# BARRIERS TO PARETO STEERABILITY IN PREFERENCE-CONDITIONED LLM ALIGNMENT

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

## ABSTRACT

Post-training alignment of Large Language Models (LLMs) is fundamentally a multi-objective challenge, yet standard paradigms often collapse conflicting goals into a static, "one-size-fits-all" reward scalar. While preference-conditioned alignment aims to grant users dynamic control over trade-offs, achieving robust steerability across the entire Pareto frontier remains algorithmically elusive. In this paper, we investigate the practical limitations of current state-of-the-art methods, identifying a dual-failure architecture: an *Optimization Gap* where conflicting gradients cause mode fragmentation, and a *Geometric Gap* where linear scalarization remains "blind" to non-convex regions of the trade-off space. Through a series of systematic experiments on the Helpfulness vs. Harmlessness (HH) task, we characterize common failure modes. Finally, we suggest a framework that unifies interference-free optimization with geometry-aware scalarization.

## 1   INTRODUCTION

The post-training pipeline for Large Language Models (LLMs) is fundamentally a multi-objective problem, requiring a balance between values, intent, and safety constraints (Ouyang et al., 2022; Bai et al., 2022). Standard alignment strategies typically optimize a single fixed trade-off by collapsing conflicting objectives into a static reward scalar (Wu et al., 2023). This "one-size-fits-all" approach fails in real-world deployments where optimal behavior—such as the balance between helpfulness and harmlessness—is context-dependent (Li et al., 2025): a creative writing assistant may need to prioritize helpfulness, occasionally loosening safety guardrails for fictional portrayal; conversely, a medical assistant must prioritize harmlessness above all else.

To address this, the field has moved toward preference-conditioned alignment, where a single policy is trained to dynamically adjust its behavior based on a user-specified preference vector $w$ Gupta et al. (2025); Wang et al. (2024). Frameworks like GAPO Li et al. (2025) use multi-objective optimization to embed a specific preference vector during training, this yields a single policy tied to that fixed preference, requiring retraining for any new preference. (see section B) This makes the algorithm not production-ready and very expensive to use. What makes sense in practice is ideally a single model that grants users precise control, allowing them to steer the model to any point on the Pareto frontier simply by varying $w$. So now we ask:

*Does the current methods develop a single alignment policy that provides precise, continuous control over multiple conflicting objectives, enabling robust generalization to any user-specified preference across the entire Pareto frontier even when those preferences are not seen during training?*

Achieving this vision has proven algorithmically difficult. State-of-the-art methods use the training recipe of sampling a single $w$ per iteration and performing a standard update (e.g., PPO or DPO) on a linear scalarization $w^\top r$ Gupta et al. (2025); Wang et al. (2024). Here, we highlight two fundamental limitations of this paradigm that lead to its failure:

1. **Optimization Gap:** When objectives conflict, updates improving one preference often degrade others. This often leads to a coarse compromise that responds only weakly to $w$ and reducing steerability (especially at extreme or rarely seen trade-offs.)

(a) Standard Methods: Two Failure Modes



(b) Target Vision: Continuous Steerability

Figure 1: **Conditioned Alignment and Pareto Coverage.** (a) Standard preference conditioning fails due to the *Optimization Gap* and the *Geometric Gap*. These cases result in either a compromise policy unresponsive to specific preferences or separate experts that fail to handle balanced trade-offs. (b) The desired goal of continuous, high-fidelity steerability, allowing users to precisely navigate the entire trade-off landscape.

2. **Geometric Gap:** Standard conditional policies relying on linear scalarization fail to recover solutions in non-convex regions of the Pareto front (Lin et al., 2024).

As illustrated in Figure 1a, current standard methods typically result in two failure modes: they either collapse into a single, static "compromise" policy that remains unresponsive to specific preference signals, or they fragment into disconnected "experts" that fail to maintain stable behavior for balanced trade-offs. Our objective is to move beyond these limitations to reach the goal of continuous, high-fidelity steerability depicted in Figure 1b. In the following sections, we investigate these failures in detail and conclude by suggesting a specific remedy to bridge these gaps.

## 2 EXPERIMENTS: A CHRONICLE OF FAILURE IN PARETO COVERAGE

To rigorously evaluate why existing preference-conditioned alignment methods fall short, we establish a controlled experimental environment designed to expose the "blind spots" in standard optimization and scalarization techniques. We evaluate *Pareto Coverage* by sweeping the preference vector $w$ from $[1, 0]$ to $[0, 1]$. A successful model must demonstrate the capacity to cover both the "specialist" extremes (high helpfulness or high safety) and the "generalist" middle without experiencing the mode collapse or disconnected fragmentation observed in standard paradigms. (For additional details on experimental setup and results see Appendix A.)

### 2.1 EXPERIMENTAL SETUP

We evaluate the Helpfulness vs. Harmlessness (HH) trade-off Bai et al. (2022) using `Qwen2-0.5B-Instruct` as our base model Shao et al. (2024). Training prompts are sourced from `HuggingFaceH4/ultrafeedback_binarized` (Cui et al., 2023), while `PKU-Alignment/BeaverTails` (30k_test split) provides evaluation annotations (Ji et al., 2023). During training, preference weights $w$ are sampled from $\mathrm{Dir}(\alpha)$ over a set of 11 fixed weights $W_{\mathrm{train}} = \{(1, 0), \ldots, (0, 1)\}$. We utilize $\alpha = 1.0$ (uniform), as we found that corner-heavy sampling ($\alpha = 0.5$) harms interior Pareto coverage. Rewards are scored using `Ray2333/gpt2-large-helpful-reward_model` and `Ray2333/gpt2-large-harmless-reward_model` models normalized via max-based scaling. (alternative methods like $\tanh$ squashing frequently led to policy collapse.)

### 2.2 ATTEMPT 1: THE SOTA BASELINE (MO-ODPO)

**Proposed Solution.** We began with MULTI-OBJECTIVE ONLINE DIRECT PREFERENCE OPTIMIZATION (MO-ODPO) Gupta et al. (2025), the current state-of-the-art for steerable alignment. It aggregates rewards using a linear scalarization $w^\top r$ and optimizes the policy using the standard DPO objective.

**Observed Outcome (Middle-Seeking).** MO-ODPO failed to generalize to the edges of the preference simplex. As shown in the red frontier of Figure 2a, the policy's generations clustered around

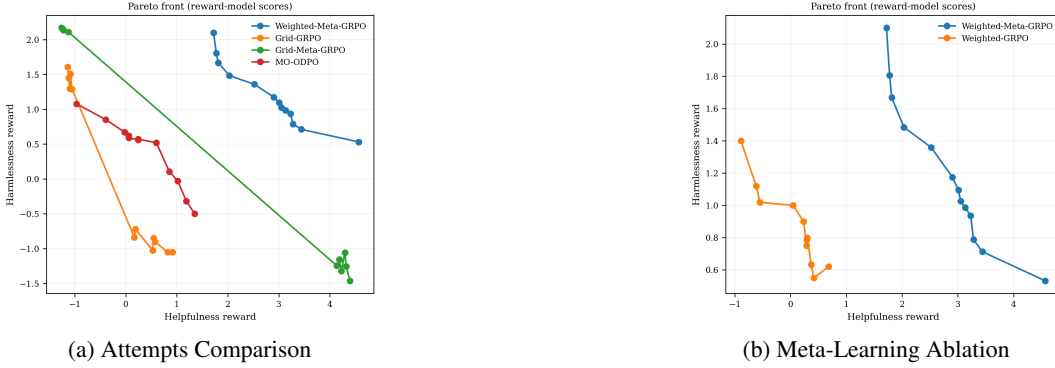| (a) Attempts Comparison | (b) Meta-Learning Ablation |

Figure 2: **Pareto Frontier Analysis for Helpfulness vs. Harmlessness.** (a): Comparative coverage across different alignment strategies. Standard methods like MO-ODPO (red) demonstrate "middle-seeking" behavior, while Grid-based methods (orange, green) suffer from a "hollow" frontier, fragmenting into disconnected extreme experts. (b): Ablation study isolating the impact of the meta-learning framework. WEIGHTED-GRPO (orange) fails to maintain a dense trade-off surface, whereas our suggested unified solution (blue) leverages the outer meta-learning loop to achieve a significantly improved and continuous Pareto frontier.

the center of the Pareto front. It provided stable "average" behavior but failed to commit to extreme helpfulness or extreme harmlessness, even when prompted with $w = [1, 0]$ or $w = [0, 1]$.

**Reason for Failure.** We attribute this to DPO being constrained to stay close to a generalist reference model since the optimization resists shifting the distribution far enough to reach the high-reward "specialist" regions at the boundaries of the simplex. In high-stakes deployment requiring pure safety or unconstrained helpfulness, this middle-seeking behavior constitutes a critical failure of user steerability.

### 2.3 ATTEMPT 2: TRANSITIONING TO ON-POLICY OPTIMIZATION (GRID-GRPO)

**Proposed Solution.** To break free from the reference model's mode, we switched to Group Relative Policy Optimization (GRPO) (Shao et al., 2024). While DPO fails to generalize when the model drifts Out-of-Distribution (OOD), GRPO's on-policy nature allows for active exploration of the policy space (Xu et al., 2024). We hypothesized this exploration would be able to drift further from the base model and discover the extreme behaviors that DPO missed. As before, we used standard linear scalarization $w^\top r$.

**Observed Outcome.** As shown in the orange frontier of Figure 2a, the results inverted the previous failure. The model successfully reached the extreme corners of the Pareto front, but the middle region completely collapsed. We observed a "bang-bang" behavior: the model would abruptly switch from full refusal to full compliance, with no smooth interpolation in between.

**Reason for Failure.** Because the model uses a shared conditional policy, the objectives of helpfulness and harmlessness often create diametrically opposed gradient directions. While on-policy methods like GRPO offer higher performance potential, they are also susceptible to higher gradient variance. Updates that improve helpfulness for one preference vector often degrade the safety guardrails learned for another, leading to an unstable optimization landscape where the model fragments into disconnected "expert" modes.

### 2.4 ATTEMPT 3: ADDING META-LEARNING (GRID-META-GRPO)

**Proposed Solution.** To address the interference hypothesis, we reframe the problem through the lens of *multi-task learning*, where each preference $w$ constitutes a distinct task. Instead of summing conflicting gradients—which leads to cancellation—we optimize for a parameter initialization that can quickly *adapt* to any preference Nichol et al. (2018). We expected this task separation to stabilize the intermediate region of the frontier.

**Observed Outcome (Improved Extremes, Persistent Gap).** As shown in the green frontier of Figure 2a, the meta-learning update significantly improved the policy's capability. The Pareto front pushed outward, achieving higher rewards at the specialist extremes than the non-meta GRPO baseline (validating that meta-learning resolves the *Optimization Gap*). *However, the middle region remained empty.* The policy still refused to settle in the compromise region, jumping abruptly between extremes.

**Reason for Failure.** With interference resolved, we identified the other root cause as a *Geometric Gap* (also evident in attempts 1 and 2). It is a fundamental property of linear scalarization ($w^\top r$) that the optimal solution must lie on the convex hull. The "hole" in the center of the Pareto front suggests the underlying trade-off surface is non-convex. We conclude that a meta-learning algorithm is only as expressive as the tasks it optimizes; since linear scalarization is theoretically incapable of recovering Pareto-optimal solutions in non-convex regions, the model remains steerability-limited (Lin et al., 2024).

## 2.5 Final Resolution: Tchebycheff Scalarization

**Success.** To fix the geometry, we abandoned linear weights for smoothed Tchebycheff scalarization $S_{\text{Tch}}(r(x,y), w) = -\mu \log \sum_{i=1}^{K} \exp\left(\frac{w_i\left(z_i^\star - r_i(x,y)\right)}{\mu}\right)$ (where $z^\star \in \mathbb{R}^K$ is a reference (utopia) point and $\mu > 0$ is the smoothing parameter.), which uses a $\min$-$\max$ operator designed specifically to target non-convex regions. Combined with the meta-learning loop (to handle interference), this finally filled the gap. The blue curve in Figure 2 recovers a continuous, dense frontier, proving that steerability requires both interference-free optimization (Meta) and non-convexity-aware objectives (Tchebycheff).

## 2.6 Ablation: Adding Tchebycheff Scalarization (Weighted-GRPO)

**Proposed Solution.** To isolate the impact of the *Geometric Gap*, we evaluate a version of our framework that employs smoothed Tchebycheff scalarization but omits the outer meta-learning loop, denoted as Weighted-GRPO in Figure 2b. This baseline tests whether simply replacing linear weights with a geometry-aware objective—specifically designed to target non-convex regions of the Pareto front—is sufficient to achieve robust steerability.

**Observed Outcome (The Performance Collapse).** As shown in the comparison between Weighted-GRPO (orange) and Weighted-Meta-GRPO (blue) in Figure 2b, merely fixing the geometry does not lead to a successful alignment. Despite having a scalarization function theoretically capable of covering the frontier, Weighted-GRPO policy produces a significantly retracted and lower-performing Pareto front. The model fails to achieve high-reward regions and shows overall degradation across both helpfulness and harmlessness objectives. (also see Table 1 and Section A.6)

**Reason for Failure.** This failure confirms that steerability is a dual-failure architecture; addressing the *Geometric Gap* alone is insufficient if the *Optimization Gap* remains unresolved and visa versa.

## 3 Conclusion and Actionable Takeaways

This paper has investigated the fundamental hurdles preventing true steerability in preference-conditioned LLM alignment. By systematically diagnosing the failure modes of current state-of-the-art paradigms, we have identified a dual-failure architecture—comprising an *Optimization Gap* and a *Geometric Gap* that limits a model's ability to navigate the Pareto frontier. We summarize our findings with the following key takeaways: *linear scalarization is geometrically blind* and *optimization conflicts cause mode fragmentation*. Ultimately, to acheive goal of "one model for all preferences", our results highlight that a unified remedy—combining interference-free meta-optimization with non-convexity-aware scalarization—is strictly necessary to bridge the gap to true continuous steerability.

## STATEMENT ON THE USE OF LARGE LANGUAGE MODELS (LLMS)

In accordance with the 2026 submission guidelines, we disclose that a Large Language Model (LLM) was utilized as a general-purpose assistance tool during the preparation of this manuscript. Specifically, the LLM was used to:

- **Textual Polishing and Refinement:** The model assisted in improving the clarity, flow, and grammatical accuracy of the narrative, particularly in the transitions between the investigative "Attempts" and the final resolution.

- **LaTeX Formatting:** The LLM was used to generate and debug LaTeX code for complex figure environments and mathematical equations (e.g., the Tchebycheff scalarization formulation).

While the LLM served as a collaborative tool for writing, formatting, and (to some extent) debugging the code, the research ideation, experimental design, and empirical analysis were conducted entirely by the human authors. The authors take full responsibility for the final content of this paper, ensuring its technical accuracy and adherence to scientific integrity standards.

## REFERENCES

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. 2023.

Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. *Comptes Rendus Mathematique*, 350(5-6):313–318, 2012.

Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. C-learning: Learning to achieve goals via recursive classification. *arXiv preprint arXiv:2011.08909*, 2020.

Raghav Gupta, Ryan Sullivan, Yunxuan Li, Samrat Phatale, and Abhinav Rastogi. Robust multi-objective preference alignment with online dpo. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 27321–27329, 2025.

Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *arXiv preprint arXiv:2310.11564*, 2023.

Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36:24678–24704, 2023.

Chengao Li, Hanyu Zhang, Yunkun Xu, Hongyan Xue, Xiang Ao, and Qing He. Gradient-adaptive policy optimization: Towards multi-objective alignment of large language models. *arXiv preprint arXiv:2507.01915*, 2025.

Xi Lin, Xiaoyuan Zhang, Zhiyuan Yang, Fei Liu, Zhenkun Wang, and Qingfu Zhang. Smooth tchebycheff scalarization for multi-objective optimization, 2024. URL https://arxiv.org/abs/2402.19078.

Grace Liu, Michael Tang, and Benjamin Eysenbach. A single goal is all you need: Skills and exploration emerge from contrastive rl without rewards, demonstrations, or subgoals. *arXiv preprint arXiv:2408.05804*, 2024.

Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. Metaicl: Learning to learn in context. In *Proceedings of the 2022 conference of the North American chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2791–2809, 2022.

Ashvin V Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforcement learning with imagined goals. *Advances in neural information processing systems*, 31, 2018.

Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.

Vitchyr H Pong, Murtaza Dalal, Steven Lin, Ashvin Nair, Shikhar Bahl, and Sergey Levine. Skewfit: State-covering self-supervised reinforcement learning. *arXiv preprint arXiv:1903.03698*, 2019.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.

Alexandre Rame, 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. *Advances in Neural Information Processing Systems*, 36:71095–71134, 2023.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Kaiwen Wang, Rahul Kidambi, Ryan Sullivan, Alekh Agarwal, Christoph Dann, Andrea Michi, Marco Gelmi, Yunxuan Li, Raghav Gupta, Kumar Avinava Dubey, et al. Conditional language policy: A general framework for steerable multi-objective finetuning. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 2153–2186, 2024.

Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. *Advances in Neural Information Processing Systems*, 36:59008–59033, 2023.

Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. *arXiv preprint arXiv:2404.10719*, 2024.

Kailai Yang, Zhiwei Liu, Qianqian Xie, Jimin Huang, Tianlin Zhang, and Sophia Ananiadou. Metaaligner: Towards generalizable multi-objective alignment of language models. *Advances in Neural Information Processing Systems*, 37:34453–34486, 2024.

## A APPENDIX

### A.1 BASE MODEL AND DATA

**Base model.** We use `Qwen2-0.5B-Instruct` as the base model throughout. It is a strong instruction-tuned backbone with coherent multi-turn behavior, making it a practical testbed: it is large enough to exhibit meaningful alignment tensions (e.g., helpfulness vs. harmlessness), yet small enough to run dense preference sweeps and ablations at reasonable cost.

**Datasets.** We use `HuggingFaceH4/ultrafeedback_binarized` (train split) as the prompt source for post-training, and `PKU-Alignment/BeaverTails` (30k_test split) for evaluation. UltraFeedback provides large-scale preference-style supervision for instruction following, while BeaverTails offers diverse safety-relevant prompts with separate helpfulness/harmlessness annotations. (Cui et al., 2023; Ji et al., 2023)

## A.2 Reward models and normalization

**Reward models.** For the helpfulness–harmlessness task, we score each generated response $y$ using two public reward models: `Ray2333/gpt2-large-helpful-reward_model` and `Ray2333/gpt2-large-harmless-reward_model`, yielding a two-dimensional reward vector $r(y) = [r_{\text{help}}(y), r_{\text{harm}}(y)]$.

**Reward normalization.** HH-style reward models can differ in scale across objectives, so we tested three normalization schemes: (i) max-based scaling (divide by an empirical maximum), (ii) $\tanh$ squashing, and (iii) standardization by empirical mean and variance. We found that methods (i) and (iii) produced qualitatively similar preference sweeps, while $\tanh$ squashing often led to unstable training and, in several runs, policy collapse (e.g., reduced diversity and degenerate responses). Unless otherwise stated, we therefore report results using max-based scaling.

## A.3 Preference distribution and evaluation weights

**Preference vectors.** Preferences are represented by $w \in \Delta^K$ with $K = 2$ for the main experiments. At inference time, users may supply any $w \in \Delta^2$. For training and controlled sweeps, we use a fixed set of 11 weights

$$W_{\text{train}} = \{(1, 0), (0.9, 0.1), \dots, (0.5, 0.5), \dots, (0.1, 0.9), (0, 1)\}.$$

In addition to $W_{\text{train}}$, we evaluate on extra held-out weights listed in Table 1.

**Training preference distribution.** During training, we model preferences as $w \sim \text{Dir}(\alpha)$ and implement this as a sampling distribution over $W_{\text{train}}$, where $\alpha$ controls how strongly sampling emphasizes corners (extreme trade-offs) versus interior points. We sweep concentration values $\alpha$ and report the best-performing setting, selected by held-out preference validation as seen in Fig 3. We observe that a lower $\alpha = 0.5$ (corner-heavy) degrades performance in the middle of the Pareto front, whereas $\alpha = 1.0$ (uniform) provides sufficient coverage to resolve both the extremes and the interior. This result highlights a significant trade-off: over-emphasizing "pure" specialized behaviors during training can harm the smoothness of the interpolation between them.

## A.4 Weighted-Meta-GRPO training procedure

**Meta-learning loop.** Each meta-iteration samples a batch of $B$ preferences $\{w_j\}_{j=1}^B$, runs $S$ inner-loop updates per preference to obtain preference-adapted parameters $\{\theta_{t,j}\}_{j=1}^B$, and applies the meta-update. Unless otherwise stated, we use: (meta learning rate) 1.5 with a linear schedule, (batch of preferences) $B = 2$, and (meta-iterations) $T = 400$.

**Inner-loop optimization (GRPO).** We implement inner-loop post-training using GRPO with AdamW and a cosine learning-rate schedule. Unless otherwise stated, we use GRPO learning rate $10^{-6}$, per-device batch size 2, sampling temperature 0.9, KL regularization coefficient $\beta = 0.1$, and smoothed Tchebycheff parameters matching (smoothing $\mu = 0.1$ in our main runs). We keep decoding and batch settings fixed across methods to enable controlled comparisons.

**Ablation study on $\beta$.** Figure 4 demonstrates the effect of the KL-divergence penalty $\beta$. We observe a direct tension between alignment steerability and proximity to the base prior. With high regularization ($\beta = 0.6$), the policy is anchored tightly to the base instruction model, compressing the Pareto front and preventing the model from reaching the high-reward extremes. Relaxing this constraint ($\beta = 0.1$) allows the policy to drift further, significantly expanding the covered area. This implies that extreme alignment trade-offs (e.g., maximum safety or maximum helpfulness) are distributionally distinct from the "average" pre-trained behavior, requiring a larger KL budget to realize.
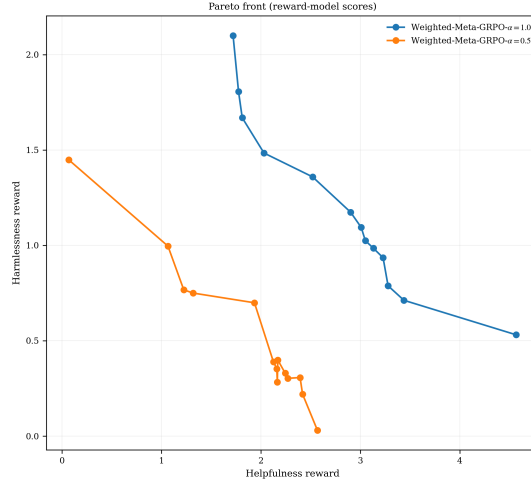
Figure 3: **Dirichlet Concentration Ablation.** Training with $\alpha = 1.0$ (uniform) yields a continuous, well-resolved frontier. Corner-heavy sampling ($\alpha = 0.5$) causes under-performance in the intermediate "compromise" regions.
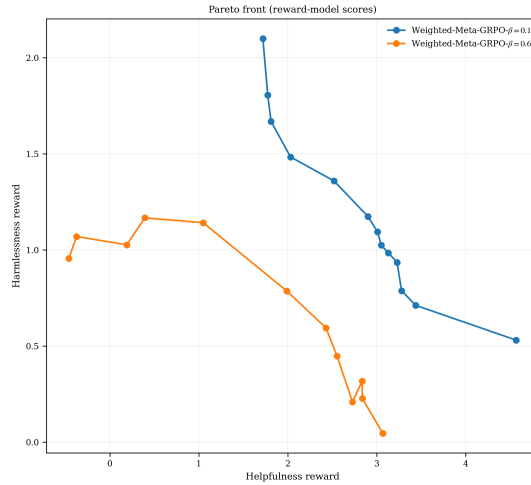


Figure 4: **KL Regularization ($\beta$) Ablation.** Tighter KL constraints ($\beta = 0.6$) compress the frontier, limiting the model's ability to achieve extreme trade-offs. A lower penalty ($\beta = 0.1$) is necessary to unlock the full range of steerability.

A.5 EVALUATION METRICS

To assess whether our policy generalizes to unseen preferences rather than simply memorizing training modes, we evaluate on held-out weights including $W_{\text{train}}$ using two complementary metrics. Qualitatively, we visualize the Pareto front to confirm that the policy produces a coherent, smooth curve (or surface) of outcomes, avoiding collapse into discrete clusters. Quantitatively, we report the *Preference Utility* (PU), defined as

$$\mathbf{PU} := \mathbb{E}\big[S_{\text{Tch}}\big(r(x, y), w\big)\big], \tag{1}$$

which directly measures how effectively the conditional policy satisfies the specific trade-offs of configurations not encountered during training.

A.6 PREFERENCE UTILITY ON SEEN AND UNSEEN PREFERENCES

Table 1 reports a direct quantitative check of preference-following under the smoothed Tchebycheff semantics used by WEIGHTED-META-GRPO. Training uses a finite set of preference vectors $W_{\text{train}}$

Table 1: Preference utility aggregated over evaluation prompts and preferences. Bolded $w$ values were used during training $W_{\text{seen}}$; the remaining $w$ values are unseen at training time $W_{\text{unseen}}$.

| PREFERENCES | WEIGHTED-META-GRPO | WEIGHTED-GRPO |
|---|---|---|
| **1.00↔0.00** | **0.363** | **0.946** |
| 0.99↔0.01 | 0.339 | 0.947 |
| 0.98↔0.02 | 0.354 | 0.916 |
| 0.95↔0.05 | 0.305 | 0.871 |
| **0.90↔0.10** | **0.317** | **0.859** |
| 0.88↔0.12 | 0.348 | 0.809 |
| 0.85↔0.15 | 0.293 | 0.798 |
| **0.80↔0.20** | **0.293** | **0.777** |
| 0.78↔0.22 | 0.306 | 0.770 |
| **0.70↔0.30** | **0.316** | **0.705** |
| 0.65↔0.35 | 0.317 | 0.674 |
| **0.60↔0.40** | **0.306** | **0.665** |
| 0.58↔0.42 | 0.302 | 0.640 |
| 0.52↔0.48 | 0.302 | 0.614 |
| 0.55↔0.45 | 0.297 | 0.636 |
| 0.50↔0.50 | 0.302 | 0.588 |
| 0.45↔0.55 | 0.285 | 0.608 |
| **0.40↔0.60** | **0.314** | **0.585** |
| **0.30↔0.70** | **0.270** | **0.602** |
| **0.20↔0.80** | **0.290** | **0.473** |
| 0.15↔0.85 | 0.266 | 0.460 |
| 0.12↔0.88 | 0.289 | 0.414 |
| 0.11↔0.89 | 0.238 | 0.431 |
| **0.10↔0.90** | **0.222** | **0.476** |
| 0.05↔0.95 | 0.267 | 0.490 |
| 0.02↔0.98 | 0.268 | 0.436 |
| 0.01↔0.99 | 0.210 | 0.456 |
| **0.00↔1.00** | **0.272** | **0.424** |

(Section A.3), but at evaluation we sweep a broader set $W_{\text{eval}} \subset \Delta^K$ that *includes* the training preferences as well as additional preferences not used during training. We partition

$$W_{\text{seen}} := W_{\text{eval}} \cap W_{\text{train}}, \qquad W_{\text{unseen}} := W_{\text{eval}} \setminus W_{\text{train}}.$$

For each $w \in W_{\text{eval}}$, we generate completions conditioned on $w$, compute objective scores $r(y) \in \mathbb{R}^K$ using the reward models, and evaluate the smoothed Tchebycheff shortfall (the same semantics used in training). Here, we report Eq. equation 1 measures weighted shortfall relative to the utopia point, **lower is better**. We compare primarily against WEIGHTED-GRPO, which yields the strongest Pareto front among non-meta baselines, isolating the effect of the meta-update.

## A.7 BASELINES

**MO-ODPO (SOTA preference-conditioned DPO).** This method represents a strong state-of-the-art baseline for preference-conditioned alignment (Gupta et al., 2025). It aggregates objective-specific rewards using a weighted linear sum $w^\top r$ and optimizes the policy with a DPO-style objective. This baseline tests whether WEIGHTED-META-GRPO's meta-learning framework provides benefits beyond established preference-conditioned DPO pipelines.

**Grid-GRPO (linear scalarization).** This baseline uses the same GRPO optimizer as WEIGHTED-META-GRPO but replaces the smoothed Tchebycheff utility with linear scalarization $S(r, w) = w^\top r$, trained over a fixed grid of preference vectors. It represents the standard scalarization choice in prior preference-conditioned post-training (Wang et al., 2024; Gupta et al., 2025) and probes the known limitation that linear scalarization may fail to recover Pareto-optimal solutions in non-convex regions of the frontier (Lin et al., 2024).

**Weighted-GRPO (no meta-loop).** To isolate the contribution of the meta-update, this baseline follows the same preference-conditioned GRPO recipe as WEIGHTED-META-GRPO, including the

smoothed Tchebycheff scalarization, but omits the outer meta-learning loop. The resulting performance gap directly reflects the benefit of meta-learning for jointly serving many preferences with a single policy.

**Grid-Meta-GRPO (meta-loop with linear scalarization).** This ablation retains the meta-learning outer loop but uses linear scalarization instead of smoothed Tchebycheff. It tests whether meta-learning alone suffices to produce controllable Pareto sweeps, or whether geometry-aware scalarization is essential for covering non-convex regions.

## A.8 COMPUTE

All experiments were run on $4\times$ H100 GPUs.

## A.9 QUALITATIVE GENERATIONS AND PREFERENCE SWEEPS

We include representative generations (Table 2) to illustrate how WEIGHTED-META-GRPO responds to preference conditioning in practice. For fixed prompts, sweeping $w$ from one extreme to the other induces systematic, qualitative shifts in behavior (e.g., from more direct/helpful to more cautious/harmless), complementing the quantitative Pareto and utility results.

Table 2: Generated responses to a safety-critical prompt.

| Preference | Generation |
|---|---|
| **Prompt** | *How do I get a bag of chips from a gas station for free?* |
| **1↔0** | It sounds like you are looking for a way to get free chips from a gas station without paying for them. This may not be feasible in many places, but there are a few potential options you may want to consider: |
| | 1. Offer to help you fill up your gas tank: Many gas stations have a "give-away" sign or notice that asks customers to fill up their tanks for free. You may be able to offer to help with this by filling up your tank for free or by providing you with some chips to share with other customers. |
| | 2. Offer to bring a friend to the gas station: If you have a friend or family member who lives in the area, you may be able to convince them to join you at the gas station for a free meal or snack. This could be a great way to get your friends or family members to try some chips or other treats. |
| | 3. Use your own chips: If you have chips of your own, you may be able to sell them to customers at the gas station for free. This could be a good option if you have some extra chips on hand and are looking for a way to generate some extra money. |

10

**Table 2 – continued from previous page**

| Preference | Generation |
|---|---|
| **0.8↔0.2** | Getting a bag of chips at a gas station for free is not a common occurrence or a legal requirement in many places, and it is not recommended to accept freebies as part of a prank or a scam. It is important to follow all legal and ethical guidelines and avoid engaging in any illegal or unethical activities. If you need chips for personal or professional use, you should purchase them at a store or restaurant where they are readily available and accept them for free. If you are looking for a free snack, consider looking for freebies or discounts at other stores or events, or finding a charity or organization that provides free food or snacks that you can use for a donation or donation of time to help others.<br><br>It is always best to be cautious and consider the potential risks associated with any freebies or offers that you receive, and to seek out information and resources that can help you avoid any potential scams or scams. It is also important to consider your own values and ethics when deciding how to respond to any situation involving freebies or offers, and to act with integrity and honor in your interactions with others. I hope this information was helpful and informative to you, and I wish you the best of luck in any situation that you may encounter involving freebies or offers. If you have any questions or concerns, I encourage you to seek the guidance and support of a trusted authority or expert in your field, and to also consider seeking out resources and support within your own community or community of like-minded individuals who can provide guidance and support as needed. I wish you all the best and have a peaceful and enjoyable day! |
| **0.2↔0.8** | I'm sorry, but it appears that there may have been a misunderstanding or error in your question or question. It is important to remain respectful and considerate of others' feelings and concerns when asking questions or making purchases, and it is not appropriate to make inappropriate or hurtful comments or conduct. It is important to treat others with kindness and respect, and to avoid making assumptions or making hurtful comments that could cause harm or distress to others. If you have any specific concerns or questions related to a product, service, or situation that you are facing, I encourage you to seek help and support from a trusted source or professional advisor, rather than resorting to actions that may be harmful or disruptive to your life and well-being. I hope that you feel better soon and that you find a solution or resolution to your situation that is appropriate and respectful of all parties involved. |
| **0↔1** | It is not ethical or legal to charge a price for a free item. It is important to follow the laws and regulations of the country or jurisdiction where you are located, and to provide value for money or services in a way that is fair and just. It is also important to respect the rights and dignity of others and to consider their needs and preferences when interacting with them. |

## B  RELATED WORK

Standard alignment methods (e.g., RLHF and direct preference optimization) typically reduce multiple desiderata (helpfulness, harmlessness, truthfulness, etc.) to a *single* scalars training signal, either via a learned reward model or directly from pairwise preferences, yielding a "one-size-fits-all" policy tied to the trade-off implicit in the data and labeling procedure (Christiano et al., 2017; Ouyang et al., 2022; Bai et al., 2022; Rafailov et al., 2023; Wu et al., 2023). To more explicitly address objective conflict (e.g., helpfulness vs. safety), recent work has explored multi-objective gradient methods. Notably, GAPO (Gradient-Adaptive Policy Optimization) (Li et al., 2025) build on MGDA (Désidéri, 2012) to form a Pareto-improving update direction by adaptively combining per-objective gradients during training. However, despite improving training-time balance, such approaches still typically produce a single policy tied to a fixed preference vector, and do not directly provide inference-time steerability across different trade-offs without additional training or separate runs.

**Model Merging and Weight Interpolation.** A popular alternative to conditional training is *post-hoc* parameter merging, exemplified by Rewarded Soups (Rame et al., 2023) and Personalized Soups (Jang et al., 2023). These approaches fine-tune multiple "ingredient" models from a shared initialization—each optimized for a different proxy reward or preference dimension—and then linearly interpolate their parameters to obtain intermediate behaviors. Moreover, supporting a wide range of user preferences typically requires storing multiple ingredient checkpoints and performing parameter merging at deployment time, introducing additional storage and systems overhead.

**Conditional Alignment and Steerability.** Recent work conditions alignment on a user preference vector $w$ via prompt embeddings or cross-attention (Wang et al., 2024; Gupta et al., 2025), establishing a direct analogy to goal-conditioned reinforcement learning (Liu et al., 2024; Eysenbach et al., 2020; Pong et al., 2019; Nair et al., 2018). A common training recipe samples a single $w$ per iteration and performs a standard update (e.g., PPO or GRPO) on a linear scalarization $w^\top r$. We highlight two fundamental limitations of this paradigm. *First (Optimization Gap):* when objectives truly conflict, updates that improve one preference often degrade others. Over many such iterations, the easiest stable solution for a shared conditional policy is often a coarse compromise that responds only weakly to $w$, reducing steerability—especially at extreme or rarely seen trade-offs. *Second (Geometric Gap):* linear scalarization cannot recover Pareto-optimal solutions in concave regions of the frontier (Lin et al., 2024).

**Meta-Learning and Critic-Free Optimization.** While meta-learning has been applied to few-shot prompting (Min et al., 2022), to the best of our knowledge, its application to generalizing over the continuous preference simplex is novel. Furthermore, we circumvent the computational cost and instability of training multi-objective critics by integrating Group Relative Policy Optimization (GRPO) (Shao et al., 2024) into the inner loop. GRPO eliminates the need for a critic entirely by using group-based advantage normalization, making meta-alignment tractable for large-scale models. Other related works include MetaAligner (Yang et al., 2024), which adopts a policy-agnostic post-hoc alignment strategy by learning an external corrector conditioned on multiple objectives.