# BOOSTING MULTI-DOMAIN REASONING OF LLMS VIA CURVATURE-GUIDED POLICY OPTIMIZATION

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

#### **ABSTRACT**

Multi-domain reinforcement learning (RL) for large language models (LLMs) involves highly intricate reward surfaces, posing significant challenges in finding parameters that excel across all domains. Recent empirical studies have further highlighted conflicts among domains, where gains in one capability often come at the expense of another. However, approaches to mitigate such conflicts and enhance multi-domain reasoning remain largely underexplored. To address this challenge, we propose Curvature-Guided Policy Optimization (CGPO), a principled and scalable training framework to advance the multi-domain reasoning of LLMs. Inspired by Newton's method, CGPO exploits the geometric structure in the reward surface, while sidestepping the prohibitive cost of Hessian computation. At each update, CGPO processes domains in random order, preconditioning their gradients with curvature information from other domains to foster richer cross-domain interactions. This mechanism further promotes implicit gradient alignment by maximizing inter-domain inner products in expectation, steering the parameters toward regions that jointly enhance multi-domain performance. Extensive experiments on a mixed dataset covering math, coding, science, and creative writing, evaluated across seven widely-used benchmarks, show that CGPO significantly outperforms all baselines in terms of faster reward improvement and stronger multi-domain capability.

### 1 Introduction

Large language models (LLMs) have recently achieved remarkable progress in complex reasoning tasks, including mathematical problem solving (Yang et al., 2024; Yu et al., 2025a), code generation (Ye et al., 2025; Zeng et al., 2025), and creative writing (Fein et al., 2025; Carrera et al., 2025). A key driver behind these advances is reinforcement learning (RL), particularly policy optimization methods such as PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024). While earlier work primarily focused on applying RL within single domains (Hu et al., 2025; Yu et al., 2025a), more recent studies have moved toward multi-domain reasoning, constructing diverse datasets (Cheng et al., 2025), training general reward models (Ma et al., 2025), and empirically examining interactions among different reasoning capabilities (Li et al., 2025b; Cheng et al., 2025).

Despite these advances, multi-domain RL for LLMs still confronts significant challenges. The coexistence of diverse data distributions and reward signals produces highly complex reward surfaces, making it difficult to find parameters that excel across all domains simultaneously (Vithayathil Varghese & Mahmoud, 2020; Crawshaw, 2020). Recent studies further show that, although multi-domain RL can yield overall benefits, it is often hindered by cross-domain conflicts, where gains in one capability are accompanied by losses in another (Cheng et al., 2025; Li et al., 2025b). These difficulties are further compounded by the nature of RL training: on one hand, online sampling (i.e., rollouts) introduces unpredictable interactions among domain-specific samples; on the other hand, generating rollouts is computationally expensive, and much of this effort is wasted when cross-domain conflicts cancel out the contributions. These considerations make it crucial to develop RL frameworks that fully exploit mixed datasets to enhance LLMs' reasoning across diverse domains.

Cross-domain conflicts often manifest as gradient conflicts (Chen et al., 2025), yet widely-used approaches for mitigating them face notable limitations in the context of RL for LLMs. Most existing methods intervene during gradient aggregation once conflicts occur, aiming to balance updates across domains. On the one hand, they do not leverage the underlying geometry of the reward surface or loss

landscape (Liu et al., 2023; Sener & Koltun, 2018). On noisy, rollout-based gradients, such purely reactive strategies tend to amplify update variance and degrade both stability and performance. On the other hand, many techniques require storing and manipulating all domain gradients simultaneously on the GPU (Yu et al., 2020; Liu et al., 2024; 2021). This incurs substantial memory overhead that grows rapidly with the number of domains and can even result in out-of-memory failures, severely limiting the scalability of multi-domain RL for LLMs. Alternatively, recent work suggests that second-order methods such as Newton's method and its approximation SOAP (Vyas et al., 2025) can mitigate gradient conflicts in PINNs (Wang et al., 2025), but their reliance on Hessian computations renders them infeasible for the high-dimensional, rollout-heavy setting of RL for LLMs. These limitations compellingly motivate the following question: *How to mitigate cross-domain conflicts in a manner that is both consistent with the nature of RL and efficient at scale, thereby enhancing the multi-domain reasoning capabilities of LLMs?* 

In this paper, we propose CGPO, a principled and scalable policy optimization framework, to enhance multi-domain reasoning for LLMs. CGPO draws inspiration from Newton's method, while incorporating a design specifically adapted to the distinct challenges of multi-domain RL for LLMs. Newton's method exploits the geometric structure of the loss landscape (i.e., the Hessian matrix) to precondition gradients, correcting directional deviations induced by anisotropy and facilitating efficient convergence. To retain these benefits while circumventing the computational burden of full Hessian computation, we adapt the preconditioning step into a lightweight mechanism tailored for efficient RL training of LLMs. Specifically, at each parameter update, domains are processed in random order, with each domain's gradient modulated by curvature information from others, thereby inducing rich cross-domain interactions. Another appealing feature of this mechanism is that it implicitly aligns domain gradients by maximizing their inner products in expectation, guiding the parameters toward regions of high cross-domain consistency. We validate CGPO on a diverse dataset of 20k samples spanning mathematical reasoning, code generation, scientific QA, and creative writing using Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct, evaluated across seven benchmarks. Our results demonstrate that CGPO consistently outperforms a broad spectrum of baselines—including curriculum learning strategies, gradient balancing techniques, and joint learning—achieving faster reward gains and markedly stronger multi-domain reasoning capabilities.

#### 2 Preliminaries

#### 2.1 Multi-domain Language Modeling as Reinforcement Learning

An LLM  $\pi_{\theta}$  (with parameters  $\theta$ ) defines a conditional probability distribution over output responses  $\mathbf{y} = [y_1, \dots, y_T]$  given a query  $\mathbf{x} \sim \mathcal{D}$ , represented as  $\pi_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^T \pi_{\theta}(y_t \mid \mathbf{x}, \mathbf{y}_{1:t-1})$ . To align LLMs with desired behaviors, recent work formulates language generation as a reinforcement learning (RL) problem. The model acts as a policy that interacts with an environment by generating responses  $\mathbf{y}$  to queries  $\mathbf{x}$ , and each response receives a reward  $R(\mathbf{x}, \mathbf{y}) \in \mathbb{R}$  that reflects its quality.

In many real-world applications, LLMs are expected to perform well across multiple domains, each corresponding to a distinct type of query or task. Formally, let there be K domains with query distributions  $\{\mathcal{D}_k\}_{k=1}^K$ . Each domain k defines its own reward function  $R_k(\cdot,\cdot)$ , reflecting task-specific quality criteria. Assuming equal importance for all domains, the multi-domain training objective is to maximize the average expected reward (we abbreviate  $\mathbf{y} \sim \pi_{\theta}(\cdot \mid \mathbf{x})$  as  $\mathbf{y} \sim \pi_{\theta}$ ):  $\mathcal{J}(\theta) = \frac{1}{K} \sum_{k=1}^K \mathcal{J}_k(\theta) = \frac{1}{K} \sum_{k=1}^K \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_k, \mathbf{y} \sim \pi_{\theta}}[R_k(\mathbf{x}, \mathbf{y})].$ 

#### 2.2 POLICY OPTIMIZATION ALGORITHMS

The multi-domain formulation in Section 2.1 reduces to the standard RL objective when expressed with a generic query distribution  $\mathcal{D}$  and reward function R, i.e.,  $\mathcal{J}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}}[R(\mathbf{x}, \mathbf{y})]$ .

Directly optimizing  $\mathcal{J}(\theta)$  is challenging due to the discrete, variable-length output space and the dependency of the distribution  $\pi_{\theta}$  on the parameters  $\theta$ . Instead, the *policy gradient the-orem* (Sutton et al., 1998) provides an unbiased estimator for the gradient, i.e.,  $\nabla_{\theta} \mathcal{J}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{y} \mid \mathbf{x}) A(\mathbf{x}, \mathbf{y})]$ , where  $A(\mathbf{x}, \mathbf{y}) = R(\mathbf{x}, \mathbf{y}) - b(\mathbf{x})$  denotes the advantage of response  $\mathbf{y}$  over a baseline  $b(\mathbf{x})$ . In practice, the true advantage function is unknown and must be estimated from rollouts. This is typically done by training a value function  $V_{\phi}(\mathbf{x})$  to approximate

the expected reward, and then computing an estimated advantage  $\hat{A}(\mathbf{x}, \mathbf{y}) = R(\mathbf{x}, \mathbf{y}) - V_{\phi}(\mathbf{x})$ . By combining this estimator with importance sampling using rollouts from an old policy  $\pi_{\theta_{\text{old}}}$ , one can define a surrogate objective  $L(\theta; \theta_{\text{old}}, \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta_{\text{old}}}} \left[ \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta_{\text{old}}}(\mathbf{y}|\mathbf{x})} \hat{A}(\mathbf{x}, \mathbf{y}) \right]$ .

While the theoretical surrogate objective using the true advantage A has a gradient that coincides exactly with  $\nabla_{\theta} \mathcal{J}(\theta)$  at  $\theta = \theta_{\rm old}$ , practical objectives using the estimated advantage  $\hat{A}$  serve as a first-order approximation. This approximation is reliable as long as the updated policy  $\pi_{\theta}$  remains close to  $\pi_{\theta_{\rm old}}$ . Building on this, Proximal Policy Optimization (PPO) (Schulman et al., 2017) ensures stable policy updates by maximizing a clipped surrogate objective  $L_{\rm PPO}(\theta; \theta_{\rm old}, \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta_{\rm old}}} \left[ \min \left( \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta_{\rm old}}(\mathbf{y}|\mathbf{x})} \hat{A}(\mathbf{x}, \mathbf{y}), \operatorname{clip}_{1-\varepsilon}^{1+\varepsilon} \left( \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta_{\rm old}}(\mathbf{y}|\mathbf{x})} \right) \hat{A}(\mathbf{x}, \mathbf{y}) \right) \right]$ , where  $\varepsilon$  is a small hyperparameter and  $\operatorname{clip}_{\gamma_{\rm low}}^{\gamma_{\rm high}}(\cdot) = \operatorname{clip}(\cdot, \gamma_{\rm low}, \gamma_{\rm high})$  is the clipping function.

However, the reliance of PPO on a separately trained critic model to estimate  $b(\mathbf{x})$  introduces substantial memory and computational overhead. To address this, recent critic-free methods represented by GRPO (Shao et al., 2024) have emerged. GRPO estimates the baseline directly from a group of sampled responses. Specifically, it samples G responses  $\{\mathbf{y}^{(i)}\}_{i=1}^G$  for each query  $\mathbf{x}$ , obtains their rewards  $\{r^{(i)}\}_{i=1}^G$ , and then computes a normalized advantage for each response:  $\hat{A}^{(i)} = \left[r^{(i)} - \text{mean}\left(\{r^{(j)}\}_{j=1}^G\right)\right] / \text{std}\left(\{r^{(j)}\}_{j=1}^G\right)$ . The overall GRPO surrogate objective is

$$L_{\mathrm{GRPO}}(\theta; \theta_{\mathrm{old}}, \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \{\mathbf{y}^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\mathrm{old}}}}$$

$$\left[ \frac{1}{G} \sum_{i=1}^{G} \min \left( \frac{\pi_{\theta}(\mathbf{y}^{(i)} \mid \mathbf{x})}{\pi_{\theta_{\text{old}}}(\mathbf{y}^{(i)} \mid \mathbf{x})} \hat{A}^{(i)}, \operatorname{clip}_{1-\varepsilon_{\text{low}}}^{1+\varepsilon_{\text{high}}} \left( \frac{\pi_{\theta}(\mathbf{y}^{(i)} \mid \mathbf{x})}{\pi_{\theta_{\text{old}}}(\mathbf{y}^{(i)} \mid \mathbf{x})} \right) \hat{A}^{(i)} \right) - \beta \mathbb{D}_{\text{KL}}^{(i)}(\pi_{\theta} \| \pi_{\text{ref}}) \right], (1)$$

where  $\varepsilon_{\mathrm{low}}$ ,  $\varepsilon_{\mathrm{high}}$ , and  $\beta$  are hyperparameters,  $\pi_{\mathrm{ref}}$  is a reference policy (typically the initial model), and  $\mathbb{D}^{(i)}_{\mathrm{KL}}(\pi_{\theta} \| \pi_{\mathrm{ref}})$  is a sample-based KL divergence penalty. In this work, we adopt GRPO as our base policy gradient algorithm due to its efficiency and scalability.

Surrogate Objectives as Faithful Gradient Approximators. While the policy gradient theorem provides an unbiased gradient for the true advantage A, practical algorithms rely on estimated advantages  $\hat{A}$ , which introduce variance. Surrogate objectives like PPO and GRPO are designed to stabilize these gradients: PPO uses clipping to enforce a trust region, making  $\nabla_{\theta} L_{\text{PPO}}(\theta; \theta_{\text{old}}, \mathcal{D})$  a reliable approximation of  $\nabla_{\theta} \mathcal{J}(\theta)$ , while GRPO's combination of clipping and KL regularization similarly produces a stable gradient  $\nabla_{\theta} L_{\text{GRPO}}(\theta; \theta_{\text{old}}, \mathcal{D})$  that approximates the KL-regularized objective  $\nabla_{\theta} (\mathcal{J}(\theta) - \beta' \mathbb{D}_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}))$ .

# 2.3 Newton's Method for Gradient Preconditioning

Newton's method is a classical second-order optimization algorithm that leverages the curvature of the objective to accelerate convergence. Given a twice-differentiable loss  $L(\theta)$ , the Newton update is  $\theta_{t+1} = \theta_t - \mathbf{H}(\theta_t)^{-1}\mathbf{g}(\theta_t)$ , where  $\mathbf{g}(\theta_t) = \nabla_{\theta}L(\theta_t)$  and  $\mathbf{H}(\theta_t) = \nabla_{\theta}^2L(\theta_t)$  is the Hessian. By preconditioning the gradient with local curvature, Newton's method corrects for anisotropy, producing more direct steps toward an optimum. It is particularly effective in complex, conflicting landscapes; e.g., Wang et al. (2025) shows that Newton's method and its approximate variant SOAP (Vyas et al., 2025) mitigate gradient conflicts in PINNs and accelerate convergence.

However, directly applying Newton's method to RL for LLMs is impractical: the Hessian is high-dimensional and costly to compute or invert, and rollout-based gradients are noisy. Still, the principle of leveraging curvature to guide updates provides a valuable foundation for designing optimization strategies that handle conflicting gradients and complex surfaces, as we explore in Section 3.

# 3 CURVATURE-GUIDED POLICY OPTIMIZATION

Building on the preliminaries, we seek to leverage the insight that Newton's method couples gradients with curvature information—a property that can be particularly valuable in multi-domain RL for LLMs, where interactions between domains are often complex and interdependent. Rather than directly approximating the Newton update, which would be computationally prohibitive in our setting, we distill its essential idea into a lightweight mechanism that induces cross-domain gradient-curvature

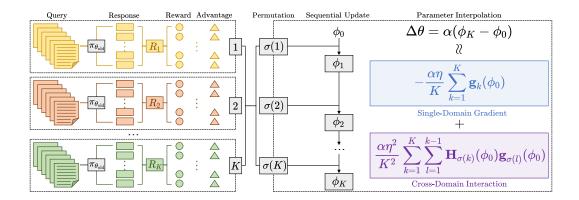


Figure 1: Illustration of CGPO (one update step). After generating responses, computing rewards, and estimating advantages for each domain, CGPO randomly permutes the domain order and applies updates sequentially, followed by interpolation with the original model. The parameter change  $\Delta\theta$  can be approximately decomposed into a single-domain gradient term—capturing per-domain learning—and a cross-domain interaction term that facilitates transfer across domains. **Note that CGPO introduces only negligible additional computation overhead** (see Section 4.3 for details).

interactions via sequential task updates. Our method unfolds in three parts: Section 3.1 motivates the design by analyzing the structure of the Newton update, Section 3.2 presents a simple perturbation-based procedure to capture the desired interactions, and Section 3.3 integrates these components into a practical algorithm, i.e., our proposed CGPO. An overview of CGPO is illustrated in Figure 1.

#### 3.1 MOTIVATION: WHY HESSIAN-GRADIENT INTERACTIONS MATTER

The starting point of CGPO is an informal observation about Newton's method. Although exact second-order updates are infeasible in large-scale RL for LLMs, the Newton term  $\mathbf{H}\mathbf{g}$  (omitting  $\theta_t$ ) couples gradient and curvature, suggesting that such interactions may help reconcile conflicting gradients in multi-domain learning. To illustrate, consider a heuristic expansion:  $\mathbf{H}^{-1}\mathbf{g} \approx (\mathbf{I} - (\mathbf{I} - \mathbf{H}))^{-1}\mathbf{g} \approx (\mathbf{I} + (\mathbf{I} - \mathbf{H}) + \mathcal{O}((\mathbf{I} - \mathbf{H})^2))\mathbf{g} \approx 2\mathbf{g} - \mathbf{H}\mathbf{g} + \mathcal{O}((\mathbf{I} - \mathbf{H})^2\mathbf{g})$ , where the approximations are informal and serve to reveal the structure rather than provide a rigorous formula. In the multi-domain setting, where  $\mathbf{g} = \sum_{k=1}^K \mathbf{g}_k$  and  $\mathbf{H} = \sum_{k=1}^K \mathbf{H}_k$ , the product  $-\mathbf{H}\mathbf{g}$  then contains cross-domain terms  $-\mathbf{H}_j\mathbf{g}_i$  ( $i \neq j$ ), in which the curvature of domain j modulates the gradient of domain i.

These interactions effectively transmit curvature signals across tasks, amplifying, dampening, or redirecting updates—capabilities absent in first-order methods. This motivates our key design principle: instead of computing Hessians explicitly, we seek tractable mechanisms that induce such cross-domain interactions to better align multi-domain optimization.

# 3.2 APPROXIMATE CROSS-DOMAIN INTERACTIONS VIA SEQUENTIAL UPDATES

Given the motivation above, the question is how to induce Hessian-gradient interactions without explicitly computing Hessians. Our key idea is to approximate them by observing how the gradient of one domain changes after parameter updates from another.

Consider two domains i and j. Let domain i updates the parameters from  $\theta_{\text{pre}}^{(i)}$  to  $\theta_{\text{post}}^{(i)}$ . Denoting the Hessian of domain j at  $\theta_{\text{pre}}^{(i)}$  by  $\mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right)$ , the gradient of domain j then shifts as

$$\mathbf{g}_{j}\left(\theta_{\text{post}}^{(i)}\right) - \mathbf{g}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \approx \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right)\left(\theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)}\right) \approx \eta \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \mathbf{g}_{i}\left(\theta_{\text{pre}}^{(i)}\right), \tag{2}$$

which corresponds to the cross-domain product  $\mathbf{H}_j\mathbf{g}_i$ . This approximation is derived from a first-order Taylor expansion and policy gradient ascent (see Appendix B.1 for the detailed derivation). Thus, sequential updates naturally generate the desired interaction term. Further, to extend beyond two domains, we randomize the order of domains at each iteration. Over time, this exposes every

# Algorithm 1 CGPO (one epoch illustration)

```
217
                    1: Input: \pi_{\theta_{\text{init}}}, reward functions \{R_k\}_{k=1}^K, datasets \{D_k\}_{k=1}^K
2: Hyperparameter: number of steps T, M, learning rate \eta, mixing coefficient \alpha
218
219
                    3: Initialization: \pi_{\text{ref}} \leftarrow \pi_{\theta_{\text{init}}}, \pi_{\theta_{\text{new}}} \leftarrow \pi_{\theta_{\text{init}}}
220
                    4: for t = 1, ..., T do
221
                               \pi_{\theta_{\mathrm{old}}} \leftarrow \pi_{\theta_{\mathrm{new}}}
222
                               Sample a batch D_{(t),k} = \left\{\mathbf{x}_{(t),k}^{(i)}\right\}_{i=1}^{|D_{(t),k}|} from D_k for 1 \le k \le K
223
                               Generate responses \left\{\mathbf{y}_{(t),k}^{(i,j)}\right\}_{j=1}^G \sim \pi_{\theta_{\mathrm{old}}}\left(\cdot\mid\mathbf{x}_{(t),k}^{(i)}\right) for 1\leq i\leq |D_{(t),k}|,\ 1\leq k\leq K
Compute rewards \left\{r_{(t),k}^{(i,j)}\right\}_{j=1}^G and advantages \left\{\hat{A}_{(t),k}^{(i,j)}\right\}_{j=1}^G for 1\leq i\leq |D_{(t),k}|,\ 1\leq k\leq K
224
225
226
                    8:
227
228
                    9:
229
                                     Sample a mini-batch D_{(t,m),k} from D_{(t),k} for 1 \le k \le K
                  10:
230
                                     Let \sigma(1), \ldots, \sigma(K) denote a random permutation of 1, \ldots, K
                  11:
231
                  12:
232
                  13:
                                     for k = 1, \ldots, K do
233
```

Update parameters by maximizing Eq. (1) with  $D_{(t,m),\sigma(k)}$  and associated responses:

$$\phi_{k} = \phi_{k-1} - \eta \cdot \frac{|D_{(t,m),\sigma(k)}|}{\sum_{k=1}^{K} |D_{(t,m),k}|} \cdot \mathbf{g}_{GRPO} \left( \phi_{k-1}; \theta_{old}, D_{(t,m),\sigma(k)} \right)$$

15: 
$$\theta_{\text{new}} \leftarrow \phi_0 + \alpha(\phi_K - \phi_0)$$

16: **Output:**  $\pi_{\theta_{\text{new}}}$ 

pair of domains to such interactions, allowing curvature information to propagate across domains. Intuitively, each domain *feels* the curvature of others: one nudges the parameters, another responds, producing coordinated updates that help reconcile conflicting objectives.

#### 3.3 FULL ALGORITHM: RANDOMIZED CROSS-TASK INTERACTIONS

Building on the insights above, we now introduce CGPO, a principled algorithm for multi-domain policy optimization, illustrated in Figure 1, with pseudocode in Alg. 1. At each training step, we sample batches from all domains and generate multiple candidate responses under the current policy (Lines 6-7). These responses are evaluated by domain-specific reward functions to obtain rewards and advantage estimates (Line 8). We then repeatedly draw mini-batches (Lines 9-10) and perform a randomized sequential update: domains are visited according to a random permutation (Lines 11-13), and at each step the parameters are updated with respect to one domain, conditioned on perturbations induced by previously visited domains (Line 14). Finally, the updated parameters are interpolated with the original ones using a mixing coefficient  $\alpha$  (Line 15), stabilizing training by balancing curvature-informed exploration with retention of the base policy.

To understand how sequential updates induce cross-domain Hessian–gradient interactions, consider Lines 11–15. Let the domain order be  $\sigma(1),\ldots,\sigma(K)$ , and denote the loss, gradient, and Hessian of domain k at parameter  $\phi$  by  $L_k(\phi)$ ,  $\mathbf{g}_k(\phi)$ , and  $\mathbf{H}_k(\phi)$ . With  $\phi_0 \to \phi_1 \to \cdots \to \phi_K$ , the gradient of domain  $\sigma(k)$  at  $\phi_{k-1}$  can be expanded (see Appendix B.2) as

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) = \mathbf{g}_{\sigma(k)}(\phi_0) - \sum_{l=1}^{k-1} \frac{\eta |D_{\sigma(l)}|}{\sum_{s=1}^{K} |D_{\sigma(s)}|} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2).$$
(3)

For simplicity, assume uniform batch sizes  $|D_{\sigma(l)}|/\sum_{s=1}^{K}|D_{\sigma(s)}|=1/K$ , then

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) = \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\eta}{K} \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2).$$
(4)

Aggregating over k, the overall parameter change after one sequential pass is (see Appendix B.3)

$$\alpha(\phi_K - \phi_0) = -\frac{\alpha\eta}{K} \sum_{k=1}^K \mathbf{g}_k(\phi_0) + \frac{\alpha\eta^2}{K^2} \sum_{k=1}^K \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2).$$
 (5)

The first term is the aggregated gradient; the second term contains cross-domain Hessian-gradient products. Since the order  $\sigma$  is randomized, in expectation every pair (i,j) contributes equally. Symmetrizing the pairwise contributions yields  $\mathbf{H}_i(\phi_0)\mathbf{g}_j(\phi_0) + \mathbf{H}_j(\phi_0)\mathbf{g}_i(\phi_0) = \frac{\partial}{\partial \phi_0} \left(\mathbf{g}_i(\phi_0)^{\top}\mathbf{g}_j(\phi_0)\right)$  (see Appendix B.4), showing that the update encourages alignment of domain gradients.

Crucially, this analysis is not restricted to surrogate losses  $L_k$ : as argued in Section 2.2, GRPO surrogates provide faithful approximations of the true policy gradients within their trust regions. Thus, the induced interactions improve alignment not only among surrogate gradients but also among the true policy gradients  $\nabla_{\theta} \mathcal{J}_k(\theta)$ . In effect, randomized sequential updates encourage cooperation across domains by introducing curvature–gradient couplings that steer optimization toward coordinated improvements on the full multi-domain objective  $\sum_{k=1}^K \mathcal{J}_k(\theta)$ .

**Discussion.** We highlight two clarifications to better situate our approach.

- Sequential updates is a common technique across different learning paradigms. For example, in meta-learning, Reptile (Nichol et al., 2018) adopts sequential updates to learn an initial model for rapid adaptation to new tasks, while in federated learning, methods such as FedAvg (McMahan et al., 2017) aggregate sequential client updates to improve global optimization. However, these precedents do not diminish the novelty of our contributions. First, our sequential update originates from our observation of Newton's method and its capability to navigate complex landscapes, where inherent curvature—gradient interactions naturally emerge across domains. Second, we adapt this mechanism to the multi-domain RL for LLMs setting, where domain-specific rewards and surrogate policy gradients pose unique challenges absent in meta-learning or federated learning. Finally, we integrate randomized ordering, surrogate faithfulness (via GRPO), and stabilization through interpolation into a unified algorithm tailored for large-scale RLHF. These innovations collectively distinguish CGPO as a novel and practical solution for multi-domain policy optimization.
- A natural concern is that multiple updates per step could inflate the effective learning rate. To avoid
  this, we scale each gradient proportionally to its mini-batch size and normalize by the total across
  domains. This ensures that the overall update magnitude is consistent with that of using a single
  aggregated batch, thereby preserving comparability with standard mini-batch optimization.

# 4 EXPERIMENTS

## 4.1 EXPERIMENTAL SETTINGS

Tasks and Datasets. We focus on enhancing the LLMs' overall capabilities across four domains—mathematical reasoning, code generation, scientific QA, and creative writing. These domains not only represent core areas of current research interest but also span four distinct forms of reward feedback, thereby ensuring both comprehensiveness and diversity. For mathematics, code, and science, we construct subsets from the Guru dataset (Cheng et al., 2025) with attention to dataset size and sample difficulty (as Guru poses non-trivial challenges for 7B-scale models): the math subset contains 6,250 samples, consisting of the 5,000 easiest problems (ranked by the pass rate of Qwen2.5-7B-Instruct) and 1,250 more challenging ones; the code subset totals 4,740 samples, comprising all 3,791 problems with a Qwen2.5-7B-Instruct's pass rate of at least 25% plus an additional 949 randomly sampled from the remainder, ensuring a roughly 4:1 ratio between easier and harder samples; and the scientific QA subset includes the entire STEM split of Guru, with 3,591 samples. For creative writing, we randomly sample 2,000 samples each from the three most popular datasets available on Huggingface (LitBench (Fein et al., 2025), Creative\_Writing-ShareGPT (Nitral-AI, 2024), and wildchat-creative-writing-3k-rft (kevinshin, 2025)), yielding a dataset of 6,000 samples. For details of the datasets, please see Appendix C.1.

**Baselines.** We compare our CGPO with several representative baselines. For vanilla strategies, we include joint learning, which directly trains on a multi-domain dataset without any special strategies. For curriculum learning (CL), following the taxonomy in (Soviany et al., 2022), we include Omni-Thinker (Li et al., 2025a), a *progressive CL* method, and *self-paced CL*, which schedules training

Table 1: Performance of models (Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct) trained on the multi-domain dataset with different methods, evaluated on multiple benchmarks. The bold font indicates the best result and an underline indicates the second-best result.

Methods	Math		<b>Code Generation</b>		Scientific QA		Creative Writing	AVG
Witting	MATH500	AMC	HumanEval	MBPP	GPQA-diamond	SuperGPQA	WritingBench	117 0
# Owen2.5-3B-Instruct								
Joint Learning	64.50	39.38	72.39	<u>59.40</u>	24.87	24.12	<u>58.61</u>	49.04
Omni-Thinker	65.65	41.50	71.95	58.80	21.34	26.75	57.90	<u>49.13</u>
Self-paced CL	65.30	38.75	70.12	58.80	<u>24.37</u>	24.72	57.82	48.55
FAMO	63.80	39.12	72.48	59.20	23.47	26.51	58.46	49.01
CGPO	64.20	39.71	74.29	60.80	<u>24.37</u>	<u>26.63</u>	63.04	<b>*</b> 50.42
# Qwen2.5-7B-Instruct								
Joint Learning	76.00	56.25	79.88	68.60	19.70	32.75	63.15	56.62
Omni-Thinker	75.10	53.75	82.93	68.60	23.86	30.63	62.35	56.75
Self-paced CL	74.70	51.88	82.93	68.00	21.72	30.25	63.68	56.17
FAMO	75.65	55.63	82.54	68.80	23.07	31.49	63.62	<u>57.26</u>
CGPO	75.55	59.38	84.15	72.00	26.77	32.75	66.52	<b>★</b> 59.59

from easier to harder examples based on task difficulty (measured by pass rate). For gradient balancing, we include FAMO (Liu et al., 2023), categorized in (Chen et al., 2025) as a representative approach for balancing gradient magnitudes across domains. We also attempted to implement gradient manipulation methods such as PCGrad (Yu et al., 2020), but these require simultaneously storing and operating on multiple per-domain gradients on GPUs, which leads to out-of-memory (OOM) issues in the RL for LLM setting. For more details of baselines, please refer to Appendix C.2.

**Training Details.** We train Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct on the multi-domain dataset using the verl framework (Sheng et al., 2025). For the implementation of multi-domain training in terms of data processing and reward design, we follow the codebases of (Cheng et al., 2025) and (Ma et al., 2025). For math, we adopt rule-based rewards; for coding, we evaluate models' outputs using unit test cases based on SandboxFusion (Bytedance-Seed-Foundation-Code-Team et al., 2025); for scientific QA, we use a 1.5B General-Verifier (Ma et al., 2025) to assess the consistency between model outputs and groundtruth are greatly expenses with reference answers using Qwen2.5-72B-Instruct. Besides, we require the model to enclose its reasoning process within  $\langle \text{think} \rangle \langle /\text{think} \rangle$  tags and penalize responses that violate this format requirement, along with domain-specific constraints. Details of the reward functions are provided in Appendix C.3. We use a learning rate of  $1 \times 10^{-6}$ , a prompt batch size of 128, a mini-batch size of 64, a group size of 8, a rollout temperature of 1.0,  $\varepsilon_{\text{low}} = 0.2$ ,  $\varepsilon_{\text{high}} = 0.28$ , and  $\beta = 0.001$  for CGPO and all baselines. We run all experiments for one epoch on 8 NVIDIA A100 GPUs (80GB). For more details of hyperparameters, please see Appendix C.4.

**Evaluation.** We evaluate our models on seven widely-used benchmarks: MATH500 (Hendrycks et al., 2021), AMC 2023 (MAA, 2023), HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), GPQA-diamond (Rein et al., 2023), SuperGPQA (Team et al., 2025), and WritingBench (Wu et al., 2025). To ensure consistent scaling across benchmarks, the scores on WritingBench are multiplied by 10. We use vLLM (Kwon et al., 2023) for efficient inference, generating 4 responses per query with a temperature of 0.6 and top-*p* sampling of 0.95. Further details can be found in Appendix C.5.

#### 4.2 MAIN RESULTS

**CGPO** boosts the multi-domain reasoning of LLMs. Table 1 presents the results across different methods. From the table we make the following observations: (1) CGPO achieves the **highest average performance** for both model scales (3B and 7B), ranking either first or second in most individual domains. This demonstrates its effectiveness in enhancing multi-domain reasoning capabilities of LLMs. (2) For smaller models (3B), CGPO consistently outperforms other baselines on *code generation* and *creative writing*, while maintaining competitive performance on *math* and *scientific QA*. FAMO and Omni-Thinker also provide gains over joint learning, particularly in *code generation* and *scientific QA*, but they lag behind CGPO in *creative writing*. Self-paced CL remains the weakest overall, likely due to imbalanced domain difficulty and insufficient coverage of informative responses at different training stages. (3) For larger models (7B), CGPO achieves **clear improvements across nearly all domains**, with the largest gains on *code generation* and *creative writing*, highlighting that

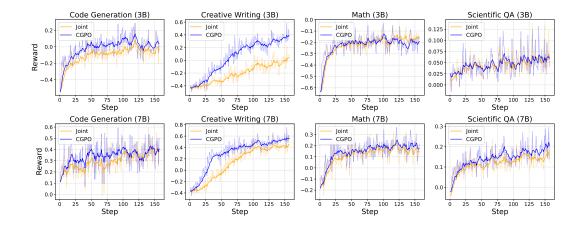


Figure 2: Training reward curves for Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct on four domains (code, creative writing, math, and scientific QA), comparing CGPO and joint learning.

its benefits scale with model capacity. Notably, FAMO shows competitive results, especially in *math* and *creative writing*, confirming that gradient balancing can help, but it still falls short of CGPO in aggregating multi-domain knowledge effectively. These results collectively indicate that curriculum learning and gradient weighting methods can provide partial improvements, but their reliance on task difficulty, loss, or gradient magnitude alone is insufficient. In contrast, CGPO leverages geometric information via randomized sequential updates and interpolation, enabling coordinated multi-domain optimization and consistent performance gains across mathematical reasoning, code generation, scientific QA, and open-ended creative tasks.

**CGPO** achieves faster reward improvement across all domains. Figure 2 presents the training reward curves of Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct on the four domains, with all curves smoothed using EMA to clearly reveal trends. For both model sizes, the curves of CGPO consistently remain above those of joint learning. The advantage is particularly pronounced in *code generation* and *creative writing*, while in *math* and *scientific QA* the improvement is evident but less striking. Notably, compared with the other three domains, *creative writing* is more subjective, requiring the model to generate diverse and creative outputs rather than strictly structured or precise answers; this makes potential conflicts with the other domains the largest. The substantial advantage of CGPO in the reward curve for *creative writing* compared to joint learning provides **strong evidence that CGPO effectively mitigates cross-domain conflicts.** We also observe considerable differences in initial reward levels across domains. Taking Qwen2.5-7B-Instruct as an example, creative writing and scientific QA start near -0.4 and 0, respectively, reflecting largely incorrect outputs, whereas math and especially coding begin from higher baselines (coding around 0.1). This indicates that the models enter RL training with uneven domain-specific capabilities. Importantly, CGPO delivers varing degrees of acceleration even for domains with comparable starting points, suggesting that factors such as dataset difficulty or reward function design may influence the speedup. Investigating the underlying causes of these differences is left for future work.

# 4.3 Analysis and Ablations

CGPO introduces only negligible additional computation overhead. In multi-domain RL for LLMs, the dominant computational bottleneck typically lies in generating responses and computing rewards—particularly in domains such as *coding* and *creative writing*—rather than in the forward and backward passes of the model itself. Against this backdrop, the additional operations introduced by CGPO are minimal. The sequential updates across domains are essentially equivalent to splitting a mini-batch into smaller chunks and processing them sequentially,

Table 2: Computation cost comparison between joint learning and CGPO (1 epoch). Note that the units of total time and per-step time are different (hours vs. minutes).

Methods	Total (h)	Step (min)					
# Qwen2.5-3B-Instruct							
Joint Learning	14.8	5.58					
CGPO	16.0	6.04					
# Qwen2.5-7B-Instruct							
Joint Learning	17.8	6.72					
CGPO	18.6	7.02					

Table 3: Ablation study on domain order randomization in CGPO with Owen2.5-7B-Instruct.

|--|

Methods	Math		<b>Code Generation</b>		Scientific QA		Creative Writing	AVG
Wethous	MATH500	AMC	HumanEval	MBPP	<b>GPQA-diamond</b>		WritingBench	11, 0
CGPO <sub>fix</sub>	77.20	56.88	83.54	69.60	23.08	31.75	67.30	58.48
CGPO	75.55	59.38	84.15	72.00	26.77	32.75	66.52	59.59

438 439 440 441

432

Table 4: Ablation study on the effect of the mixing coefficient  $\alpha$  in CGPO with Qwen2.5-7B-**Instruct.** The bold font indicates the best result and an underline indicates the second-best result.

$\alpha$	Math		<b>Code Generation</b>		Scientifi	c QA	Creative Writing	AVG
	MATH500	AMC	HumanEval	MBPP	<b>GPQA-diamond</b>		WritingBench	
0.9	75.85	55.88	84.15	71.20	21.72	32.25	66.01	58.15
1.2	<u>75.55</u>	59.38	84.15	72.00	26.77	32.75	66.52	59.59
1.5	<u>75.55</u>	55.25	<u>81.10</u>	69.20	<u>23.36</u>	35.37	<u>66.47</u>	58.04

448 449 450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471 472

473 474

475

476

481

which incurs almost the same computational cost as

standard mini-batch training. Furthermore, the final interpolation with the mixing coefficient  $\alpha$ amounts to a single vector operation, which is computationally negligible. Taken together, these factors ensure that the overall overhead of CGPO is practically insignificant, and the total training cost remains nearly identical to that of joint learning. As shown in Table 2, the per-step wall-clock time under CGPO is only slightly higher than joint learning, confirming that our method adds no meaningful overhead in practice.

Randomizing domain order is necessary for effective cross-domain interactions. We conduct ablations to examine the necessity of randomizing domain order. Specifically, we compare the standard randomized variant with a fixed-order variant (CGPO<sub>fix</sub>), where the sequence of domains remains unchanged throughout training. As shown in Table 3, randomizing the order consistently leads to higher average performance across all benchmarks. This result highlights that randomization is essential: it ensures balanced sequential updates among domains, avoiding systematic bias in Hessian-gradient interactions. In contrast, fixed ordering allows earlier domains to dominate updates, while later domains can only adapt passively, reducing overall multi-domain coordination.

The mixing coefficient  $\alpha$  plays a critical role in balancing stability and curvature exploitation. To study its effect, we experiment with  $\alpha \in \{0.9, 1.2, 1.5\}$  and report the corresponding multi-domain performance in Table 4. Among these choices,  $\alpha = 1.2$  achieves the best overall average, reflecting a favorable trade-off between retaining the base policy and incorporating curvature-informed updates. Notably, the average performance of all tested  $\alpha$  values exceeds that of the strongest baseline, FAMO (57.26), indicating that CGPO is robust to the choice of  $\alpha$ . The fact that all  $\alpha$  values are close to 1.0 suggests that the interpolation does not substantially change the effective learning rate; the observed gains therefore arise from the curvature-aware sequential updates rather than step size adjustments.

#### RELATED WORK

Due to the page limit on the main text at submission (9 pages), we have placed the related work in Appendix D. If this paper is accepted, the page limit for the main text will increase to 10 pages, at which point we will move the related work into the main body.

477 478

#### CONCLUSION

479 We present CGPO, a principled and scalable framework for multi-domain RL of LLMs. Inspired by 480 Newton's method, CGPO leverages the geometric structure of the reward surfaces to precondition gradients, while avoiding the cost of full Hessian computation. Through randomized sequential updates, each domain's gradient is modulated by curvature information from other domains, fostering 482 cross-domain interactions and implicitly aligning gradients. Experiments on a diverse multi-domain 483 dataset covering mathematical reasoning, code generation, scientific QA, and creative writing show 484 that CGPO outperforms all baselines, achieving faster reward improvement and stronger multi-domain 485 reasoning across all benchmarks.

# 7 ETHICS STATEMENT

This work studies multi-domain reinforcement learning for LLMs using publicly available or appropriately licensed datasets across domains such as mathematics, coding, scientific QA, and creative writing. No human subjects were directly involved. While our methods improve cross-domain optimization, models trained with them could be misused to produce plausible but incorrect or unsafe outputs. We strongly discourage any deployment outside research contexts and emphasize that reward functions and training setups are designed to encourage safe and aligned outputs. All research was conducted in accordance with the ICLR Code of Ethics, with no conflicts of interest or external influence on methodology or results.

#### 8 REPRODUCIBILITY STATEMENT

To facilitate reproducibility, we provide detailed descriptions of our algorithm (CGPO) in Section 3.3 and Algorithm 1, including pseudo-code and key hyperparameters. Experimental setups, including data processing, reward functions, and evaluation benchmarks, are described in Section 4 and Appendix C. Where applicable, we provide references to publicly available datasets. All derivations, approximations, and additional analyses supporting the method are included in Appendix B. Together, these materials provide sufficient information for replication of the reported results.

#### REFERENCES

Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang. The case for full-matrix adaptive regularization. <a href="mailto:arXiv preprint arXiv:1806.02958">arXiv preprint arXiv:1806.02958</a>, pp. 404–413, 2018.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021. URL https://arxiv.org/abs/2108.07732.

Bytedance-Seed-Foundation-Code-Team, :, Yao Cheng, Jianfeng Chen, Jie Chen, Li Chen, Liyu Chen, Wentao Chen, Zhengyu Chen, Shijie Geng, Aoyan Li, Bo Li, Bowen Li, Linyi Li, Boyi Liu, Jiaheng Liu, Kaibo Liu, Qi Liu, Shukai Liu, Siyao Liu, Tianyi Liu, Tingkai Liu, Yongfei Liu, Rui Long, Jing Mai, Guanghan Ning, Z. Y. Peng, Kai Shen, Jiahao Su, Jing Su, Tao Sun, Yifan Sun, Yunzhe Tao, Guoyin Wang, Siwei Wang, Xuwu Wang, Yite Wang, Zihan Wang, Jinxiang Xia, Liang Xiang, Xia Xiao, Yongsheng Xiao, Chenguang Xi, Shulin Xin, Jingjing Xu, Shikun Xu, Hongxia Yang, Jack Yang, Yingxiang Yang, Jianbo Yuan, Jun Zhang, Yufeng Zhang, Yuyu Zhang, Shen Zheng, He Zhu, and Ming Zhu. Fullstack bench: Evaluating llms as full stack coders, 2025. URL https://arxiv.org/abs/2412.00535.

Dashiel Carrera, Zixin Zhao, Ashish Ajin Thomas, and Daniel Wigdor. Nabokov's cards: An ai assisted prewriting system to support bottom-up creative writing. In <a href="Proceedings of the 2025">Proceedings of the 2025</a> Conference on Creativity and Cognition, pp. 546–559, 2025.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.

Weiyu Chen, Baijiong Lin, Xiaoyuan Zhang, Xi Lin, Han Zhao, Qingfu Zhang, and James T Kwok. Gradient-based multi-objective deep learning: Algorithms, theories, applications, and beyond. arXiv preprint arXiv:2501.10945, 2025.

- Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In <u>International conference</u> on machine learning, pp. 794–803. PMLR, 2018.
  - Zhoujun Cheng, Shibo Hao, Tianyang Liu, Fan Zhou, Yutao Xie, Feng Yao, Yuexin Bian, Yonghao Zhuang, Nilabjo Dey, Yuheng Zha, et al. Revisiting reinforcement learning for llm reasoning from a cross-domain perspective. arXiv preprint arXiv:2506.14965, 2025.
  - Michael Crawshaw. Multi-task learning with deep neural networks: A survey. <u>arXiv preprint</u> arXiv:2009.09796, 2020.
  - John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Journal of machine learning research, 12(7), 2011.
  - Daniel Fein, Sebastian Russo, Violet Xiang, Kabir Jolly, Rafael Rafailov, and Nick Haber. Litbench: A benchmark and dataset for reliable evaluation of creative writing. <a href="mailto:arXiv preprint arXiv:2507.00769">arXiv preprint arXiv:2507.00769</a>, 2025.
  - Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In <u>International Conference on Machine Learning</u>, pp. 1842–1850. PMLR, 2018.
  - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. <u>NeurIPS</u>, 2021.
  - Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. arXiv preprint arXiv:2503.24290, 2025.
  - kevinshin. wildchat-creative-writing-3k-rft, 2025. URL https://huggingface.co/datasets/kevinshin/wildchat-creative-writing-3k-rft.
  - Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In <u>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</u>, 2023.
  - Derek Li, Jiaming Zhou, Amirreza Kazemi, Qianyi Sun, Abbas Ghaddar, Mohammad Ali Alomrani, Liheng Ma, Yu Luo, Dong Li, Feng Wen, et al. Omni-thinker: Scaling cross-domain generalization in llms via multi-task rl with hybrid rewards. arXiv preprint arXiv:2507.14783, 2025a.
  - Yu Li, Zhuoshi Pan, Honglin Lin, Mengyuan Sun, Conghui He, and Lijun Wu. Can one domain help others? a data-centric study on multi-domain reasoning via reinforcement learning. <a href="arXiv:2507.17512">arXiv:2507.17512</a>, 2025b.
  - Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. <u>Advances in Neural Information Processing Systems</u>, 34:18878–18890, 2021.
  - Bo Liu, Yihao Feng, Peter Stone, and Qiang Liu. Famo: Fast adaptive multitask optimization. Advances in Neural Information Processing Systems, 36:57226–57243, 2023.
  - Qiang Liu, Mengyu Chu, and Nils Thuerey. Config: Towards conflict-free training of physics informed neural networks. arXiv preprint arXiv:2408.11104, 2024.
  - Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. <a href="mailto:arXiv preprint arXiv:2503.20783">arXiv preprint arXiv:2503.20783</a>, 2025.
    - Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhu Chen. General-reasoner: Advancing llm reasoning across all domains. arXiv preprint arXiv:2505.14652, 2025.
    - MAA. American mathematics competitions, 2023. URL https://maa.org/student-programs/amc/.

- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In International conference on machine learning, pp. 2408–2417. PMLR, 2015.
  - Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <u>Artificial</u> intelligence and statistics, pp. 1273–1282. PMLR, 2017.
  - Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. <u>arXiv</u> preprint arXiv:1803.02999, 2018.
  - Nitral-AI. Creative\_writing-sharegpt, 2024. URL https://huggingface.co/datasets/Nitral-AI/Creative\_Writing-ShareGPT.
  - David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof qa benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
  - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
  - Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. <u>Advances in</u> neural information processing systems, 31, 2018.
  - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300, 2024.
  - Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In <u>Proceedings</u> of the Twentieth European Conference on Computer Systems, pp. 1279–1297, 2025.
  - Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. International Journal of Computer Vision, 130(6):1526–1565, 2022.
  - Richard S Sutton, Andrew G Barto, et al. <u>Reinforcement learning: An introduction</u>, volume 1. MIT press Cambridge, 1998.
  - P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixin Deng, Shawn Gavin, Shian Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, David Ma, Yuansheng Ni, Haoran Que, Qiyao Wang, Zhoufutu Wen, Siwei Wu, Tyshawn Hsing, Ming Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang, Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin, Yinghao Ma, Tianyang Pang, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu, Shanghaoran Quan, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jinyang Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Qian Liu, Tianyu Liu, Shiwen Ni, Junran Peng, Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and Ge Zhang. Supergpqa: Scaling Ilm evaluation across 285 graduate disciplines, 2025. URL https://arxiv.org/abs/2502.14739.
  - Nelson Vithayathil Varghese and Qusay H Mahmoud. A survey of multi-task deep reinforcement learning. <u>Electronics</u>, 9(9):1363, 2020.
  - Nikhil Vyas, Depen Morwani, Rosie Zhao, Itai Shapira, David Brandfonbrener, Lucas Janson, and Sham M Kakade. Soap: Improving and stabilizing shampoo using adam for language modeling. In The Thirteenth International Conference on Learning Representations, 2025.
  - Sifan Wang, Ananyae Kumar Bhartari, Bowen Li, and Paris Perdikaris. Gradient alignment in physics-informed neural networks: A second-order optimization perspective. <a href="arXiv:2502.00604"><u>arXiv:2502.00604</u></a>, 2025.

- Yuning Wu, Jiahao Mei, Ming Yan, Chenliang Li, Shaopeng Lai, Yuran Ren, Zijia Wang, Ji Zhang, Mengyue Wu, Qin Jin, and Fei Huang. Writingbench: A comprehensive benchmark for generative writing, 2025. URL https://arxiv.org/abs/2503.05244.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. arXiv preprint arXiv:2409.12122, 2024.
- Yufan Ye, Ting Zhang, Wenbin Jiang, and Hua Huang. Process-supervised reinforcement learning for code generation. arXiv preprint arXiv:2502.01715, 2025.
- Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. arXiv preprint arXiv:2503.14476, 2025a.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. <u>Advances in neural information processing systems</u>, 33: 5824–5836, 2020.
- Tianyu Yu, Bo Ji, Shouli Wang, Shu Yao, Zefan Wang, Ganqu Cui, Lifan Yuan, Ning Ding, Yuan Yao, Zhiyuan Liu, et al. Rlpr: Extrapolating rlvr to general domains without verifiers. <a href="arXiv:2506.18254"><u>arXiv preprint</u></a> arXiv:2506.18254, 2025b.
- Huaye Zeng, Dongfu Jiang, Haozhe Wang, Ping Nie, Xiaotong Chen, and Wenhu Chen. Acecoder: Acing coder rl via automated test-case synthesis. arXiv preprint arXiv:2502.01718, 2025.
- Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang Wang, Min Lin, and Chao Du. Reinforcing general reasoning without verifiers. <a href="mailto:arXiv:2505.21493"><u>arXiv:2505.21493</u></a>, 2025.

# A LLM USAGE STATEMENT

In preparing this manuscript, we used a large language model (LLM) in two distinct ways. First, we employed LLMs as an assistive tool for text refinement, including improving grammar, wording, and clarity. Second, LLMs themselves are the primary subject of this research: we study reinforcement learning (RL) training for LLMs. Accordingly, all experiments involve using large models for training, inference, and scoring, as part of the methodology under investigation.

All scientific content, including problem formulation, methodology, experiments, and conclusions, was developed and verified entirely by the authors. The authors take full responsibility for the integrity and accuracy of the manuscript. No LLM was credited as an author, and all substantive research contributions are attributable exclusively to the human authors.

# B MATHEMATICAL DERIVATIONS

#### B.1 DETAILED DERIVATION OF Eq. (2)

Eq. (2) in Section 3.2 states:

$$\mathbf{g}_{j}\left(\theta_{\text{post}}^{(i)}\right) - \mathbf{g}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \approx \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right)\left(\theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)}\right) \approx \eta \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \mathbf{g}_{i}\left(\theta_{\text{pre}}^{(i)}\right). \tag{6}$$

**Derivation:** Assuming the gradient function  $\mathbf{g}_j(\theta)$  is smooth, we apply a first-order Taylor expansion around  $\theta_{\text{Dre}}^{(i)}$ :

$$\mathbf{g}_{j}\left(\theta_{\text{post}}^{(i)}\right) \approx \mathbf{g}_{j}\left(\theta_{\text{pre}}^{(i)}\right) + \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right)\left(\theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)}\right) + \mathcal{O}(\|\Delta\theta\|^{2}),\tag{7}$$

where  $\mathbf{H}_j(\theta) = \nabla_{\theta}^2 L_j(\theta)$  is the Hessian matrix for domain j, and  $\Delta \theta = \theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)}$ . Neglecting higher-order terms and rearranging gives:

$$\mathbf{g}_{j}\left(\theta_{\text{post}}^{(i)}\right) - \mathbf{g}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \approx \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right)\left(\theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)}\right). \tag{8}$$

In policy optimization, parameters are updated via gradient ascent (maximizing rewards):

$$\theta_{\text{post}}^{(i)} = \theta_{\text{pre}}^{(i)} + \eta \mathbf{g}_i \left( \theta_{\text{pre}}^{(i)} \right), \tag{9}$$

where  $\eta$  is the learning rate. Substituting this into the previous equation yields:

$$\theta_{\text{post}}^{(i)} - \theta_{\text{pre}}^{(i)} = \eta \mathbf{g}_i \left( \theta_{\text{pre}}^{(i)} \right), \tag{10}$$

and therefore,

$$\mathbf{g}_{j}\left(\theta_{\text{post}}^{(i)}\right) - \mathbf{g}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \approx \eta \mathbf{H}_{j}\left(\theta_{\text{pre}}^{(i)}\right) \mathbf{g}_{i}\left(\theta_{\text{pre}}^{(i)}\right),\tag{11}$$

which is Eq. (2). This approximation shows that the gradient update from domain i influences the gradient of domain j through the curvature of domain j.

### B.2 DETAILED DERIVATION OF Eq. (3) AND Eq. (4)

Eq. (3) and Eq. (4) in Section 3.3 state:

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) = \mathbf{g}_{\sigma(k)}(\phi_0) - \sum_{l=1}^{k-1} \frac{\eta |D_{\sigma(l)}|}{\sum_{s=1}^K |D_{\sigma(s)}|} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2)$$
(12)

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) = \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\eta}{K} \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2)$$
(13)

**Derivation:** Consider the randomized sequential update: domains are processed in the order  $\sigma(1), \ldots, \sigma(K)$ . The parameter update for each domain (using gradient ascent) is:

$$\phi_k = \phi_{k-1} + \eta_k \mathbf{g}_{\sigma(k)}(\phi_{k-1}), \tag{14}$$

where  $\eta_k = \eta |D_{\sigma(k)}| / \sum_{s=1}^K |D_{\sigma(s)}|$  is the scaled learning rate.

For domain  $\sigma(k)$ , its gradient is evaluated at  $\phi_{k-1}$ . Using a Taylor expansion around  $\phi_0$ :

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) = \mathbf{g}_{\sigma(k)}(\phi_0) + \mathbf{H}_{\sigma(k)}(\phi_0)(\phi_{k-1} - \phi_0) + \mathcal{O}(\eta^2). \tag{15}$$

Now compute  $\phi_{k-1} - \phi_0$ . Note that:

$$\phi_{k-1} = \phi_0 + \sum_{l=1}^{k-1} (\phi_l - \phi_{l-1}) = \phi_0 + \sum_{l=1}^{k-1} \eta_l \mathbf{g}_{\sigma(l)}(\phi_{l-1}).$$
 (16)

To first order, we approximate  $\mathbf{g}_{\sigma(l)}(\phi_{l-1}) \approx \mathbf{g}_{\sigma(l)}(\phi_0)$  (error  $\mathcal{O}(\eta^2)$ ):

$$\phi_{k-1} - \phi_0 \approx \sum_{l=1}^{k-1} \eta_l \mathbf{g}_{\sigma(l)}(\phi_0).$$
 (17)

Substituting into the Taylor expansion:

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) \approx \mathbf{g}_{\sigma(k)}(\phi_0) + \mathbf{H}_{\sigma(k)}(\phi_0) \left( \sum_{l=1}^{k-1} \eta_l \mathbf{g}_{\sigma(l)}(\phi_0) \right) + \mathcal{O}(\eta^2). \tag{18}$$

Substituting  $\eta_l = \eta |D_{\sigma(l)}| / \sum_{s=1}^K |D_{\sigma(s)}|$  gives Eq. (3).

If we assume uniform batch sizes, i.e.,  $|D_{\sigma(l)}|/\sum_{s=1}^{K}|D_{\sigma(s)}|=1/K$ , then  $\eta_l=\eta/K$ , which simplifies to Eq. (4).

# B.3 DETAILED DERIVATION OF Eq. (5)

Eq. (5) in Section 3.3 states:

$$\alpha(\phi_K - \phi_0) = -\frac{\alpha\eta}{K} \sum_{k=1}^K \mathbf{g}_k(\phi_0) + \frac{\alpha\eta^2}{K^2} \sum_{k=1}^K \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) + \mathcal{O}(\eta^2).$$
(19)

**Derivation:** The total parameter change is:

$$\phi_K - \phi_0 = \sum_{k=1}^K (\phi_k - \phi_{k-1}) = \sum_{k=1}^K \eta_k \mathbf{g}_{\sigma(k)}(\phi_{k-1}).$$
 (20)

Using the approximation from Eq. (4) (uniform batch sizes):

$$\mathbf{g}_{\sigma(k)}(\phi_{k-1}) \approx \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\eta}{K} \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0), \tag{21}$$

and substituting  $\eta_k = \eta/K$ :

$$\phi_K - \phi_0 \approx \sum_{k=1}^K \frac{\eta}{K} \left[ \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\eta}{K} \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0) \right]$$
$$= \frac{\eta}{K} \sum_{k=1}^K \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\eta^2}{K^2} \sum_{k=1}^K \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0). \tag{22}$$

Multiplying by the mixing coefficient  $\alpha$ :

$$\alpha(\phi_K - \phi_0) \approx \frac{\alpha \eta}{K} \sum_{k=1}^K \mathbf{g}_{\sigma(k)}(\phi_0) - \frac{\alpha \eta^2}{K^2} \sum_{k=1}^K \sum_{l=1}^{k-1} \mathbf{H}_{\sigma(k)}(\phi_0) \mathbf{g}_{\sigma(l)}(\phi_0).$$
 (23)

Note that  $\sum_{k=1}^{K} \mathbf{g}_{\sigma(k)}(\phi_0) = \sum_{k=1}^{K} \mathbf{g}_k(\phi_0)$  (permutation invariant), yielding Eq. (5).

#### B.4 DERIVATION OF GRADIENT ALIGNMENT SYMMETRIZATION

In Section 3.3, it is mentioned that after randomization, the cross-term expectation symmetrizes as:

$$\mathbf{H}_{i}(\phi_{0})\mathbf{g}_{j}(\phi_{0}) + \mathbf{H}_{j}(\phi_{0})\mathbf{g}_{i}(\phi_{0}) = \frac{\partial}{\partial\phi_{0}} \left(\mathbf{g}_{i}(\phi_{0})^{\mathsf{T}}\mathbf{g}_{j}(\phi_{0})\right). \tag{24}$$

**Derivation:** The key mathematical insight is the following identity concerning the gradient of the inner product between two gradients.

Consider the inner product  $S(\phi_0) = \mathbf{g}_i(\phi_0)^{\top} \mathbf{g}_j(\phi_0)$ . The gradient of this scalar function S with respect to  $\phi_0$  is given by:

$$\nabla_{\phi_0} S = \nabla_{\phi_0} \left( \mathbf{g}_i(\phi_0)^\top \mathbf{g}_j(\phi_0) \right) = \mathbf{H}_i(\phi_0) \mathbf{g}_j(\phi_0) + \mathbf{H}_j(\phi_0) \mathbf{g}_i(\phi_0), \tag{25}$$

where we have used the product rule and the symmetry of the Hessian matrices,  $\mathbf{H}_j = \mathbf{H}_j^{\top}$ . This result can be seen by noting that the derivative of  $\mathbf{g}_i^{\top}\mathbf{g}_j$  w.r.t.  $\phi_0$  is  $(\partial \mathbf{g}_i/\partial \phi_0)^{\top}\mathbf{g}_j+\mathbf{g}_i^{\top}(\partial \mathbf{g}_j/\partial \phi_0)=\mathbf{H}_i\mathbf{g}_j+\mathbf{g}_i^{\top}\mathbf{H}_j$ . Since  $\mathbf{g}_i^{\top}\mathbf{H}_j$  is a row vector, its transpose is  $\mathbf{H}_j\mathbf{g}_i$ . The gradient (as a column vector) is therefore  $\mathbf{H}_i\mathbf{g}_j+\mathbf{H}_j\mathbf{g}_i$ .

Under a randomized ordering  $\sigma$ , the expectation of the cross-term involving  $\mathbf{H}_{\sigma(k)}\mathbf{g}_{\sigma(l)}$  for k>l will involve pairs (i,j) symmetrically. The update term derived from the second-order expansion is proportional to  $\mathbf{H}_i\mathbf{g}_j$ . The symmetric form  $\mathbf{H}_i\mathbf{g}_j + \mathbf{H}_j\mathbf{g}_i$  appearing in the gradient of the inner product  $\nabla_{\phi_0}(\mathbf{g}_i^{\mathsf{T}}\mathbf{g}_j)$  indicates that, in expectation, the update encourages an increase in the inner product between the gradients of different domains, thus promoting their alignment.

#### C MORE DETAILS OF EXPERIMENTS

#### C.1 TASKS AND DATASETS

We focus on enhancing LLMs' overall capabilities across four domains—mathematical reasoning, code generation, scientific QA, and creative writing. These domains not only represent **core areas of current research interest** but also **span four distinct forms of reward feedback**, thereby ensuring both **comprehensiveness** and **diversity**.

- **Mathematics**: we construct a subset of 6,250 samples from the Guru dataset (Cheng et al., 2025). This includes the 5,000 easiest problems (ranked by the pass rate of Qwen2.5-7B-Instruct) and 1,250 more challenging ones, ensuring a balance between accessible and difficult problems.
- Code generation: we select a total of 4,740 samples from Guru. Specifically, we take all 3,791 problems with a Qwen2.5-7B-Instruct's pass rate of at least 25% and add 949 problems randomly sampled from the remainder, yielding an approximate 4:1 ratio between easier and harder samples.
- **Scientific QA**: we include the entire STEM split of Guru, resulting in 3,591 samples. This preserves the full coverage of science-related reasoning tasks while maintaining consistency with prior benchmarks.
- Creative writing: we randomly sample 2,000 samples each from three popular Hugging-face datasets—LitBench (Fein et al., 2025), Creative\_Writing-ShareGPT (Nitral-AI, 2024), and wildchat-creative-writing-3k-rft (kevinshin, 2025)—to construct a dataset of 6,000 samples, ensuring stylistic variety and broad coverage of open-ended writing abilities.

# C.2 BASELINES

We compare our CGPO against four representative baselines: joint learning, Omni-Thinker (Li et al., 2025a), Self-Paced CL, and FAMO (Liu et al., 2023).

• **Joint learning.** Joint learning is the most basic paradigm in MTL. It aggregates the loss functions of all tasks into a single objective, enabling simultaneous optimization. As a straightforward training strategy without any task-specific adjustments, joint learning serves as a reference point for evaluating improvements brought by more advanced methods.

- Omni-Thinker. Omni-Thinker belongs to *progressive CL* methods as categorized in (Soviany et al., 2022). It introduces the backward transfer (BWT) metric to quantify the extent of catastrophic forgetting across domains. Based on BWT analysis, Li et al. (2025a) proposes a fixed training order—*code* → *math* → *scientific QA* → *creative writing*—with the goal of minimizing forgetting induced by multi-domain learning.
- Self-paced CL. Self-paced CL enables the model to adaptively select training samples according to its learning state. In our implementation, we employ Qwen2.5-7B-Instruct to rank samples by winrate from easy to difficult, and train sequentially following this order. This curriculum reduces the risk of being misled by difficult samples in the early stages, thereby improving stability and promoting better generalization.
- FAMO. FAMO is a gradient-balancing approach for MTL. It adjusts loss weights to maximize the improvement rate of the task that progresses the slowest, ensuring that all tasks advance at a comparable pace. This balanced optimization strategy suppresses task dominance and guides the model toward solutions that are both fairer across tasks and stronger in overall performance. FAMO approximates weight updates using historical loss values instead of explicitly computing multi-task gradients, reducing per-iteration time and memory complexity to  $\mathcal{O}(1)$ . This efficiency makes it particularly suitable for large-scale LLM training.

#### C.3 REWARD FUNCTIONS

For all domains, we require the model to enclose its reasoning process within <think></think> tags. The reward functions for the four domains are as follows.

• Math. We adopt a rule-based reward function:

$$r_{\mathrm{math}}(o,a) = \begin{cases} 1.0, & \text{if } o \text{ has a valid format and verify}_{\mathrm{math}}(o_{\mathrm{ans}},a) = \mathrm{true}, \\ -0.5, & \text{if } o \text{ has a valid format but verify}_{\mathrm{math}}(o_{\mathrm{ans}},a) = \mathrm{false}, \\ -1.0, & \text{if } o \text{ has an invalid format}, \end{cases}$$

where  $o_{\rm ans}$  denotes the predicted answer extracted from structured tags (e.g., <answer></answer>) in the model output o, and verify  $_{\rm math}(\cdot,\cdot)$  checks symbolic equivalence between  $o_{\rm ans}$  and the ground-truth answer a via a deterministic parser (e.g., handling equivalent forms of expressions or equations).

• Code generation. We adopt a sandbox-based unit test reward:

$$r_{\text{code}}(o, \text{test\_case}) = \begin{cases} 1.0, & \text{if } o \text{ has a valid format and } \text{exec}(o_{\text{ans}}) \models \text{unittest}(o_{\text{ans}}, \text{test\_case}), \\ -0.5, & \text{if } o \text{ has a valid format but } \text{exec}(o_{\text{ans}}) \not\models \text{unittest}(o_{\text{ans}}, \text{test\_case}), \\ -1.0, & \text{if } o \text{ has an invalid format (syntactically invalid),} \end{cases}$$

where  $o_{ans}$  is the generated code, executed in a sandbox and validated against the unit tests associated with the sample;  $\models$  denotes logical satisfaction.

• Scientific QA. We employ a 1.5B General-Verifier<sup>1</sup> (Cheng et al., 2025) to assess consistency between the model's output and the ground-truth answer:

$$r_{\rm qa}(o,a) = \begin{cases} 1.0 - 0.05 \cdot \min(||o_{\rm ans}| - |a|| \,, 10) \,, & \text{if $o$ has a valid format and $o_{\rm ans} = a$,} \\ 0, & \text{if $o$ has a valid format but $o_{\rm ans} \neq a$,} \\ -1.0, & \text{if $o$ has an invalid format,} \end{cases}$$

where  $o_{ans}$  is the extracted answer content. Here, "valid format" means the response adheres to QA conventions (e.g., no garbled text, complete sentences).

• Creative writing. We adopt an LLM-as-a-Judge strategy, scoring the model's output o against a reference  $o_{\text{ref}}$  via pairwise comparison:

$$r_{\mathrm{writing}}(o,o_{\mathrm{ref}}) = \begin{cases} 1.0, & \text{if } o \text{ has a valid format and } o \succ o_{\mathrm{ref}}, \\ 0.25, & \text{if } o \text{ has a valid format and } o \sim o_{\mathrm{ref}}, \\ -0.5, & \text{if } o \text{ has a valid format and } o \prec o_{\mathrm{ref}}, \\ -1.0, & \text{if } o \text{ has an invalid format}, \end{cases}$$

where  $o \succ o_{\text{ref}}$  (preferred),  $o \sim o_{\text{ref}}$  (tie), and  $o \prec o_{\text{ref}}$  (worse) are determined by a fixed evaluator (Qwen2.5-72B-Instruct) serving as the judge.

<sup>1</sup>https://huggingface.co/TIGER-Lab/general-verifier

# C.4 HYPERPARAMTERS

We use a learning rate of  $1 \times 10^{-6}$ , a prompt batch size of 128, a mini-batch size of 64, a group size of 8, a rollout temperature of 1.0,  $\varepsilon_{\text{low}} = 0.2$ ,  $\varepsilon_{\text{high}} = 0.28$ , and  $\beta = 0.001$  for CGPO and all baselines. All methods are trained for one epoch. For the mixing coefficient  $\alpha$ , we tune it within the range of 0.5–1.5, and provide an ablation study on  $\alpha$  in Section 4.3.

#### C.5 EVALUATION

To comprehensively evaluate cross-domain capabilities, we adopt authoritative benchmarks spanning four domains: **Math**, **Coding**, **Scientific QA**, and **Creative Writing**. The evaluation settings are detailed below:

#### · Math domain

- MATH500 (Hendrycks et al., 2021): A set of 500 challenging problems sampled from the full MATH dataset, covering seven areas: elementary algebra, algebra, geometry, number theory, combinatorics, probability, and calculus. Problems are presented in open-ended form and require precise solutions. This benchmark is widely adopted for assessing LLMs' mathematical reasoning and problem-solving abilities.
- AMC 2023 (MAA, 2023): A set of 50 questions taken from the AMC 12A and 12B (2023) contests, spanning algebra, geometry, number theory, combinatorics, and probability. Multiple-choice options are removed, requiring models to directly output the final answer. This benchmark focuses on higher-order reasoning, problem analysis, and accurate calculation.

#### Coding domain

- HumanEval (Chen et al., 2021): Consisting of 164 human-written Python programming tasks, ranging from basic algorithms to medium-level function implementations. It evaluates whether models can generate correct and executable code from natural language descriptions.
- MBPP (Austin et al., 2021): A collection of 974 beginner-level Python problems designed to test the ability to synthesize short programs from natural language instructions. It is a standard benchmark for fundamental code generation.

#### · Scientific QA domain

- GPQA (diamond split) (Rein et al., 2023): Graduate-level QA items written and verified by domain experts across physics, chemistry, biology, and earth sciences. The diamond split represents the most difficult and highest-quality subset, specifically constructed to prevent shallow memorization or pattern matching. To ensure consistent evaluation, we reconstruct ordered option lists using randomized indexing.
- SuperGPQA (Team et al., 2025): Comprising 285 interdisciplinary graduate-level reasoning problems, curated to prevent direct solutions via search engines. To reduce computational cost, we use random seed 42 to sample 200 problems, ensuring both representativeness and reliable measurement of deep reasoning ability.

# · Creative Writing domain

- WritingBench (Wu et al., 2025): A benchmark of 1000 real-world writing tasks spanning 6 domains and 100 sub-themes, covering diverse styles, task types, and difficulty levels. It evaluates generated text on quality, coherence, creativity, and task alignment through a structured scoring framework. For efficiency, we sample 200 requests using random seed 42, and apply the official critic model WritingBench-Critic-Model-Qwen-7B<sup>2</sup> for automated scoring, striking a balance between evaluation cost and representativeness.

# D RELATED WORK

**Multi-domain RL for LLMs.** The application of RL in LLMs receives widespread attention (Schulman et al., 2017; Shao et al., 2024; Yu et al., 2025a; Liu et al., 2025). However, RL strategies that simultaneously and steadily enhance the capabilities of LLMs across multiple domains remain an

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/AQuarterMile/WritingBench-Critic-Model-Qwen-7B

 open challenge. A key difficulty in this area lies in designing reward functions that work effectively across diverse domains. Some researchers develop reward computation methods that are broadly applicable across multiple domains. For example, Zhou et al. (2025) simplify the binary reward function by leveraging properties of the ground truth. RLPR (Yu et al., 2025b) constructs its reward based on the probability of generating correct outputs. Other researchers create distinct reward computation methods tailored to specific domains. For instance, Li et al. (2025a) propose a hybrid reward system that employs rule-based, sandbox-based, and LLM-as-a-Judge frameworks, customized for different types of data. Another challenge lies in appropriately handling interactions among multiple domains. Cheng et al. (2025) study the effects of single-domain training on other domains. Li et al. (2025b) further examine interactions across several domains, including math, coding, and puzzle solving. Existing approaches mainly rely on experimental and qualitative observations, while a deeper understanding of cross-domain interactions remains largely unexplored.

Mitigating Gradient Conflicts. Gradient conflicts pose a major challenge in machine learning, leading to slow learning and wasted computation (Chen et al., 2025). Much work in multi-task learning addresses this by balancing or projecting gradients to reduce interference, such as GradNorm (Chen et al., 2018), which adjusts each task's gradient according to its relative loss, PCGrad (Yu et al., 2020), which projects away conflicting directions, MGDA (Sener & Koltun, 2018), which seeks Pareto-optimal updates, and ConFIG (Liu et al., 2024) or CAGrad (Liu et al., 2021), which optimize updates under constraints to ensure conflict-free directions. While effective in standard MTL, these approaches face key limitations in RL for LLMs: they generally either require storing all domain gradients on the GPU, which quickly becomes memory-intensive and can often cause out-of-memory failures, or act reactively without leveraging the underlying geometry of the reward landscape, which usually makes them prone to high variance on noisy, rollout-based gradients. These challenges motivate scalable, memory-efficient methods that can mitigate cross-domain conflicts while supporting multi-domain RL training, such as our proposed CGPO.

**Second-Order Optimization Methods.** The loss landscapes of deep neural networks are often highly complex, posing significant challenges for first-order optimization algorithms, such as gradient descent, which rely solely on local gradient information. Without insights into the geometric structure of the landscape, first-order methods can easily get trapped in saddle points or narrow valleys, making it difficult to reach better local optima. In contrast, second-order optimization methods, such as Newton's method, exploit geometric information like the Hessian matrix to precondition gradients according to the local curvature, offering stronger theoretical guarantees. To mitigate the computational cost of full Hessian computation, various approximate Newton methods have been proposed, including AdaGrad, K-FAC, GGT, Shampoo, and SOAP (Duchi et al., 2011; Martens & Grosse, 2015; Agarwal et al., 2018; Gupta et al., 2018; Vyas et al., 2025). Recent studies show that Newton's method and its approximate variant SOAP (Vyas et al., 2025) can alleviate gradient conflicts in physics-informed neural networks (PINNs) (Wang et al., 2025), providing inspiration for our approach. However, due to the massive parameter scale of large language models, directly applying Newton-type methods or their approximations in RL for LLMs is infeasible. Motivated by this, we distill the core idea of leveraging curvature information and develop CGPO, a principled and scalable framework for multi-domain RL in LLMs.