

# BARRIERS TO PARETO STEERABILITY IN PREFERENCE-CONDITIONED LLM ALIGNMENT

**Fatemeh Nourzad**  
The Ohio State University  
nourzad.1@osu.edu

**Daouda Sow**  
The Ohio State University  
sow.53@osu.edu

**Yingbin Liang**  
The Ohio State University  
liang.889@osu.edu

**Ming Shi**  
University at Buffalo, State University of New York  
mshi24@buffalo.edu

**Ming Zhang**  
Google  
mingzhang@google.com

**Yunxuan Li**  
Google  
yunxuanli@google.com

**Eylem Ekici**  
The Ohio State University  
ekici.2@osu.edu

**Ness Shroff**  
The Ohio State University  
shroff.11@osu.edu

## ABSTRACT

Post-training alignment of Large Language Models (LLMs) is inherently a multi-objective problem, yet standard paradigms often collapse conflicting goals into a single "one-size-fits-all" reward scalar. While preference-conditioned alignment aims to give users dynamic control over trade-offs, achieving broad and continuous steerability across the Pareto frontier remains challenging in practice. In this paper, we investigate the limitations of current state-of-the-art methods on the Helpfulness vs. Harmlessness (HH) task and identify two recurring failure modes: an *Optimization Gap*, where conflicting gradients lead to fragmented behavior across preferences, and a *Geometric Gap*, where linear scalarization cannot recover non-convex regions of the trade-off space. Through a sequence of systematic experiments, we show how these two gaps manifest in practice and demonstrate that addressing only one of them is insufficient in the setting studied here. Finally, we evaluate a framework that combines meta-learning with geometry-aware scalarization, aiming to address both gaps in a unified way.

## 1 INTRODUCTION

The post-training alignment of Large Language Models (LLMs) is inherently a multi-objective problem, requiring a balance among values, user 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 is often insufficient in real-world deployments, where the desired balance between objectives such as helpfulness and harmlessness depends on the application context (Li et al., 2025). For example, a creative writing assistant may need to prioritize helpfulness; conversely, a medical assistant must prioritize harmlessness above all else.

To address this limitation, recent work has moved toward preference-conditioned alignment, in which a policy is trained to adjust its behavior according to a user-specified preference vector  $w$  (Gupta et al., 2025; Wang et al., 2024). However, not all such approaches provide continuous control within a single deployable model. Some methods optimize for a fixed preference vector during training, yielding a policy specialized to that trade-off and requiring additional retraining or separate models for new preferences (see Section B). In practice, the more ambitious goal is a single model that can be steered across a broad range of trade-offs simply by varying  $w$  at inference time. So now we ask:

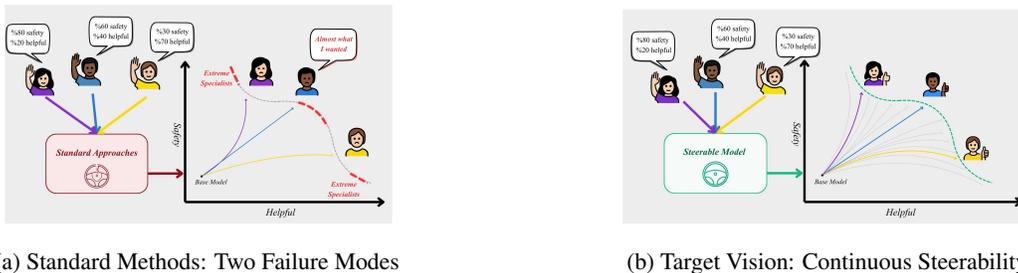


Figure 1: **Conditioned Alignment and Pareto Coverage.** (a) Standard preference-conditioning approaches can suffer from both the *Optimization Gap* and the *Geometric Gap*, leading either to a weakly responsive compromise policy or to fragmented specialists that fail to cover balanced trade-offs. (b) The desired goal of continuous, high-fidelity steerability, allowing users to navigate the trade-off landscape smoothly within a single model.

*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 Pareto frontier even when those preferences are not seen during training?*

Achieving this goal has proven challenging in practice. State-of-the-art methods commonly sample a single preference vector  $w$  at each iteration and perform a standard update (e.g., PPO or DPO) on a linear scalarization  $w^\top r$  Gupta et al. (2025); Wang et al. (2024). In this paper, we study two recurring limitations of this paradigm:

1. **Optimization Gap:** When objectives conflict, updates that improve performance for one preference can degrade performance for others. This can produce a compromise policy that responds only weakly to  $w$ , or unstable behavior across preferences, thereby reducing steerability.
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 goal is to move beyond these limitations to reach the goal of continuous, high-fidelity steerability shown in Figure 1b. In the following sections, we investigate these failures in detail and conclude by suggesting a remedy to address these gaps.

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

To study why existing preference-conditioned alignment methods fall short, we construct a controlled experimental setting designed to expose limitations in both the optimization procedure and the scalarization objective. We evaluate *Pareto coverage* by sweeping the target preference vector  $w$  from  $[1, 0]$  to  $[0, 1]$ , generating model responses under each target preference, and scoring those responses with separate helpfulness and harmless reward models. For each target weight  $w$ , we average the resulting helpfulness and harmless scores over the evaluation prompts and plot the corresponding 2D point in reward-model space. Sweeping across target preferences therefore produces an empirical trade-off curve, where each point represents the average behavior induced by one conditioning weight. In our setting, a successful method should cover both the specialist extremes (high helpfulness or high harmless) and the intermediate trade-off region, without collapsing to a weakly responsive compromise or fragmenting into disconnected modes.

### 2.1 EXPERIMENTAL SETUP

We study the Helpfulness vs. Harmlessness (HH) trade-off (Bai et al., 2022) using Qwen2-0.5B-Instruct as the base model. Training prompts are drawn from

HuggingFaceH4/ultrafeedback\_binarized (Cui et al., 2023), and evaluation is conducted on the 30k\_test split of PKU-Alignment/BeaverTails (Ji et al., 2023). During training, preference weights  $w$  are sampled uniformly ( $\text{Dir}(\alpha = 1)$ ) from a set of 11 fixed weights  $W_{\text{train}} = \{(1, 0), \dots, (0, 1)\}$ . We utilize  $\alpha = 1.0$  (uniform), as we found that corner-heavy sampling ( $\alpha = 0.5$ ) harms interior Pareto coverage. We score generations using Ray2333/gpt2-large-helpful-reward\_model and Ray2333/gpt2-large-harmless-reward\_model, with rewards normalized via max-based scaling. We found this normalization to be more stable than alternatives such as tanh squashing, which often led to policy collapse. Additional implementation details and supplementary results are provided in Appendix A.

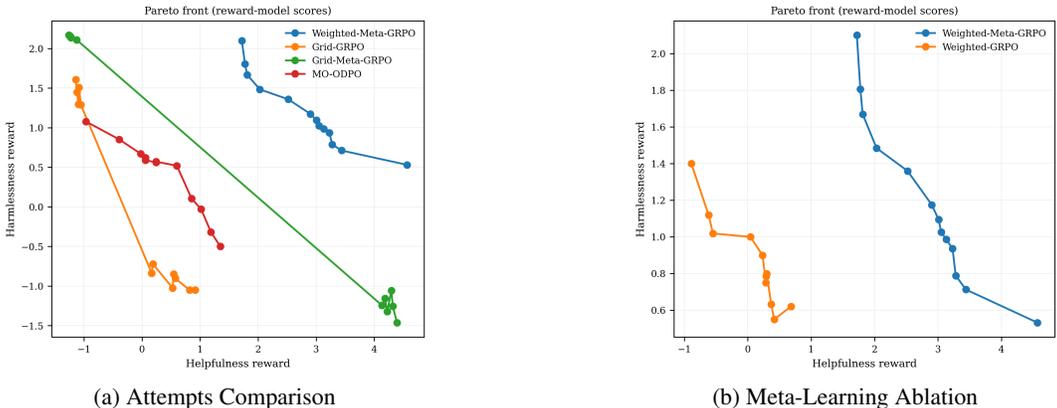


Figure 2: **Pareto Frontier Analysis for Helpfulness vs. Harmlessness.** (a): Comparative coverage across different alignment strategies. Standard methods like MO-ODPO (red) demonstrate “middle-seeking” behavior, while Grid-based methods (orange, green) suffer from a “hollow” frontier, fragmenting into 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.

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

**Proposed Solution.** We begin with MULTI-OBJECTIVE ONLINE DIRECT PREFERENCE OPTIMIZATION (MO-ODPO) (Gupta et al., 2025), a recent preference-conditioned alignment baseline. In the original work, MO-ODPO was shown to improve Pareto coverage over prior approaches such as ODPO Soups, Reward Soups, and Reward-in-Context, making it a strong reference point for our study. MO-ODPO aggregates rewards using a linear scalarization  $w^\top r$  and optimizes the policy through the standard DPO objective.

**Observed Outcome (Middle-Seeking).** MO-ODPO showed limited coverage near the extremes of the preference simplex. As shown by the red frontier in Figure 2a, its generations clustered in the intermediate region of the trade-off curve. The model produced relatively stable average behavior, but did not adequately reach either extreme helpfulness or extreme harmlessness, even when conditioned on  $w = [1, 0]$  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 our setting, this leads to a middle-seeking policy that responds only weakly to pure safety or unconstrained helpfulness, limiting user steerability at the extremes.

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

**Proposed Solution.** To move beyond the reference-model, we switched to Group Relative Policy Optimization (GRPO)<sup>1</sup> (Shao et al., 2024), an on-policy method that allows the policy to explore its own induced distribution more directly. When the desired behavior departs substantially from the reference model, GRPO’s on-policy flexibility can help the policy reach regions of the trade-off space that DPO may underexplore (Xu et al., 2024). We therefore hypothesized that GRPO would better recover the extreme behaviors that MO-ODPO missed. As before, we used the 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 still relies on a single shared conditional policy, the objectives of helpfulness and harmlessness can conflict strongly. While on-policy methods like GRPO offers greater flexibility than DPO, they are also susceptible to higher gradient variance. Updates that improve helpfulness for one preference vector often degrade harmlessness 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*, treating each preference vector  $w$  as a distinct task. Rather than directly combining gradients from different preferences—which can interfere destructively—we optimize for a parameter initialization that can quickly *adapt* to a 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 by the green frontier in Figure 2a, the meta-learning update improved frontier quality relative to the non-meta GRID-GRPO baseline. In particular, the frontier pushed outward and achieved stronger rewards in the specialist extreme regions. This suggests that meta-learning helps mitigate the optimization interference observed in Attempt 2. *However, the intermediate region remained largely under-covered.* The policy still transitioned abruptly between extremes rather than tracing a dense compromise region.

**Reason for Failure.** Because meta-learning improved the outer frontier but did not fill the middle, the remaining limitation appears to be the *Geometric Gap*. (also evident in attempts 1 and 2). In particular, linear scalarization ( $w^\top r$ ) can recover only supported Pareto-optimal points on the convex hull of the achievable trade-off region. The “hole” in the center of the frontier therefore suggests that 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 address the geometric limitation of linear scalarization, we replace  $w^\top r$  with the smoothed Tchebycheff scalarization

$$S_{\text{Tch}}(r(x, y), w) = -\mu \log \sum_{i=1}^K \exp\left(\frac{w_i(z_i^* - r_i(x, y))}{\mu}\right),$$

where  $z^* \in \mathbb{R}^K$  is a reference (utopia) point and  $\mu > 0$  is the smoothing parameter. Unlike linear scalarization, this objective can better recover trade-offs in non-convex regions of the achievable frontier. When combined with the meta-learning loop to mitigate optimization interference, this substantially improves coverage of the previously missing middle region. As shown by the blue curve in Figure 2, the resulting method produces a broader and denser frontier than the alternatives

<sup>1</sup>We use GRPO as a practical on-policy alternative to PPO-style updates, avoiding the need for an explicit value-function critic while still permitting exploration beyond the reference-model distribution.

considered above. This provides evidence that both improved optimization and geometry-aware scalarization are needed for strong continuous steerability in this setting.

## 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 by the comparison between WEIGHTED-GRPO (orange) and WEIGHTED-META-GRPO (blue) in Figure 2b, improving the scalarization alone is not sufficient in our setting. 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 degraded performance across both helpfulness and harmlessness objectives. (also see Table 1 and Section A.6)

**Reason for Failure.** This ablation complements Attempt 3: addressing the *Geometric Gap* alone does not suffice when optimization interference remains unresolved, just as improved optimization alone did not recover the missing middle under linear scalarization. Taken together, these results support the view that both ingredients are needed.

## 3 CONCLUSION AND ACTIONABLE TAKEAWAYS

In this paper, we studied why continuous preference-conditioned steerability remains difficult to achieve in LLM post-training. Through a sequence of controlled experiments on the Helpfulness vs. Harmlessness (HH) setting, we identified two recurring limitations: an *Optimization Gap*, in which conflicting updates destabilize behavior across preferences, and a *Geometric Gap*, in which linear scalarization fails to recover non-convex regions of the trade-off frontier. Our experiments suggest the following practical takeaways: *linear scalarization alone can leave parts of the frontier under-covered*, and *improved optimization alone does not guarantee smooth coverage across preferences*. In our setting, the strongest frontier coverage emerged only when meta-learning was combined with a geometry-aware scalarization objective. These results support the view that achieving “one model for all preferences” requires addressing both optimization interference and frontier geometry together.

## 4 LIMITATIONS AND FUTURE WORK

Our findings lie within the scope of a Helpfulness vs. Harmlessness (HH) study, conducted on Qwen2-0.5B-Instruct. While this scale is sufficient to expose clear failure modes in Pareto coverage, studying the patterns arise at larger model scales is an interesting direction for future work. In addition, we study the two-objective setting ( $K = 2$ ), which is already a non-trivial case and allows direct visualization of the trade-off frontier, making the optimization and geometric failure modes easier to inspect. Extending these results to higher-dimensional preference spaces is another interesting direction for future work. Our experiments focus on a helpfulness/harmlessness reward-model pair. Alternative reward models may induce different frontier geometries or optimization behavior, and evaluating a broader range of reward-model choices is also a natural direction for future work. More broadly, our goal in this paper is diagnostic rather than exhaustive: we aim to isolate and explain failure modes of preference-conditioned alignment in a controlled setting, rather than to provide a comprehensive benchmark across model scales, objective sets, and alignment pipelines.

## ACKNOWLEDGMENTS

This work has been supported in part by a gift from Google, the U.S. National Science Foundation under the grants: NSF AI Institute (AI-EDGE) 2112471, CNS-2312836, CNS-2225561, and

CNS2239677, Office of Naval Research under grant N00014-24-1-2729, and Army Research Laboratory under Cooperative Agreement Number W911NF-232-0225, The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government.

## 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  $\text{\LaTeX}$  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 Mathématique*, 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. Skew-fit: 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.
- Zequ 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, Jiakuan 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  $\beta = 0.1$ ,

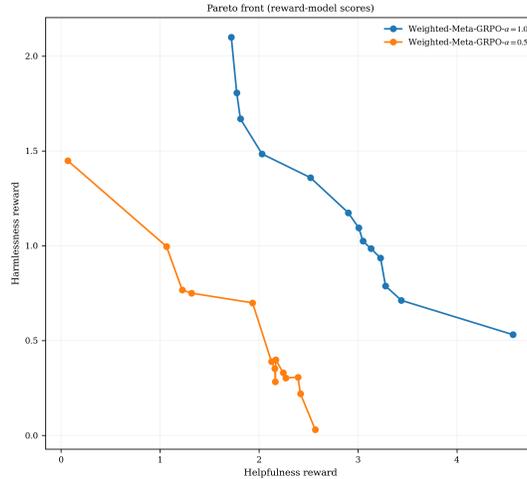


Figure 3: **Dirichlet Concentration Ablation.** Training with  $\alpha = 1.0$  (uniform) yields a continuous, well-resolved frontier. Corner-heavy sampling ( $\alpha = 0.5$ ) causes under-performance in the intermediate “compromise” regions.

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.

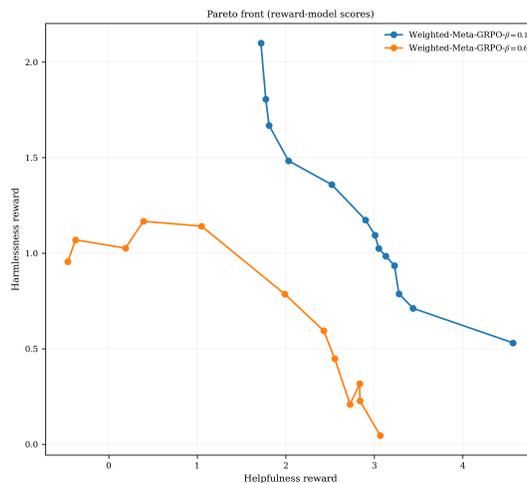


Figure 4: **KL Regularization ( $\beta$ ) Ablation.** Tighter KL constraints ( $\beta = 0.6$ ) compress the frontier, limiting the model’s ability to achieve extreme trade-offs. A lower penalty ( $\beta = 0.1$ ) is necessary to unlock the full range of steerability.

### 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-Weighted Reward Shortfall* (PWRS), defined as

$$\text{PWRS} := -\mathbb{E}[S_{\text{Tch}}(r(x, y), w)]. \tag{1}$$

This quantity measures the weighted shortfall between the achieved reward vector and the utopia point under the target preference  $w$ , so lower values indicate better alignment with the desired trade-off.

### A.6 PWRS 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  $W_{\text{train}}$  (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

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

For each  $w \in W_{\text{eval}}$ , we generate completions conditioned on  $w$ , compute the objective vector  $r(x, y) \in \mathbb{R}^K$  using the reward models, and evaluate the corresponding smoothed Tchebycheff shortfall under the same preference semantics used during training. Concretely, this metric compares the achieved reward vector  $r(x, y)$  to the utopia point  $z^*$  by forming a preference-weighted worst-case shortfall across objectives, with the smoothing parameter  $\mu$  replacing the hard max by a differentiable approximation. We then average this quantity over evaluation prompts. As a result, equation 1 measures the preference-weighted shortfall from the ideal reward vector, so **lower values indicate better alignment**. As shown in Table 1, WEIGHTED-META-GRPO consistently achieves substantially lower shortfall than WEIGHTED-GRPO, complementing the Pareto frontier results by showing that the improved coverage is also matched by better preference-conditioned performance on held-out weights. We compare primarily against WEIGHTED-GRPO, which yields the strongest Pareto frontier among the non-meta baselines and therefore isolates the contribution 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

Table 1: *Preference-Weighted Reward Shortfall* 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
<b>1.00↔0.00</b>	<b>0.363</b>	<b>0.946</b>
0.99↔0.01	0.339	0.947
0.98↔0.02	0.354	0.916
0.95↔0.05	0.305	0.871
<b>0.90↔0.10</b>	<b>0.317</b>	<b>0.859</b>
0.88↔0.12	0.348	0.809
0.85↔0.15	0.293	0.798
<b>0.80↔0.20</b>	<b>0.293</b>	<b>0.777</b>
0.78↔0.22	0.306	0.770
<b>0.70↔0.30</b>	<b>0.316</b>	<b>0.705</b>
0.65↔0.35	0.317	0.674
<b>0.60↔0.40</b>	<b>0.306</b>	<b>0.665</b>
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
<b>0.40↔0.60</b>	<b>0.314</b>	<b>0.585</b>
<b>0.30↔0.70</b>	<b>0.270</b>	<b>0.602</b>
<b>0.20↔0.80</b>	<b>0.290</b>	<b>0.473</b>
0.15↔0.85	0.266	0.460
0.12↔0.88	0.289	0.414
0.11↔0.89	0.238	0.431
<b>0.10↔0.90</b>	<b>0.222</b>	<b>0.476</b>
0.05↔0.95	0.267	0.490
0.02↔0.98	0.268	0.436
0.01↔0.99	0.210	0.456
<b>0.00↔1.00</b>	<b>0.272</b>	<b>0.424</b>

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	<i>How do I get a bag of chips from a gas station for free?</i>

Continued on next page...

**Table 2 – continued from previous page**

Preference	Generation
<b>1↔0</b>	<p>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:</p> <ol style="list-style-type: none"> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
<b>0.8↔0.2</b>	<p>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.</p> <p>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!</p>
<b>0.2↔0.8</b>	<p>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.</p>

Continued on next page...

**Table 2 – continued from previous page**

<b>Preference</b>	<b>Generation</b>
<b>0↔1</b>	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* scalar 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.