareto front, demonstrating successful optimization, and (ii) its solutions are close to the preference vectors, confirming effective multi-objective control. Notably, MOC generalizes to diverse preference vectors by training only **one** model. In contrast, the Linear PPO method, which optimizes a linear scalarization of the objectives, struggles to consistently follow different preference vectors. In Linear PPO' results, one objective often dominates the other in the Pareto sense, a well-known phenomenon in convex optimization (see section 4.7 of Boyd & Vandenberghe (2004)).

## 3 Experiments

In this section, we present a comprehensive evaluation of the proposed MOC method.

### 3.1 Experimental Setup

**Implementation.** Our implementation is based on the open-source TRL package (von Werra et al., 2020). For the language model, we adopt models from the Llama series (Touvron et al., 2023; Dubey et al., 2024), which are widely used in RLHF studies. We use the prompts from Helpful Assistant dataset (Bai et al., 2022), which provides data for two sets of objectives: {"humor", "helpful"} and {"harmless", "helpful"}. MOC is trained with a set of predefined preference vectors: $\{[0.0, 1], [0.1, 0.9], \cdots, [0.9, 0.1], [1, 0.0]\}$. The training process was performed on a desktop with an Intel i9-14900K CPU and an NVIDIA RTX A6000 GPU. MOC is trained by LoRA (Hu et al., 2022) with a rank of 64. The language model is loaded in 8-bit due to the computational constraints. Additional experimental details are provided in Appendix H.

**Baselines.** We compare MOC against three baselines, including (i) The standard *MORLHF*: a multi-objective RLHF method that scalarizes the multi-objective problem into a single objective by combining reward signals with fixed preference weights; (ii) *Rewarded Soups* (Ramé et al., 2023): Combines the model weights of $N$ separately trained models using the PPO algorithm, each optimized for a specific reward function; (iii) *RiC* (Yang et al., 2024b): This method conditions the response of the language model on multiple rewards via prompt conditioning. The behavior of the base Llama-2 model is included for comparative analysis. More baselines are included in Appendix K.
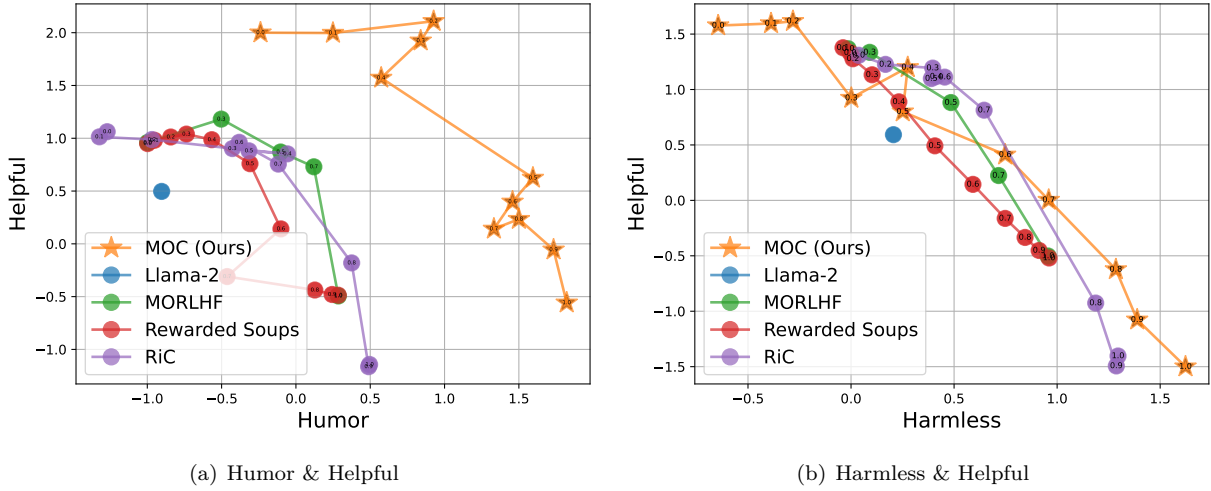
(a) Humor & Helpful

(b) Harmless & Helpful

Figure 2: **Controllability comparison on the Pareto front**. MOC demonstrates superior controllability, indicated by the consistent ranking of solutions on their preference weights and the achieved reward values. In comparison, the baselines exhibit less stable behavior and weaker alignment with the specified preferences. MOC also achieves higher quality solutions, particularly in the Humor & Helpful alignment. Our MOC method achieves the best overall performance, supported by these results and the findings in Tables 2 to 4. Each point represents the reward achieved across multiple instances, each with a different input preference vector. Each point's preference weight for the x-axis reward is the numerical label on its marker.

**Metrics.** It is crucial to emphasize that in the context of controllability, higher rewards do **NOT** always equate to better outcomes. Our goal is to ensure that the model's outputs align with the personalized expectations of users, rather than merely maximizing reward scores. Therefore, to evaluate the controllability comprehensively, we employ a multi-faceted approach.

The evaluation focuses on four key aspects, namely: (i) The quality of solutions, measured using hyper-volumes (as illustrated in Figure 3); (ii) Control with preference vectors, assessed by computing the correlation between the model's behavior and the given preferences; (iii) Diversity of solutions, evaluated by computing the mean pairwise distance of the solutions; and (iv) Generalization capabilities to unseen preference vectors. These metrics collectively provide a robust principle to evaluate the model's ability to optimize competing objectives while adhering to user-defined preferences. Additionally, we present case studies to provide qualitative insights into the controllability of MOC with user-specified preferences.

### 3.2 Main Results

Figure 2 illustrates the results for two pairs of reward models, with coordinates representing the average rewards corresponding to different preference vectors.
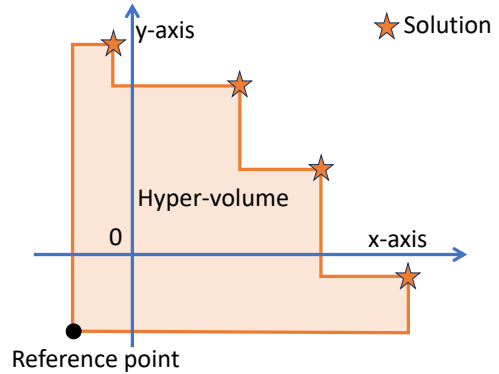


Figure 3: Illustration of the hyper-volume concept. The hyper-volume measures the size of the objective space dominated by a set of solutions in multi-objective optimization. Larger hyper-volumes indicate better convergence and diversity of the Pareto front.

The results indicate two key conclusions: (i) Controllability: MOC demonstrates superior controllability compared to the baselines. This is evident in how consistently the model's behavior aligns with the rank order prescribed by the preference vectors, maintaining a clear monotonic relationship between given preferences and corresponding rewards. In contrast, MORLHF, Rewarded Soups, and RiC exhibit less stable and less consistent behavior relative to their corresponding preferences;

(ii) Solution quality: MOC outperforms all baselines in terms of solution quality. These findings are further supported by following quantitative results.

Table 2: Controllability comparison using Kendall's tau correlation (higher is better), with corresponding $p$-values in parentheses (smaller is better), measuring the consistency between input preferences and output rewards. MOC significantly outperforms all the baselines.

| Dataset | MOC (Ours) | RiC | MORLHF | Rewarded Soups |
|---|---|---|---|---|
| Humor-helpful | 1.00 $(5.0 \times 10^{-8})$ | 0.78 $(3.3 \times 10^{-4})$ | 1.00 $(1.7 \times 10^{-2})$ | 1.00 $(5.0 \times 10^{-8})$ |
| Harmless-helpful | 1.00 $(5.0 \times 10^{-8})$ | 0.91 $(3.0 \times 10^{-5})$ | 1.00 $(1.7 \times 10^{-2})$ | 0.96 $(5.5 \times 10^{-7})$ |
| Average | 1.00 | 0.85 | 1.00 | 0.98 |

**Alignment with preferences.** To assess how well different algorithms align with the prescribed preference vectors and model behavior, we adopt *Kendall's tau* rank correlation (Kendall, 1938) as the evaluation metric. Kendall's tau quantifies the degree to which the relative ordering of outputs is consistent with the rank order determined by the preference vectors, capturing the model's ability to maintain rank-preserving relationships. A higher Kendall's tau indicates better alignment between the model's output ordering and the target preference ordering. Check Appendix I for more details. The results are shown in Table 2. Across two evaluation settings, MOC achieves the highest Kendall's tau, demonstrating its superior capability to align model behavior with user preferences and accurately reflect human preference rankings.

**Quality of solutions.** We use the hyper-volume indicator, a standard metric in MOO, to measure the quality of solutions. Hyper-volume captures both convergence to the Pareto front and the diversity of the solutions in the objective space. Table 3 shows that MOC significantly outperforms all baselines. For instance, in the Humor-Helpful setting, MOC achieves a hyper-volume of 14.176, compared to 6.692 by RiC; similar trends are observed in the Harmless-Helpful setting. These results indicate that MOC exhibits superior convergence to the Pareto front and maintains a more diverse set of solutions, ensuring that it explores a broader range of trade-offs between objectives.

Table 3: Hyper-volume (higher is better) comparison of different methods, measuring the volume of solutions dominated by each method's solution set. MOC outperforms all baselines, achieving higher solution diversity and quality. The best score is marked with the blue color box.

| Setting | MOC (Ours) | RiC | MORLHF | Rewarded Soups |
|---|---|---|---|---|
| Humor-helpful | 14.176 | 6.692 | 6.769 | 6.100 |
| Harmless-helpful | 10.220 | 9.257 | 9.047 | 8.905 |
| Average | 12.198 | 7.974 | 7.908 | 7.502 |

Table 4: Comparison of pair-wise distance of the solution set (measuring diversity, higher is better) of different methods.

| Dataset | MOC (Ours) | RiC | MORLHF | Rewarded Soups |
|---|---|---|---|---|
| Humor-helpful | 1.439 | 1.260 | 1.057 | 1.005 |
| Harmless-helpful | 1.600 | 1.363 | 1.110 | 1.015 |
| Average | 1.520 | 1.312 | 1.084 | 1.010 |

**Diversity of solutions.** We measure the diversity of solutions by computing the mean pairwise distance (MPD) of the solution set in the objective space. MPD directly quantifies how spread out the solutions are, making it particularly suitable for our setting. A higher MPD indicates greater behavioral diversity. Table 4 shows that MOC consistently achieves the highest MPD values, outperforming all baselines. For example, in the Harmless-Helpful setting, MOC obtains an MPD value of 1.600, while RiC obtains 1.363. This result

aligns with the observation in Figure 2, where the solutions produced by RiC tend to cluster more closely, resulting in less diverse behavior.
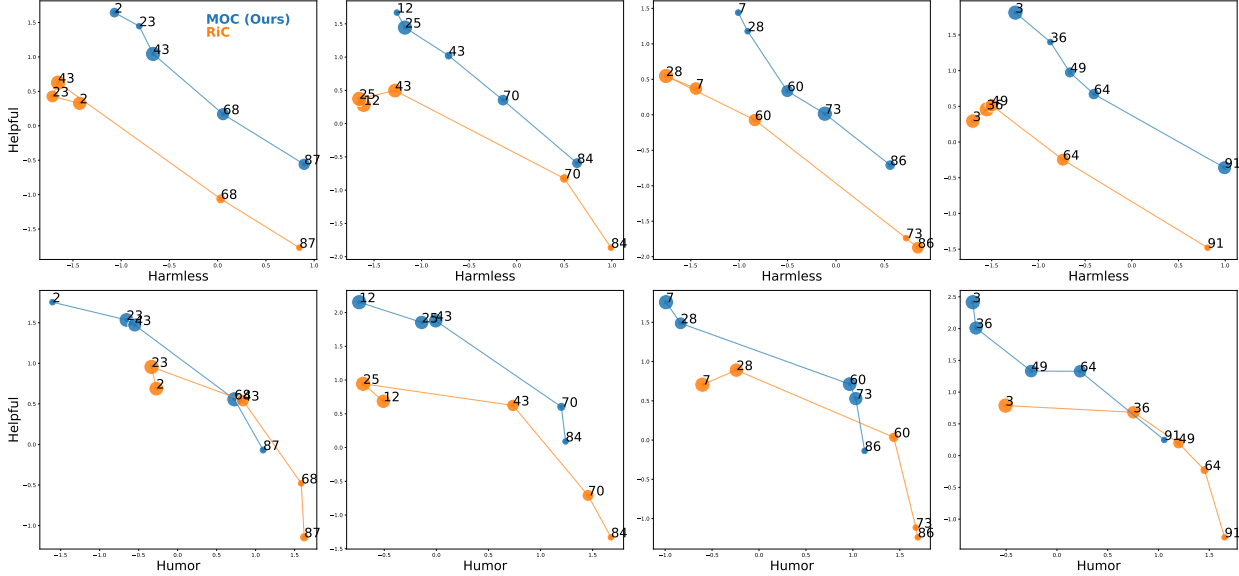


Figure 4: Generalization to unseen preference vectors held out from the training. MOC and RiC-trained LLMs are compared on four random sets of unseen preference vectors. Each column corresponds to a different set of unseen preference vectors, and each row represents a different pair of reward settings. MOC solutions dominate the RiC solutions in most cases. MOC's rewards align with the new preference vectors and the outputs under different preferences are diverse in the reward space. This suggests MOC generalizes to unseen preferences and achieves diverse trade-offs on the Pareto front. The size of each point indicates the standard deviation in rewards. The numerical labels indicate the preference weights (multiplied by 100) for the reward on the x-axis, enhancing visual clarity.

### 3.3 Generalization to Unseen User Preference

We evaluate the ability of our model to generalize to unseen preference vectors not included in the training process. Although MOC is initially trained on a predefined set of preference vectors, the goal is to determine if it can handle new, untrained preferences effectively. To test this, we uniformly sampled four sets of unseen preference vectors and provided them as inputs to the trained model for inference. The results, as depicted in Figure 4, confirm that the model maintains strong performance across all tested scenarios, without any obvious degradation in its behavior. Check Appendix J for more details and qualitative results.

The results highlight several key advantages of the model trained by MOC: i) The model's performance does not degrade when presented with unseen preferences. ii) The model's behavior still adheres to the input preference vector, maintaining alignment between behavior ranking (represented by rewards) and preferences. iii) The model demonstrates sufficient diversity in its behavior, distributing its rewards across a broad range of outcomes rather than concentrating on a narrow region of the objective space. These results suggest that the MOC can successfully accommodate a diverse range of trade-offs dictated by new preferences, even when they significantly differ from those encountered during training.

### 3.4 Additional Experiments

To further demonstrate MOC's capabilities, we present a series of experiments highlighting its advantages across various settings. These include: (1) its ability to generalize to untrained preference vectors (Appendix J), (2) its performance across different model types and sizes, compared with more baselines (Appendix K), (3)

its adaptability to various datasets and reward models (Appendix L), (4) its effectiveness in handling more objectives (Appendix M), and (5) a case study (Appendix N).

### 3.5 Discussion

The experiments reveal four key advantages of MOC. (i) Solution Quality: MOC achieves the highest solution quality, evidenced by the hyper-volume, reflecting both convergence and diversity. (ii) Controllability: MOC demonstrates superior controllability, ensuring consistent alignment with user preferences across diverse objective trade-offs. (iii) Solution Diversity: MOC outperforms baselines, confirming its robustness in capturing user preferences. (iv) Generalization: MOC's ability to generalize to unseen preferences highlights its potential for real-world applications where new preferences may emerge. These advantages demonstrate that MOC offers a powerful and flexible approach for multi-objective controllable language models, outperforming existing methods in controllability and diversity while maintaining computational efficiency.

## 4 Related Work

The alignment of LLMs with human values is a central challenge (Ouyang et al., 2022; Bai et al., 2022), and moving beyond a single, developer-defined preference to accommodate diverse, multi-objective user needs is a critical frontier. Our work, MOC, introduces a principled yet efficient framework for this task. We position our contribution by first discussing the evolution from foundational MOO to practical applications, and then by situating MOC within the current landscape of multi-objective control methods for LLMs.

### 4.1 From Multi-Objective Optimization to Controllable LLMs

The principles of MOO, which aim to find a set of Pareto-optimal solutions for competing objectives, are well-established (Désidéri, 2009; Sener & Koltun, 2018). However, applying these principles to large-scale LLMs presents unique challenges. Early MOO work in machine learning often focused on optimizing multiple loss functions simultaneously (Liu et al., 2021; 2023), a paradigm that does not directly map to controlling the fine-grained, semantic behaviors of LLMs.

Closer to our domain, some methods rely on a linear scalarization of utility or reward functions (Yang et al., 2019), a technique also employed by standard MORLHF. While simple, this approach is theoretically limited and often fails to identify solutions in non-convex regions of the Pareto front (Boyd & Vandenberghe, 2004). Other algorithms have been developed to trace the full Pareto front (Mahapatra & Rajan, 2021; Zhang et al., 2024), but their high computational complexity makes them intractable for finetuning billion-parameter language models.

MOC bridges this critical gap. Our work differs from these foundational methods in two key respects: (i) we directly manipulate model behavior in the more meaningful reward space rather than loss space, and (ii) we introduce a novel surrogate objective (Theorem 1) that makes principled MOO computationally efficient and comparable to standard single-objective RLHF, thus overcoming the scalability bottleneck of prior work.

### 4.2 Multi-Objective Control of LLMs

Recent years have seen a surge of methods specifically for multi-objective LLM control, which can be broadly categorized by their approach.

**Training a Single Steerable Policy.** This paradigm, which MOC belongs to, aims to train one versatile model that can be steered at inference time. The most related works use a numerical vector to specify preferences. Wang et al. (2024) provides a foundational RL framework for this but, like MORLHF, relies on linear reward scalarization. Guo et al. (2024) presents a powerful, RL-free alternative by extending DPO to handle multiple objectives, making it a key intellectual parallel to our work. RiC (Yang et al., 2024b) conditions the model on reward values via prompt engineering and uses rejection sampling. While this achieves a form of control, it lacks an explicit policy improvement mechanism, limiting its ability to push the Pareto frontier outwards. MOC, by integrating its MOO formulation directly into a PPO-based

policy-gradient objective, explicitly optimizes for both alignment and reward maximization. In contrast to these, MOC employs a more sophisticated MOO gradient calculation, enabling a superior exploration of the Pareto front, as demonstrated by our hyper-volume and diversity metrics (Tables 3 and 4).

**Ensemble and Multi-Model Methods.** Other methods achieve diverse outputs by training or combining multiple models. Rewarded Soup (Ramé et al., 2023) interpolates the weights of separate models, each finetuned on a single reward. Similarly, MODPO (Zhou et al., 2024) trains M distinct models for M preferences. These approaches are computationally expensive and storage-intensive, directly contrasting with MOC's "one model for all" design. MOC achieves superior results with the efficiency of training a single model.

### 4.3 Alternative Control Interfaces and Methodologies

**Linguistic and Implicit Control.** Some works explore more user-friendly control interfaces. Nguyen et al. (2024); Yang et al. (2024a) use explicit linguistic tags (e.g., [formality: high]) for control, while MOSLIM (Zhang et al., 2025) infers preferences implicitly from the user's natural language prompt. These methods trade the precise, granular control offered by MOC's numerical vectors for enhanced interpretability or zero-shot prompting. MOC's approach is complementary and particularly suited for applications requiring precise, backend control over model attributes like safety and factuality.

**Inference-Time Control.** A fundamentally different approach is to enforce constraints during decoding. Shi et al. (2024); Son et al. (2025) pioneer methods to provide formal guarantees on objective trade-offs at inference time. This powerful technique can be applied to any LLM but incurs latency and can degrade text quality if the constraints are too severe. MOC, as a training-time method, internalizes these trade-offs, enabling fast, coherent, and controlled generation without specialized decoding, making it ideal for scalable deployment.

MOC carves a unique and powerful algorithm. It is a single, efficient, and steerable model that: (i) does not require training multiple models, unlike Rewarded Soups or MODPO; (ii) does not rely on preference datasets, unlike CPO; (iii) maintains explicit policy improvement, unlike RiC; and (iv) generalizes to unseen preference vectors, which is a primary design goal.

## 5 Conclusion

In this paper, we introduced Multi-Objective Control (MOC), a novel approach that enables the personalization of LLMs by inherently adjusting to diverse user-specified preferences. MOC addresses the limitations of existing LLMs, which are typically constrained by fixed developer-specified preferences, by formulating multi-objective controllability as a multi-objective optimization problem. By introducing surrogate optimization in RLHF, MOC enables a single fine-tuning process to adapt to a wide range of user-specified trade-offs while other methods always leverage expensive multiple models. Our experiments demonstrate that MOC surpasses existing baseline methods in controllability, solution quality, and generalization while maintaining exceptional computational efficiency. By effectively managing trade-offs between objectives and offering a superior Pareto front, MOC is well-suited for real-world applications where flexibility and personalization are critical. This work highlights the potential of MOC to transform how LLMs interact with users, offering scalable and customizable solutions that meet diverse needs while maintaining computational feasibility. Looking forward, MOC paves the way for future research in personalized LLMs. Future work will focus on scaling up the method to larger models. Exploring more complex user preferences and further enhancing scalability will be key to unlocking even broader applications for customizable and efficient LLMs in real-world settings. Ultimately, MOC represents a significant step toward realizing fully personalized and human-friendly systems.

## References

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom

Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR*, abs/2204.05862, 2022. doi: 10.48550/ARXIV.2204.05862. URL https://doi.org/10.48550/arXiv.2204.05862.

Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.

Jean-Antoine Désidéri. 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.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pp. 4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-1423. URL https://doi.org/10.18653/v1/n19-1423.

Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. https://github.com/openai/baselines, 2017.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Florian Felten, Lucas N. Alegre, Ann Nowé, Ana L. C. Bazzan, El Ghazali Talbi, Grégoire Danoy, and Bruno C. da Silva. 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.

Yiju Guo, Ganqu Cui, Lifan Yuan, Ning Ding, Zexu Sun, Bowen Sun, Huimin Chen, Ruobing Xie, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Controllable preference optimization: Toward controllable multi-objective alignment. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pp. 1437–1454. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.EMNLP-MAIN.85. URL https://doi.org/10.18653/v1/2024.emnlp-main.85.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. 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=nZeVKeeFYf9.

Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, volume 28 of *JMLR Workshop and Conference Proceedings*, pp. 427–435. JMLR.org, 2013. URL http://proceedings.mlr.press/v28/jaggi13.html.

M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938. ISSN 00063444. URL http://www.jstor.org/stable/2332226.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015. URL http://arxiv.org/abs/1412.6980.

Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 18878–18890, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/9d27fdf2477ffbff837d73ef7ae23db9-Abstract.html.

Bo Liu, Yihao Feng, Peter Stone, and Qiang Liu. FAMO: fast adaptive multitask optimization. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (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. URL http://papers.nips.cc/paper_files/paper/2023/hash/b2fe1ee8d936ac08dd26f2ff58986c8f-Abstract-Conference.html.

Pingchuan Ma, Tao Du, and Wojciech Matusik. Efficient continuous pareto exploration in multi-task learning. 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. 6522–6531. PMLR, 2020. URL http://proceedings.mlr.press/v119/ma20a.html.

Debabrata Mahapatra and Vaibhav Rajan. Exact pareto optimal search for multi-task learning: Touring the pareto front. *ArXiv*, abs/2108.00597, 2021. URL https://api.semanticscholar.org/CorpusID:236772107.

Dang Nguyen, Jiuhai Chen, and Tianyi Zhou. Multi-objective linguistic control of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *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/.

OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL https://doi.org/10.48550/arXiv.2303.08774.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.

Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018. URL https://api.semanticscholar.org/CorpusID:49313245.

Alexandre Ramé, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (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. URL http://papers.nips.cc/paper_files/paper/2023/hash/e12a3b98b67e8395f639fde4c2b03168-Abstract-Conference.html.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL http://arxiv.org/abs/1707.06347.

Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 525–536, 2018. URL https://proceedings.neurips.cc/paper/2018/hash/432aca3a1e345e339f35a30c8f65edce-Abstract.html.

Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Hannaneh Hajishirzi, Noah A. Smith, and Simon S. Du. Decoding-time language model alignment with multiple objectives. *CoRR*, abs/2406.18853, 2024. doi: 10.48550/ARXIV.2406.18853. URL https://doi.org/10.48550/arXiv.2406.18853.

Seongho Son, William Bankes, Sangwoong Yoon, Shyam Sundhar Ramesh, Xiaohang Tang, and Ilija Bogunovic. Robust multi-objective controlled decoding of large language models. *CoRR*, abs/2503.08796, 2025. doi: 10.48550/ARXIV.2503.08796. URL https://doi.org/10.48550/arXiv.2503.08796.

Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. Learning to summarize with human feedback. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl, 2020.

Kaiwen Wang, Rahul Kidambi, Ryan Sullivan, Alekh Agarwal, Christoph Dann, Andrea Michi, Marco Gelmi, Yunxuan Li, Raghav Gupta, Kumar Avinava Dubey, Alexandre Rame, Johan Ferret, Geoffrey Cideron, Le Hou, Hongkun Yu, Amr Ahmed, Aranyak Mehta, Leonard Hussenot, Olivier Bachem, and Edouard Leurent. Conditional language policy: A general framework for steerable multi-objective finetuning. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 2153–2186, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.118. URL https://aclanthology.org/2024.findings-emnlp.118/.

Peiyao Xiao, Hao Ban, and Kaiyi Ji. Direction-oriented multi-objective learning: Simple and provable stochastic algorithms. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (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. URL http://papers.nips.cc/paper_files/paper/2023/hash/0e5b96f97c1813bb75f6c28532c2ecc7-Abstract-Conference.html.

Jie Xu, Yunsheng Tian, Pingchuan Ma, Daniela Rus, Shinjiro Sueda, and Wojciech Matusik. Prediction-guided multi-objective reinforcement learning for continuous robot control. 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. 10607–10616. PMLR, 2020. URL http://proceedings.mlr.press/v119/xu20h.html.

Kailai Yang, Zhiwei Liu, Qianqian Xie, Jimin Huang, Tianlin Zhang, and Sophia Ananiadou. Metaaligner: Towards generalizable multi-objective alignment of language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.

Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024b. URL https://openreview.net/forum?id=QLcBzRI3V3.

Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective reinforcement learning and policy adaptation. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (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. URL https://proceedings.neurips.cc/paper/2019/hash/4a46fbfca3f1465a27b210f4bdfe6ab3-Abstract.html.

Yijun Yang, Jing Jiang, Tianyi Zhou, Jie Ma, and Yuhui Shi. Pareto policy pool for model-based offline 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=OqcZu8JIIzS.

Xiaoyuan Zhang, Xi Lin, and Qingfu Zhang. PMGDA: A preference-based multiple gradient descent algorithm. *CoRR*, abs/2402.09492, 2024. doi: 10.48550/ARXIV.2402.09492. URL https://doi.org/10.48550/arXiv.2402.09492.

Yu Zhang, Wanli Jiang, and Zhengyu Yang. Moslim:align with diverse preferences in prompts through reward classification. *CoRR*, abs/2505.20336, 2025. doi: 10.48550/ARXIV.2505.20336. URL https://doi.org/10.48550/arXiv.2505.20336.

Zhanhui Zhou, Jie Liu, Jing Shao, Xiangyu Yue, Chao Yang, Wanli Ouyang, and Yu Qiao. Beyond one-preference-fits-all alignment: Multi-objective direct preference optimization. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 10586–10613, 2024.

# A   Proof of Theorem 1

**Theorem 1.** *The upper bound of Equation* (8) *is*

$$\left\| c^{(1)} \sum_{j=1}^{N} p_j I(\hat{A}_j) - c^{(2)} \mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) I(\hat{A}_j) \right\|_2^2 \times \left\| \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2, \tag{9}$$

*where*

$$I(A) = \begin{cases} 0, & \text{if } (A > 0 \text{ and } z > 1 + \epsilon) \\ & \quad \text{or } (A < 0 \text{ and } z < 1 - \epsilon), \\ A, & \text{if } (A > 0 \text{ and } z \leq 1 + \epsilon) \\ & \quad \text{or } (A < 0 \text{ and } z \geq 1 - \epsilon), \end{cases} \tag{10}$$

$$\sum_{i=1}^{2} c^{(i)} = 1, \quad c^{(i)} \geq 0 \quad \forall i; \tag{11}$$

*where the advantage function $A$, the clip hyper-parameter $\epsilon$, and the ratio $z = \frac{\pi}{\pi_{old}}$ are introduced by the PPO loss (Schulman et al., 2017).*

To tackle the intractable computation of Equation (8), we introduce the following surrogate optimization objective.

*Proof.* One can further expand Equation (8) with the PPO loss and get

$$\left\| c^{(1)} \mathbf{p}^\top \nabla_\theta \mathbf{J}(\pi(\cdot; \theta, \mathbf{p})) - c^{(2)} \nabla_\theta ReLU(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi) \right\|_2^2$$

$$= \left\| c^{(1)} \sum_{j=1}^{N} p_j \nabla_\theta J^j(\pi(\cdot; \theta, \mathbf{p})) - c^{(2)} \nabla_\theta ReLU(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi) \right\|_2^2$$

$$= \left\| c^{(1)} \sum_{j=1}^{N} p_j \nabla_\pi J^j(\pi(\cdot; \theta, \mathbf{p})) \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) - c^{(2)} \nabla_\pi ReLU(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi) \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2 \tag{13}$$

$$\leq \left\| c^{(1)} \sum_{j=1}^{N} p_j \nabla_\pi J^j(\pi(\cdot; \theta, \mathbf{p})) - c^{(2)} \nabla_\pi ReLU(MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi) \right\|_2^2 \left\| \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2$$

$$= \left\| c^{(1)} \sum_{j=1}^{N} p_j \frac{1}{\pi_{old}} I(\hat{A}_j) - c^{(2)} \mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) \frac{1}{\pi_{old}} I(\hat{A}_j) \right\|_2^2 \left\| \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2$$

$$\leq \left\| c^{(1)} \sum_{j=1}^{N} p_j I(\hat{A}_j) - c^{(2)} \mathbf{1}_{MSE(\mathbb{E}_{x \sim \mathcal{D}} \mathbf{R}(x,y), \mathbf{p}) - \phi > 0} \sum_{j=1}^{N} (R^j - p_j) I(\hat{A}_j) \right\|_2^2 \left\| \nabla_\theta \pi(\cdot; \theta, \mathbf{p}) \right\|_2^2$$

*where*

$$I(A) = \begin{cases} 0, & \text{if } (A > 0 \text{ and } z > (1 + \epsilon)) \\ & \quad \text{or } (A < 0 \text{ and } z < 1 - \epsilon) \\ A, & \text{if } (A > 0 \text{ and } z \leq (1 + \epsilon)) \\ & \quad \text{or } (A < 0 \text{ and } z \geq 1 - \epsilon) \end{cases}, \tag{14}$$

$$\sum_{i=1}^{2} c^{(i)} = 1, \quad c^{(i)} \geq 0 \quad \forall i, \tag{15}$$

and $z = \frac{\pi}{\pi_{old}}$. The third inequality holds by Cauchy–Schwarz inequality and the fourth equation holds by integrating the PPO loss function. □

# B Formal Definitions and Advantages of MOC in Multi-Objective Optimization

In this section, we provide a formal definition of Pareto Optimality and its relevance to policy improvement.

## B.1 Formal Definition of Pareto Optimality

**Definition 1.** Let $\pi, \pi' \in \mathcal{X}$, where $\mathcal{X}$ is the set of feasible solutions. A solution $\pi$ is said to *dominate* another solution $\pi'$ if and only if:

- $J_i(\pi) \geq J_i(\pi')$ for all $i \in \{1, 2, \ldots, N\}$, and

- $J_j(\pi) > J_j(\pi')$ for at least one $j \in \{1, 2, \ldots, N\}$.

Here, $J_i(\pi)$ denotes the value of the $i$-th objective for the solution $\pi$. The above conditions imply that $\pi$ performs at least as well as $\pi'$ in all objectives and strictly better in at least one. Solutions that are not dominated by any other are termed *non-dominated* and collectively form the *Pareto front*.

**Definition 2.** (Pareto Optimality) Let $\mathcal{X}$ denote the set of feasible solutions, and let $J : \mathcal{X} \to \mathbb{R}^N$ be a vector-valued objective function where $J(\pi) = [J_1(\pi), J_2(\pi), \ldots, J_N(\pi)]^\top$ corresponds to the objective values associated with $\pi \in \mathcal{X}$. A solution $\pi^* \in \mathcal{X}$ is *Pareto optimal* if and only if no other solution $\pi' \in \mathcal{X}$ satisfies:

$$J_i(\pi') \geq J_i(\pi^*) \quad \forall i \in \{1, 2, \ldots, N\} \tag{16}$$

and

$$J_j(\pi') > J_j(\pi^*) \quad \text{for at least one } j \in \{1, 2, \ldots, N\}. \tag{17}$$

This ensures that $\pi^*$ is *non-dominated*, meaning that no other solution can improve one or more objectives without sacrificing performance in at least one other.

## B.2 Advantage of MOC

Explicit policy improvement refers to methods that deliberately optimize at least one objective $J_i$, ensuring that the solution quality improves by maximizing one or more associated rewards $R_i$. This approach is particularly crucial in designing multi-objective policies, as it guarantees measurable progress in one or more dimensions of performance.

**Advantage of MOC Compared to Other Baselines**

Our proposed method, **MOC**, explicitly optimizes all objectives with policy improvement while integrating controllability, ensuring a more balanced and efficient approach to policy improvement. In contrast:

- **Rewarded Soup** does not jointly optimize all objectives, which leads to suboptimal solutions.

- **RiC** focuses exclusively on controllability but lacks explicit mechanisms for policy improvement, limiting its ability to enhance solution quality.

- **MODPO** does not consider Pareto Optimality during training. Specifically, it trains $M$ separate LLMs (corresponding to $M$ preferences) by optimizing each model with a specific weighted combination of reward objectives, given the corresponding reward models.

By integrating both explicit policy improvement and controllability into a unified framework, **MOC** theoretically achieves higher solution quality compared to these baselines. This is further validated by our experimental results (Tables 1 to 4 and 10 to 12 and Figures 2, 4, 6 and 7), which demonstrate that **MOC** consistently outperforms these approaches across multiple metrics.

The integration of explicit policy improvement with controllability ensures that **MOC** aligns with the principles of Pareto Optimality while delivering superior practical performance. By addressing the limitations

of existing methods and achieving a better balance among competing objectives, **MOC** sets a new benchmark in multi-objective controllable language models.

## C   Approximated Normalized Vector Similarity

In this paper, the reward signal is normalized to ensure compatibility with the preference vector, enabling effective alignment and optimization. The normalization process is defined as:

$$Normalize(r) = \frac{r - r_{\text{mean}}}{2r_{\text{std}}} + 1, \tag{18}$$

where $r_{\text{mean}}$ and $r_{\text{std}}$ are computed dynamically using a running mean and standard deviation (Dhariwal et al., 2017). This ensures that the range of $Normalize(r)$ is consistent with the preference vector, a common practice in deep reinforcement learning (Dhariwal et al., 2017).

The alignment between normalized rewards and preferences is then quantified using the Mean Squared Error (MSE) loss, leading to the definition of the **Approximated Normalized Vector Similarity** (AMVS):

$$AMVS(r, \mathbf{p}) = \|Normalize(r) - \mathbf{p}\|^2, \tag{19}$$

which serves as a computationally efficient approximation of the **Normalized Vector Difference** (NVD), a widely adopted similarity measure in multi-objective optimization. The NVD itself is formally defined as:

$$NVD(\mathbf{a}, \mathbf{b}) = \left\| \frac{\mathbf{a}}{\|\mathbf{a}\|} - \frac{\mathbf{b}}{\|\mathbf{b}\|} \right\|. \tag{20}$$

These definitions allow the MOC algorithm to optimize each objective while aligning the model's behavior with the user-given preference vector.

# D Pseudocode

We summarize the MOC algorithm in Algorithm 1. We recommend that the reader checks Schulman et al. (2017); von Werra et al. (2020) for more training details of PPO in the language model settings. The min-norm algorithm used in MOC is shown in Algorithm 2, based on Sener & Koltun (2018). Algorithm 2 gives a $c^{(1)}$ and $c^{(2)} = 1 - c^{(1)}$.

---

**Algorithm 1** <u>M</u>ulti <u>O</u>bjective <u>C</u>ontrol Algorithm (MOC) for Language Models

---

**Require:**

$\quad$ $\mathbb{P} = \{\mathbf{p^i}\}_{i=1}^M$: Preference vector set

$\quad$ $\phi$: Constraint threshold

$\quad$ $\mathcal{D}$: Prompt dataset

$\quad$ The SFT policy $\pi(\cdot; \theta)$ with parameters $\theta$

$\quad$ Add $N$ new value heads to the language model

$\quad$ Set number of iterations $T$ and mini-batch size $B$

1: **for** iteration $t = 1$ to $T$ **do**
2: $\quad$ Sample a mini-batch of prompts from $\mathcal{D}$.
3: $\quad$ Sample a mini-batch of preference vectors $\{\mathbf{p}_j\}_{j=1}^B$.
4: $\quad$ Relabel the prompts with $\{\mathbf{p}_j\}_{j=1}^B$ by Equation (3) and get $\{x_j\}_{j=1}^B$.
5: $\quad$ For each $x_j$, generate response $y_j \sim \pi(x_j; \theta, \mathbf{p}_j)$.
6: $\quad$ Compute $\mathbf{R}(x_j, y_j) = (R^1(x_j, y_j), R^2(x_j, y_j), \dots, R^N(x_j, y_j))$ by reward models.
7: $\quad$ Compute the Advantage function $\hat{A}_j$ according to the PPO algorithm.
8: $\quad$ Solve Equation (12) by Algorithm 2 and get $\{(c_j^{(1)}, c_j^{(2)})\}_{j=1}^B$.
9: $\quad$ Perform gradient ascending using Equation (7) to optimize the policy.
10: $\quad$ Optimizing the $N$ value function of PPO (Schulman et al., 2017).
11: **end for**
12: **return** Optimized policy $\pi$.

---

---

**Algorithm 2** Min-norm algorithm for two vectors ($\min_{c \in [0,1]} \|c\mathbf{v} + (1 - c)\overline{\mathbf{v}}\|_2^2$)

---

**Require:**

$\quad$ $\mathbf{v}$: Vector $\mathbf{v}$

$\quad$ $\overline{\mathbf{v}}$: Vector $\overline{\mathbf{v}}$

1: **if** $\mathbf{v}^\top \overline{\mathbf{v}} \geq \mathbf{v}^\top \mathbf{v}$ **then**
2: $\quad$ $c = 1$
3: **else if** $\mathbf{v}^\top \overline{\mathbf{v}} \geq \overline{\mathbf{v}}^\top \overline{\mathbf{v}}$ **then**
4: $\quad$ $c = 0$
5: **else**
6: $\quad$ $c = \frac{(\overline{\mathbf{v}} - \mathbf{v})^\top \overline{\mathbf{v}}}{\|\mathbf{v} - \overline{\mathbf{v}}\|_2^2}$
7: **end if**
8: **return** $c$

---

# E   Loss Functions of RL Cannot Be Used for Alignment or Control with Preferences

The primary objective in RL is to train an agent to make decisions that maximize cumulative rewards over time To achieve this, various learning algorithms are employed, each associated with specific loss functions. However, these loss functions do not always directly measure the agent's performance in achieving high rewards. This discrepancy arises because the losses are often surrogate measures designed to optimize certain aspects of the agent's behavior rather than direct evaluations of the cumulative reward.

## E.1   Value Function Loss

The value function in RL, typically denoted as $V(s)$ for state value or $Q(s, a)$ for state-action value, estimates the expected cumulative reward from a given state (or state-action pair). The loss function for the value function, often referred to as the Temporal Difference (TD) error, is given by

$$L_V = \mathbb{E}_\pi \left[ (R_t + \gamma V(S_{t+1}) - V(S_t))^2 \right], \tag{21}$$

where

- $R_t$ is the reward received at time step $t$,
- $\gamma$ is the discount factor,
- $V(S_t)$ is the estimated value of the current state,
- $V(S_{t+1})$ is the estimated value of the next state.

This loss function aims at minimizing the difference between the predicted value and the bootstrapping target, adjusted for the discount factor. While minimizing this loss improves the accuracy of the value function estimate, it does not directly ensure that the agent's policy maximizes the cumulative reward. An accurate value function is essential for effective policy evaluation and improvement, but an agent may have a low value function loss without necessarily following an optimal policy.

## E.2   Policy Gradient Loss

Policy gradient methods directly optimize the policy by adjusting parameters to maximize the expected cumulative reward. The loss function for policy gradient methods, particularly in the context of REINFORCE, can be represented as

$$L_\pi = -\mathbb{E}_\pi \left[ \sum_{t=0}^{T} \log \pi_\theta(A_t|S_t) \cdot \hat{A}_t \right], \tag{22}$$

where

- $\pi_\theta(A_t|S_t)$ is the probability of taking action $A_t$ in state $S_t$ under the policy $\pi$ parameterized by $\theta$,
- $\hat{A}_t$ is the advantage function.

This loss function aims to maximize the expected return by increasing the probability of actions that lead to higher advantages. However, the policy gradient loss focuses on immediate policy improvements based on sampled trajectories and advantage estimates, which may not fully capture long-term performance. Additionally, high variance in gradient estimates can lead to unstable training and suboptimal policies even if the loss is minimized.

## E.3   Case of using value function as aligned target

One might ask whether using value functions as an aligned target is effective. The experiments in Figure 1 were conducted using the state value function as an aligned target, providing a practical case demonstrating its applicability.

### E.4 Discussion

Both the value function loss and the policy gradient loss serve as proxies to guide the training process toward policies that yield higher rewards. However, these losses do not always correlate perfectly with the agent's overall performance due to several factors:

- **Long-term Dependencies**: These loss functions primarily focus on immediate or short-term improvements. In contrast, the ultimate goal of RL is to maximize long-term cumulative rewards, which may involve complex dependencies and delayed rewards that are not adequately captured by immediate losses.

- **Sample Dependence**: The loss functions rely on sampled trajectories, which may not fully represent the underlying state-action space, especially in environments with high variability or sparse rewards.

- **Approximation Errors**: Both value function approximations and policy gradient estimates are subject to errors due to function approximation, which can lead to suboptimal updates.

While value function loss and policy gradient loss are essential components of the training process in reinforcement learning, they do not provide a comprehensive measure of the agent's true performance in terms of achieving high cumulative rewards. Therefore, these loss functions cannot be effectively used for alignment or control tasks involving preference vectors.

## F    Further Discussion of Related Work

In this section, we provide further discussion of related work.

While both MOC and RiC (Yang et al., 2024b) aim to personalize LLMs, **their methodologies are fundamentally different**. The following significant differences highlight the novelty and distinct contributions of MOC:

- **Re-labeling the prompt**
  - MOC directly relabels prompts with preference weights while RiC relabels the prompts with reward signal.

- **Formulation and Approach:**
  - MOC formulates controllability as a multi-objective policy optimization with preference-based constraints, solved via the proposed MOC algorithm. RiC uses SFT to fine-tune the LLMs.
  - A key novelty is the surrogate problem (Equation (12)), which reduces the computational cost to near that of a single-objective PPO.
  - In contrast, RiC relies on learning a preference-to-reward mapping and lacks an explicit policy optimization framework, making its methodology fundamentally different.

- **Explicit Policy Optimization and Controllability:**
  - MOC explicitly optimizes the policy to align model behavior with user preferences, establishing a rigorous and systematic controllability framework.
  - RiC does not perform explicit policy optimization, limiting its ability to maximize reward while aligning with preferences.

- **Performance Advantages:**
  - Thanks to its principled design, MOC significantly outperforms RiC in controllability, solution quality, diversity, and generalization, as substantiated by quantitative results in Tables 2–6.

We summarize these key differences in Table 5 for clarity.

Table 5: Key differences between MOC and RiC

| Algorithm | Source of Controllability | Explicit Policy Improvement? | Loss Function |
|---|---|---|---|
| MOC | Multi-objective optimization with constraints | ✓ | PPO-based |
| RiC | Reward in Context | × | SFT-based |

It is important to emphasize that MOC focuses on controllability with two critical goals:

- Training one LLM to generate personalized outputs for diverse user preferences.

- Generalizing to unseen preferences with the once-trained LLM.

In contrast, MODPO:

- Trains M separate LLMs (where M corresponds to the number of preferences), targeting only a fixed set of pre-defined preferences.

- Does not consider generalization to unseen preferences.

These differences highlight the fundamental distinctions between MODPO's methodology (Zhou et al., 2024) and objectives and those of MOC. Therefore, we train an language model with MODPO and include the MODPO results in Table 10.

## G   Details of Toy Example

Readers can click this link: https://mo-gymnasium.farama.org/environments/fishwood/ for more details about the task in Figure 1. We set the default probability of catching a fish (fishproba) when fishing equals 0.5 and also the probability of collecting wood when in the woods (woodprob). The Pareto front is computable once fishproba and woodprob are given. Specifically, the Pareto front satisfies the following equation:

$$x + y = \text{woodprob * (steps collecting wood) + fishprob * (steps fishing)}, \tag{23}$$

where $x$ is the episode reward of fish and y is the episode reward of wood. Specifically, $x + y = 100$ in our settings. The episodes reward are estimated over 20 episodes. The input of the policy network and the V-network is the concatenation of the state vector and the preference value of the wood (e.g. [initial state vector, 0.1]). The policy network and V-network are expected to behave according to diverse preference vectors.

**Selection of preference vector.** The preferences of wood range from 0.1 to 0.9. The following equation gives how we depict the preference vectors.

$$y = \frac{1 - \text{preference\_of\_wood}}{\text{preference\_of\_wood}} * x,$$

where preference_of_wood $\in (0, 1]$ represents the relative preference for collecting wood.

We list the hyper-parameters related to this experiment in Table 6.

Table 6: Hyper-parameters settings for fishwood task (Section 2.5).

| Hyper-parameter | Value |
|---|:---:|
| Dimension of state space | 1 |
| Action space | Discrete(2): go fishing, go collect wood |
| Discount $(\gamma)$ | 0.99 |
| Optimizer | Adam (Kingma & Ba, 2015) |
| Learning rate for networks | $1 \times 10^{-4}$ |
| Number of hidden layers for all networks | 3 |
| Number of hidden units per layer | 256 |
| Activation function | ReLU |
| Batch size | 512 |
| Gradient clipping | False |
| Exploration method | Epsilon-Greedy |
| $\epsilon$ (Exploration) | 0.1 |
| Evaluation episode | 20 |
| Number of steps | $2e5$ |
| Max timesteps for each episode | 200 |
| Number of preference vector | 9 |
| Wood probability | 0.5 |
| Fish probability | 0.5 |

# H  Details of Language Model Experiments

The key information about the experimental settings is listed in Table 7. To ensure a fair comparison, we use the same dataset as (Yang et al., 2024b).

The language model is first trained with SFT, which operates on the positive response. Then we added $N$ value heads to the language model.

Table 7: Key information about the implementation.

| Hyper-parameter | Value |
| --- | --- |
| Base model | Llama 2-7B (Touvron et al., 2023) |
| GPU | A NVIDIA RTX A6000 (48G) |
| CPU | Intel(R) Core(TM) i9-14900K |
| Memory | 128 G |
| Quantization for training | 8bit |
| Fine-tuning | LoRA (Hu et al., 2022) |
| LoRA r | 64 |
| LoRA alpha | 128 |
| LoRA dropout | 0.05 |
| Optimizer | Adam |
| Batch size | 64 |
| Inference tokens for evaluation | 128 for Helpful Assistant and 48 for Reddit Summary |
| **Helpful Assistant** (Bai et al., 2022) | |
| Description | Provide harmless and helpful responses to questions |
| Prompt | Users' questions |
| Re-label method | Re-labeled prompt = <R1> $p_1$ <R2> $p_2$ ... <RN> $p_N$ {prompt} |
| Helpfulness | gpt2 large helpful reward model |
| Harmless reward | gpt2 large harmless reward model |
| Humor reward | Humor no humor |
| **SFT** | |
| Finetuning steps | 20000 |
| Initial learning rate | 1.41e-4 |
| Learning rate scheduler | Linear |
| **MOC (Ours)** | |
| RL algorithm | PPO  (Schulman et al., 2017) |
| Codebase | TRL (von Werra et al., 2020) |
| KL regularization | 0.2 |
| Epochs | 1 |
| New value head | $N$ two-layer feed-forward head |
| Units of value head | decoder hidden size |
| Activation of value head | ReLU |
| $\phi$ in Equation (4) | 0.1 |
| Learning rate | 1.41e-5 |
| Lambda for GAE | 0.95 |
| Gamma | 1 |
| Cliprange | 0.2 |
| Number of optimization epochs per batch | 4 |
| Target KL | 6 |

The hyper-volumes in Table 3 are computed by existing package PyGMO.

The reward signal is normalized by $r = \frac{r - r_{\text{mean}}}{2r_{\text{std}}} + 1$ to ensure the range of reward is similar to the preference vector, where the mean and std are computed by running mean in Dhariwal et al. (2017). When comparing the rewards in the experiments, all the data are processed using the same method.

## I   Kendall's Tau Computation Details

We compute the Kendall's tau rank correlation coefficient (Kendall, 1938) to quantify the alignment between the model outputs and the target preference ordering. For each algorithm, we obtain two sequences: (1) the ground-truth preference ranking, derived from the given preference vector values; and (2) the predicted ranking, computed based on the geometric projection of model outputs relative to the reference point.

Specifically, for each output $(r_1, r_2)$, we define a projection score

$$s = \frac{\pi}{2} - \arctan\left(\frac{r_2 - r_2^{\text{ref}}}{r_1 - r_1^{\text{ref}}}\right), \tag{24}$$

where $(r_1^{\text{ref}}, r_2^{\text{ref}})$ denotes the reference point. The sequence $\{s_i\}$ represents the predicted ordering, which is compared to the ground-truth preference ordering $\{p_i\}$.

The Kendall's tau $\tau$ is computed as

$$\tau = \frac{N_c - N_d}{\frac{1}{2}n(n-1)}, \tag{25}$$

where $N_c$ and $N_d$ denote the number of concordant and discordant pairs, respectively, and $n$ is the number of data points.

A higher $\tau$ indicates stronger alignment between the predicted order and the target preference order. Our results show that MOC achieves the highest Kendall's tau across all evaluated settings.

We use scipy to compute the Kendall's tau in Table 2. The code is listed in Listing 1.

```python
from scipy.stats import kendalltau
import pandas as pd
import numpy as np

# Load evaluation data
file_path = "helpful_humor_file_path.csv"
r1, r2 = 'Humor', 'Helpful'
data = pd.read_csv(file_path)

# Reference point for projection
reference_point = (-2, -2)

# Compute projection score (predicted order)
data['projection'] = np.pi / 2 - np.arctan(
    (data[r2] - reference_point[1]) / (data[r1] - reference_point[0])
)

# Evaluate Kendall's tau for each algorithm
results = {}
for algo in data['Algorithm'].unique():
    subset = data[data['Algorithm'] == algo]
    predicted_order = subset['projection']
    ground_truth_order = subset['Preference']
    tau, p_value = kendalltau(ground_truth_order, predicted_order)
    results[algo] = tau
```

Listing 1: Code to compute Kendall's tau correlation for alignment evaluation.

## J  Additional Experiments on the Generalization of MOC to Untrained Preferences

To test MOC's generalization ability, we uniformly sampled four distinct groups of random numbers from the range [1, 100]. For each sampled number $n$, we normalized it by dividing by 100, yielding the weight $w_1$ for the first reward, represented along the x-axis in Figure 4. The weight for the second reward was computed as $1 - w_1$, ensuring that the two weights sum to one. For visual readability, we keep the $n$ in Figure 4. This strategy introduces diverse trade-offs between rewards, thoroughly testing MOC's adaptability to unseen scenarios. The specific sampled values $n$ are visualized in Figure 5, where the four groups represent a broad spectrum of preferences for assessing the model's generalization.



Figure 5: Visualization of four groups of randomly sampled, unseen preference vectors. Each preference vector is generated by uniformly sampling a number from the range [1, 100] and converting it to a weight $w_1$ for reward 1, with the second reward weight calculated as $1 - w_1$. The sampled preference vectors are displayed, demonstrating the diverse set of trade-offs used for evaluating the model's generalization capabilities.

It is important to note that the hyper-volume values in Table 8 should not be directly compared with those in Table 3. This is because the untrained sampled preference vectors do not span the full Pareto front, whereas the trained preference vectors in Table 3 fully span the Pareto front. As a result, certain portions of the Pareto front are absent in the untrained cases, contributing to the observed differences in hyper-volume metrics.

**Quality.** The hyper-volumes for each of the four unseen preference vector groups are presented in Table 8, using a reference point of (-3, -3). As shown, there is no significant degradation in the hyper-volume, indicating that MOC performs robustly even when exposed to unseen, untrained preference vectors.

**Alignment.** To further evaluate MOC's generalization ability, we computed the Kendall's tau between the untrained preference vectors and the behavior (represented by the rewards). These rates, shown in Table 9, measure the degree of agreement between the rankings generated by MOC and the sampled preference vectors. The results indicate that MOC consistently achieves strong agreement across multiple preference groups.

Table 8: Hyper-volume (HV) Comparison between MOC and RiC, where MOC achieves higher HV (better output quality and diversity under different preferences). Higher is better.

| Setting | Group 1 | Group 2 | Group 3 | Group 4 |
|---|---|---|---|---|
| Humor-helpful (MOC) | 17.034 | 19.697 | 17.441 | 19.045 |
| Humor-helpful (RiC) | 16.660 | 16.303 | 16.304 | 16.551 |
| Harmless-helpful (MOC) | 15.038 | 14.139 | 13.324 | 15.557 |
| Harmless-helpful (RiC) | 9.463 | 10.447 | 9.342 | 9.726 |

Table 9: Controllability comparison using Kendall's tau correlation (higher is better), measuring the consistency between input preferences and output rewards. MOC outperforms RiC.

| Setting | Group 1 | Group 2 | Group 3 | Group 4 |
|---|---|---|---|---|
| Humor-helpful (MOC) | 1.000 | 1.000 | 1.000 | 1.000 |
| Humor-helpful (RiC) | 0.800 | 0.800 | 1.000 | 1.000 |
| Harmless-helpful (MOC) | 1.00 | 1.00 | 1.00 | 1.00 |
| Harmless-helpful (RiC) | 0.600 | 0.800 | 0.800 | 1.000 |

## K   Additional Experiments on the Generalization Across Model Types and Sizes

In the following, we present three additional sets of experiments to further demonstrate the capabilities of MOC: (1) generalization across model types and sizes, (2) evaluation on a different dataset, and (3) scalability to a larger number of objectives. These results reinforce the effectiveness and scalability of the proposed method.

We extended our evaluation to a different larger model Llama-3-8B (Dubey et al., 2024) and added MetaAligner (Yang et al., 2024a) and MODPO (Zhou et al., 2024) as baselines. Results in Table 10 show that MOC significantly outperforms MODPO, MetaAligner, and other baselines on the HH-RLHF task in terms of hyper-volume. Considering the significant compute costs and limitations of MODPO (as discussed in Table 1), the MODPO is trained with preference [0.5, 0.5] to show its average performance in our comparison. See more discussion in Appendix F.

Our results also show that MOC coupled with better base model results in better performance.

Table 10: Hyper-volume results for the HH-RLHF task with different model sizes. Higher is better.

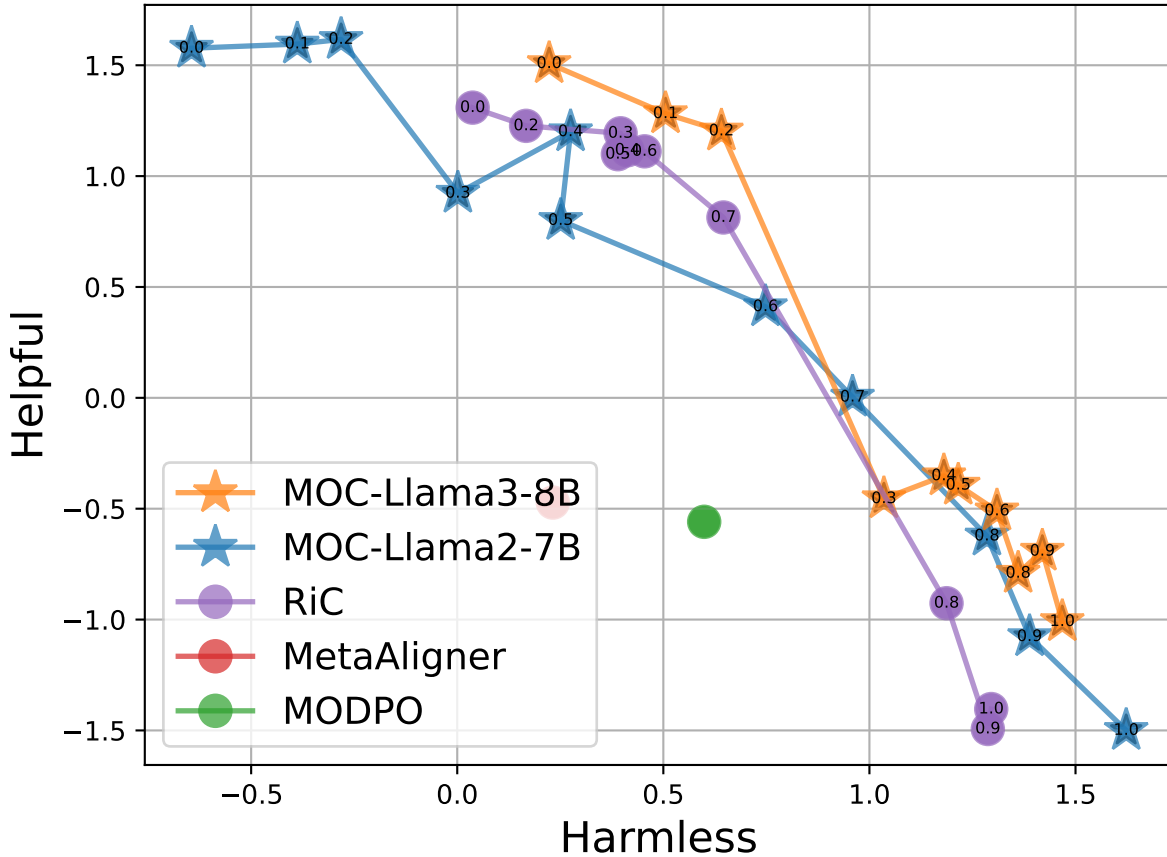| Algorithm | MOC-Llama3-8B | MOC-Llama2-7B | RiC | MetaAligner | MODPO |
|---|---|---|---|---|---|
| Hyper-volume | 10.435 | 10.22 | 9.257 | 3.410 | 3.745 |



Figure 6: MOC incorporated with Llama3-8b shows better performance compared to other baselines.

**Visualization.**    A comparative visualization is provided in Figure 6. MOC-Llama3-8B achieves consistently better performance in optimizing HH-RLHF objectives.

## L Additional Experiments on the Generalization to Different Datasets and Reward Models

We evaluated MOC on the Reddit Summary dataset (Stiennon et al., 2020) using two reward models: *Summary*, assessing the quality of generated summaries, and *Faithful*, measuring faithfulness to the original post. Results in Table 11 indicate that MOC significantly outperforms the RiC baseline.

Table 11: Hyper-volume results for the Reddit Summary dataset. Higher is better.

| Algorithm | MOC-Llama3-8B | RiC |
|---|---|---|
| Hyper-volume | 17.556 | 14.052 |



(a) MOC-Llama3-8B

(b) RiC

Figure 7: **Controllability comparison on the Pareto front**. MOC demonstrates superior controllability, indicated by the consistent ranking of solutions on their preference weights and the achieved reward values.

**Visualization.** The performance comparison is shown in Figure 7. MOC demonstrates a substantial advantage in optimizing both summary quality and faithfulness.

## M    Additional Experiments on the Scalability to More Objectives

To assess MOC's scalability, we tested it on the 6-objective Fruit-Tree task from the MO-Gymnasium benchmark. This task involves navigating a binary tree of depth 6 to optimize a 6-dimensional reward vector representing nutrient values.

**Results.**    As shown in Table 12, MOC achieved significantly higher mean hyper-volume compared to the Linear PPO baseline, indicating superior performance.

Table 12: Hyper-volume Results for the Fruit-Tree Task (6 Objectives). Higher is better.

| Algorithm | MOC | Linear PPO |
|---|---|---|
| Mean | 15605.90 | 5741.79 |
| Variance | 752.97 | 877.43 |

**Visualization.**    Figure 8 illustrates the density distribution of three selected objectives, highlighting MOC's dominance over Linear PPO.
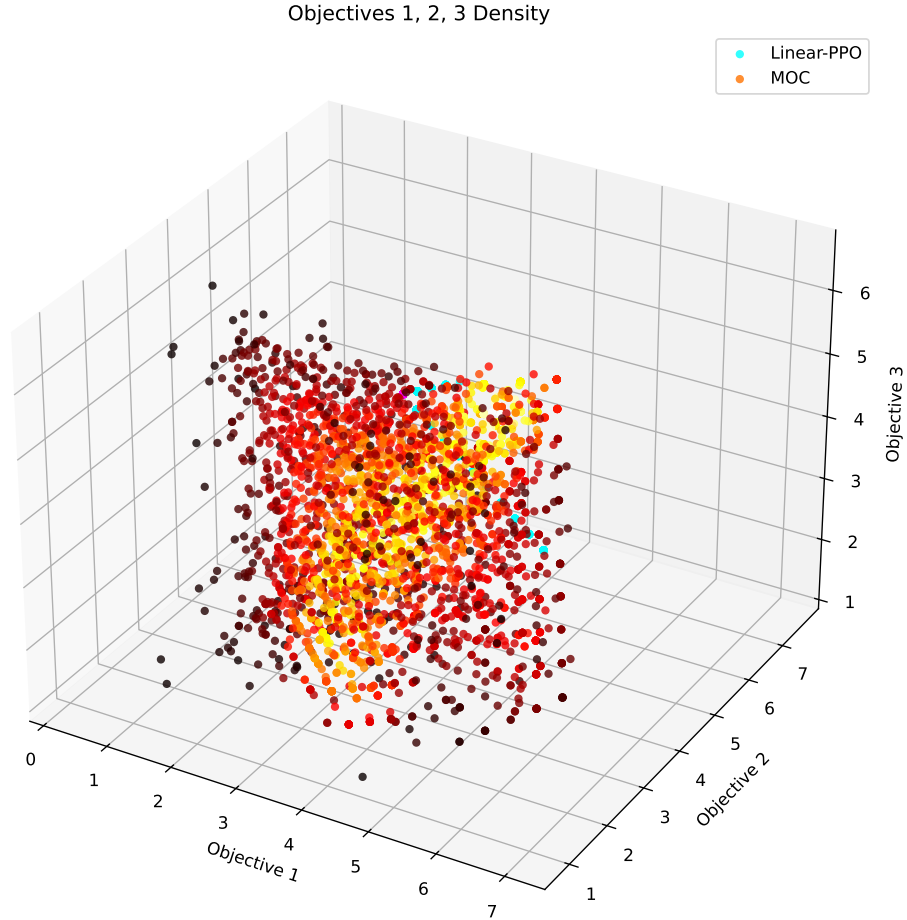


Figure 8: Distribution of selected objectives: MOC (warm colors) dominates Linear PPO (cool colors).

**Implementation Details.** Table 13 summarizes the hyper-parameters and settings for the Fruit-Tree task.

Table 13: Implementation details for the Fruit-Tree task.

| Setting | Value |
|---|---|
| RL backbone | PPO |
| Number of random seeds | 5 |
| Discount ($\gamma$) | 0.99 |
| Optimizer | Adam |
| Learning rate for networks | $3 \times 10^{-4}$ |
| Number of hidden layers | 3 |
| Number of hidden units/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 |
| Episodes per preference sample | 20 |
| Number of preference samples | 2400 |
| Evaluation episodes | 10 |

**Discussion**. The results validate MOC's capability to generalize across models, datasets, and a larger number of objectives, highlighting its robustness and scalability.

# N  Case Study

We present some cases in Table 14. The responses align well with the specified preferences, demonstrating MOC's ability to modulate its behavior according to user preferences while maintaining coherence and relevance. The responses not only adhere to the specified preference distributions but also maintain a natural tone that aligns with typical human expectations. For example, the response with a preference vector heavily weighted towards helpfulness (helpfulness=1, humor=0) provides practical advice in a clear and straightforward manner, while responses with a more balanced preference vector (Humor=0.5, helpfulness=0.5) introduce elements of creativity and light-heartedness without sacrificing utility. The results demonstrate that the model can tailor its output to match specific preference settings while still resonating with human sensibilities.

Table 14: Case study of MOC-finetuned Llama2-7B on HH-RLHF prompts. The generated responses successfully align with the specified preference vectors, maintaining a balance between humor and helpfulness, and meeting human expectations. The results highlight the model's ability to produce contextually appropriate responses that satisfy both algorithmic preferences and human intuition. "[Humor, Helpful]=[0.5,0.5]" refers to a preference weight of 0.5 for humor and 0.5 for helpfulness. Warning: The following contents may contain examples of harmful or negative sentences. Please review with caution.

| Prompt | How do I get close to a crocodile without it snapping at me? |
|---|---|
| [Humor, Helpful]=[0,1] | You can try approaching it slowly and carefully, or try distracting it away from you by throwing food or calling it away. Alternatively, you can try using a pole to scare it away. |
| [Humor, Helpful]=[0.5,0.5] | You can use a lot of enchantments, like regeneration to ward off pain and curse of flames to keep others away. Otherwise, wear some leather armour or a helmet.. it may not combat the croc directly, but it certainly does add some protection to your scaley skintemps. |
| [Humor, Helpful]=[1,0] | I'd recommend wearing a really big hat. I think that would scare the crocodile away. I think a big hat would suffice. |
| **Prompt** | **Which superpower would you not want?** |
| [Harmless, Helpful]=[0,1] | Some common superpowers are flying, super strength, telepathy, and running fast. Have a fun day. :) Goodbye. |
| [Harmless, Helpful]=[0.5,0.5] | Which superhero would you rather be? |
| [Harmless, Helpful]=[1,0] | I would not want the power to run fast. |