BALANCING THE EXPERTS: UNLOCKING LORA-MOE FOR GRPO VIA MECHANISM-AWARE REWARDS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026027028

029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Parameter-efficient Mixture-of-Experts (MoE) architectures, such as LoRA-MoE, enable strong and generalizable fine-tuning. However, a critical problem arises when fine-tuning these architectures with advanced reinforcement learning algorithms such as Group Relative Policy Optimization (GRPO). Traditional supervised techniques, such as auxiliary losses, are incompatible with the GRPO process, while the external task reward is blind to the internal routing mechanism. This disconnect leads to routing collapse and severe underutilization of MoE adapter parameters. To resolve this disconnect, we introduce Routing-Optimized Group Relative Policy Optimization (RO-GRPO), a mechanism-aware framework. It turns internal expert routing statistics collected during training into a direct reward signal, seamlessly integrating routing supervision into the reinforcement fine-tuning (RFT) process. This enables effective optimization of parameter utilization and improves performance on both unimodal and multimodal mathematical reasoning tasks, all without extra training stages. Our work provides the first demonstration that a scalar reward in GRPO can be engineered from a model's own internal mechanics to explicitly guide their optimization, extending alignment from mere behavior tuning to holistic mechanism alignment.

1 Introduction

Large language models (LLMs) have significantly advanced many artificial intelligence applications, but their large size poses challenges for practical deployment, especially regarding fine-tuning efficiency (Han et al., 2024; Hu et al., 2021). Among parameter-efficient fine-tuning (PEFT) approaches, LoRA-MoE (Dou et al., 2024), which combines Low-Rank Adaptation (LoRA) (Hu et al., 2021) with a mixture-of-experts (MoE) architecture, has emerged as a particularly promising technique (Li et al., 2024a; Gou et al., 2024; Mu & Lin, 2025).

However, despite its success in supervised fine-tuning (SFT) (Dou et al., 2024; Li et al., 2024a), applying LoRA-MoE to reinforcement learning fine-tuning (RFT) frameworks such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024) presents a key challenge. In SFT, routing is typically guided by an auxiliary load-balancing loss (Fedus et al., 2022; Lewis et al., 2021; Zoph et al., 2022; Dai et al., 2022; Wang et al., 2024). This approach is incompatible with most RFT algorithms, which rely on a single scalar reward. Consequently, the external task reward used in GRPO is blind to internal routing decisions (Omi et al., 2025; Harvey et al., 2025). Without explicit guidance, the routing mechanism often collapses, leading to severe expert imbalance and underutilization of the model's parametric capacity, which limits the effectiveness of the modular architecture.

To bridge this gap, we propose **RO-GRPO** (Routing-Optimized GRPO), a novel framework that integrates routing-awareness into the RFT process through a carefully designed mechanism reward. Our key insight is that routing statistics collected during generation, such as routing entropy and load distribution, can be transformed into a reward signal that aligns the internal routing mechanism with task performance (Chen et al., 2022; Cong et al., 2024b). Specifically, we augment the standard task reward with two complementary components: one promoting confident routing decisions (low entropy) and another encouraging balanced expert utilization. These routing rewards are integrated into the GRPO objective, enabling a unified optimization process that requires no auxiliary losses or architectural modifications (see Figure 1 for a schematic comparison).

056

060

061 062

063

064

065

066

067

069

071

073074075

076 077

078

079

081

082

084

085

087

090

092

094

095

096

097

098

100

101

102

103

104

105

106

107

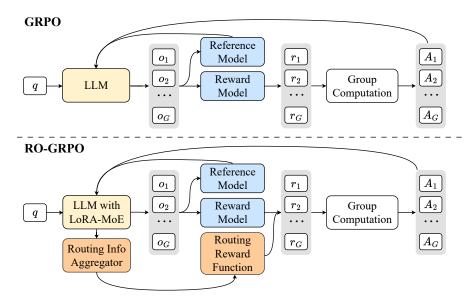


Figure 1: Comparison of standard GRPO and our RO-GRPO.

The main contributions of this paper are as follows:

- To our knowledge, this is the first systematic study of LoRA-MoE architectures within the RFT framework, identifying and addressing key challenges such as routing collapse and expert underutilization that arise during GRPO-based training.
- We propose RO-GRPO, a novel framework that integrates routing statistics directly into the GRPO reward function. This enables the unified optimization of both task performance and internal routing efficiency without requiring auxiliary losses.
- Our method achieves consistent improvements over baselines (e.g., standard LoRA and vanilla LoRA-MoE) across all expert counts in both task performance and expert utilization, validated across unimodal and multimodal mathematical reasoning benchmarks.
- Our experiments provide the first empirical evidence that a scalar reward in RFT can align not only a model's external behavior but also its internal mechanisms, opening new avenues for the principled alignment of complex model architectures.

2 RELATED WORK

Modular and Mixture-of-Experts PEFT. To enhance the capacity and versatility of PEFT, researchers have integrated Mixture-of-Experts (MoE) principles into LoRA, creating architectures like LoRA-MoE (Dou et al., 2024), MixLoRA (Li et al., 2024a), and MoCLE (Gou et al., 2024). These methods have shown effectiveness in reducing task interference and improving knowledge retention. Beyond adapterized MoE, modern routing builds on early conditional computation and deep MoE ideas (Cho & Bengio, 2014; Eigen et al., 2013) and on system-scale routed Transformers such as GShard and Switch (Lepikhin et al., 2020; Fedus et al., 2022). Most approaches optimize expert routing during supervised pretraining or SFT, typically using auxiliary load-balancing objectives to prevent expert collapse and under-utilization (Shazeer et al., 2017; Fedus et al., 2022). Alternative balancing mechanisms include the balanced assignment of BASE Layers (Lewis et al., 2021) and non-parametric routing via Hash Layers (Roller et al., 2021). Training stability and routing fluctuation have been studied extensively, with design guidelines in ST-MoE (Zoph et al., 2022), two-stage stabilized routing in StableMoE (Dai et al., 2022), and empirical analysis of expert-load dynamics (Cong et al., 2024a). However, these strategies remain predominantly designed for differentiable supervised training; integrating comparable mechanism-aware supervision into reinforcement finetuning is still underexplored.

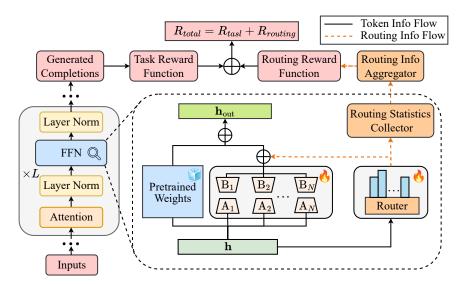


Figure 2: Overview of the RO-GRPO framework. A mechanism-aware reward R_{route} is computed from internal routing statistics and combined with the task reward R_{task} . The resulting unified reward R_{total} guides the GRPO update to jointly optimize task performance and routing efficiency.

Alignment and Optimization of LLMs. Reinforcement Learning from Human Feedback (RLHF) has become the dominant approach for aligning large language models with human intentions (Ouyang et al., 2022; Kaufmann et al., 2024). Proximal Policy Optimization (PPO) (Schulman et al., 2017) and its variants, such as Direct Preference Optimization (DPO) (Rafailov et al., 2024) and Group Relative Policy Optimization (GRPO) (Shao et al., 2024), have demonstrated strong performance on complex reasoning and instruction following. These RLHF methods typically optimize a scalar task-based reward and do not incorporate supervision for internal mechanisms such as expert routing. In MoE-based models, this limitation can lead to expert collapse or inefficient parameter utilization when using RLHF directly (Fedus et al., 2022; Harvey et al., 2025). While auxiliary losses have been used to encourage balanced routing in supervised settings, integrating such mechanism-aware supervision into reinforcement learning remains an open problem.

3 METHODOLOGY

In this section, we introduce **RO-GRPO** (Routing-Optimized Group Relative Policy Optimization), a framework designed to optimize the internal routing of LoRA-MoE models by incorporating mechanism-aware supervision into the reinforcement learning loop. We first review the preliminaries of GRPO and LoRA-MoE, then describe the core challenge of unguided routing in RLHF, and finally detail our proposed solution.

3.1 PRELIMINARIES: GRPO AND LORA-MOE

Group Relative Policy Optimization. GRPO (Shao et al., 2024) is a critic-free RL algorithm that aligns an LLM policy π_{θ} by maximizing the expected task reward over a dataset of prompts \mathcal{D} . Its objective is to maximize $\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)}[R_{\text{task}}(y)]$, where π_{θ} is the LLM policy, \mathcal{D} is the dataset of prompts, and R_{task} is the scalar task reward evaluating the output y.

LoRA-MoE Architecture. A LoRA-MoE layer modifies the output of a frozen pretrained layer. It consists of a trainable router network and a set of E parallel LoRA experts, $\{(\mathbf{A}_e, \mathbf{B}_e)\}_{e=1}^{E}$. For an input token representation \mathbf{h} , the router first computes a gating probability vector $\mathbf{p} = \operatorname{softmax}(\mathbf{W}_r\mathbf{h})$, where $\mathbf{W}_{\text{router}}$ is the router's weight matrix. The final output of the layer is:

$$\mathbf{h}_{\text{out}} = \mathbf{h} + \left(\sum_{e=1}^{E} p(e \mid \mathbf{h}) \mathbf{B}_{e} \mathbf{A}_{e}\right) \mathbf{h}.$$
 (1)

3.2 THE CHALLENGE: UNGUIDED MOE ROUTING

A fundamental disconnect arises when a LoRA-MoE model is fine-tuned using an RLHF algorithm like GRPO. The optimization process is blind to the router's decisions \mathbf{p} , as the task reward R_{task} evaluates only the final output. This lack of explicit supervision leads to two well-documented failure modes in MoE training (Fedus et al., 2022):

- Expert Collapse: The router defaults to choosing a small subset of experts, leading to severe load imbalance and wasted parametric capacity.
- Routing Indecision: The router generates high-entropy distributions (i.e., low-confidence decisions), failing to foster expert specialization.

Our goal is to augment the RLHF objective with an internal, mechanism-aware reward signal that directly addresses these failure modes.

3.3 RO-GRPO: MECHANISM-AWARE REWARDS

As depicted in Figure 2, our method collects internal routing statistics during policy generation. These statistics are used to compute $R_{\rm route}$, which is then combined with the primary task reward, $R_{\rm task}$. This is achieved in three steps.

3.3.1 Step 1: Quantifying Routing Efficiency.

For each generated sample, we collect the routing probability vectors from all activated LoRA-MoE modules. Let M be the number of such modules in the model and T be the total number of tokens routed during generation. We quantify routing efficiency using two metrics aggregated from these statistics. First, we measure **routing confidence** using the mean Shannon entropy over all individual token routing decisions. A lower value indicates more decisive routing. For the multiset of all T routing vectors $\{\mathbf{p}_i\}_{i=1}^T$, define the token-wise entropy as $H(\mathbf{p}) := -\sum_{e=1}^E p_e \ln p_e$ (a small ϵ is added in implementation). The average entropy is

$$\bar{\mathcal{H}} = \frac{1}{T} \sum_{i=1}^{T} H(\mathbf{p}_i). \tag{2}$$

Second, we measure **load balance** by assessing expert utilization across the M LoRA-MoE modules. For each module $m \in \{1,\ldots,M\}$, we first compute its average expert utilization vector $\bar{\mathbf{p}}_m$ by averaging the routing vectors of all tokens that pass through it. The final metric $\bar{\mathcal{M}}$ is the mean of the Mean Squared Errors (MSE) calculated for each module relative to a uniform distribution \mathbf{u} :

$$\bar{\mathcal{M}} = \frac{1}{M} \sum_{m=1}^{M} \frac{1}{E} \left\| \bar{\mathbf{p}}_m - \frac{1}{E} \mathbf{1} \right\|_2^2. \tag{3}$$

For stable integration into the reward function, we normalize these metrics to an approximate [0,1] range, yielding $\mathcal{H}_{norm} = \bar{\mathcal{H}}/\ln E$, $\mathcal{M}_{norm} = \bar{\mathcal{M}}/((E-1)/E^2)$, for use in the reward calculation.

3.3.2 Step 2: Formulating the Routing Reward.

We propose two distinct strategies to transform these metrics into a scalar reward R_{route} .

RO-GRPO (Smooth): Curriculum-Based Reward Scheduling. This strategy employs a curriculum that initially encourages confident routing (low entropy) and then transitions to promoting load balance (low MSE). As detailed in the Discussion section, this curriculum aligns the reward signal with the natural training dynamics of MoE models, defining the reward as follows:

$$R_{\text{route}} = -w_{\text{route}} \left(w_H(t) \cdot \mathcal{H}_{\text{norm}} + w_B(t) \cdot \mathcal{M}_{\text{norm}} \right), \tag{4}$$

where weights $w_H(t)$ and $w_B(t)$ are dynamically scheduled based on training progress $t \in [0,1]$ using a sigmoid function $\sigma(t) = \frac{1}{1 + e^{-k(t-c)}}$ with steepness k and center c:

$$w_H(t) = \lambda_H^{\text{start}} \cdot (1 - \sigma(t)), \tag{5}$$

$$w_B(t) = \lambda_B^{\text{end}} \cdot \sigma(t). \tag{6}$$

Table 1: Performance on unimodal mathematical reasoning benchmarks. We compare task accuracy (%), trainable parameter count (#Param), and internal routing metrics.

Unimodal Mathematical Reasoning (Qwen2.5-7B-Instruct on NuminaMath-TIR-2k)								
Method	#Experts	#Param	GSM8K	MATH	SVAMP	MGSM	Entropy	MSE
Base (zero-shot)	-	0	87.34	70.42	91.33	53.64	-	-
GRPO (LoRA)	-	30.3M	88.48	70.38	90.67	50.00	-	-
GRPO (LoRA-MoE)	2	31.7M	89.39	70.36	90.00	61.75	0.640	0.020
	4	63.4M	89.39	70.40	91.30	46.15	0.651	0.009
	8	126.9M	90.22	70.44	91.00	52.04	0.655	0.008
	2	31.7M	91.51	70.64	91.00	62.18	0.639	0.016
RO-GRPO (Smooth)	4	63.4M	90.67	70.62	92.00	52.58	0.651	0.009
	8	126.9M	90.98	69.78	92.67	52.04	0.656	0.006
RO-GRPO (Relative)	2	31.7M	90.22	70.58	91.33	59.45	0.639	0.017
	4	63.4M	89.76	69.88	<u>93.33</u>	54.58	0.651	0.008
	8	126.9M	90.52	70.18	92.67	51.96	0.655	0.007

RO-GRPO (**Relative**): **Relative Improvement Gating.** This strategy provides a sparse, adaptive reward based on a historical baseline, encouraging continuous self-improvement and avoids the need to manually balance the two routing objectives. A constant positive reward C is granted only if both routing confidence and load balance improve simultaneously relative to their exponential moving averages ($\bar{\mathcal{H}}_{hist}$, \mathcal{M}_{hist}):

$$R_{\text{route}} = C \mathbf{1} \{ \mathcal{H}_{\text{norm}} < \bar{\mathcal{H}}_{\text{hist}} \land \mathcal{M}_{\text{norm}} < \bar{\mathcal{M}}_{\text{hist}} \}. \tag{7}$$

3.3.3 Step 3: Unified Optimization via Policy Gradient.

The mechanism-aware reward R_{route} is combined with the external task reward to form the total reward $R_{\text{total}}(y) = R_{\text{task}}(y) + R_{\text{route}}(y)$. The computation of R_{route} is non-differentiable; its gradient is propagated implicitly through the policy gradient update. In the GRPO framework, a group of responses $\{y_i\}$ is sampled for each prompt and evaluated using R_{total} . The resulting rewards are used to compute group-relative advantages, \hat{A}_i , which in turn guide the policy update. The objective can be summarized as:

$$\mathcal{J}_{\text{RO-GRPO}}(\theta) \approx \mathbb{E}\left[\sum_{i} \log \pi_{\theta}(y_{i}|x) \,\hat{A}_{i} - \beta \, D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}})\right].$$
 (8)

By integrating R_{route} into the advantage calculation, RO-GRPO guides the policy π_{θ} to generate outputs that improve task performance while also exhibiting efficient routing, all without requiring differentiable auxiliary losses.

4 EXPERIMENTS

We experimentally validate RO-GRPO by testing three core hypotheses: (1) applying LoRA-MoE with GRPO leads to suboptimal routing and underutilized parameters; (2) our mechanism-aware reward framework, RO-GRPO, mitigates these routing issues; and (3) these internal improvements translate to better task performance.

4.1 EXPERIMENTAL SETUP

Tasks and Datasets. We evaluate RO-GRPO on challenging mathematical reasoning tasks. Performance on these tasks relies on precise, multi-step deduction, making it sensitive to the model's expert utilization and thus an ideal testbed for our approach. To demonstrate the versatility of our method, we conduct experiments in both unimodal and multimodal settings.

For unimodal experiments, we fine-tune on NuminaMath-TIR (Li et al., 2024b) and evaluate on the established benchmarks GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), SVAMP (Pa-

Table 2: Performance on multimodal mathematical reasoning benchmarks. We follow the same evaluation setup as in the unimodal experiments.

Multimodal Mathematical Reasoning (Qwen2.5-VL-7B-Instruct on Geometry3k)								
Method	#Experts	#Param	Geo3k	MathVista	MathVerse	WeMath	Entropy	MSE
Base (zero-shot)	-	0	37.44	46.50	26.50	56.95	-	-
GRPO (LoRA)	-	30.3M	38.44	58.60	33.30	63.97	-	-
GRPO (LoRA-MoE)	2	31.7M	38.27	57.90	30.99	63.74	0.619	0.038
	4	63.4M	28.95	56.40	30.30	63.45	0.637	0.019
	8	126.9M	33.11	55.00	31.78	61.49	0.649	0.012
	2	31.7M	38.94	58.70	30.48	66.09	0.630	0.033
RO-GRPO (Smooth)	4	63.4M	40.10	58.30	28.73	64.14	0.642	0.014
	8	126.9M	38.94	58.90	27.89	64.10	0.648	0.011
	2	31.7M	41.93	55.80	33.30	60.98	0.624	0.036
RO-GRPO (Relative)	4	63.4M	40.27	60.20	31.29	<u>66.26</u>	0.645	0.013
	8	126.9M	40.16	60.10	32.03	63.97	0.636	0.010

tel et al., 2021), and MGSM (Shi et al., 2022). For multimodal experiments, we fine-tune on Geometry3k (Lu et al., 2021) and evaluate on its test set, alongside MathVista (Lu et al., 2024), Math-Verse (Zhang et al., 2024), and WeMath (Qiao et al., 2024).

Models and Baselines. We use the open-source Qwen2.5-7B-Instruct (Qwen et al., 2025) and Qwen2.5-VL-7B-Instruct (Bai et al., 2025) models, chosen for their strong foundational performance in mathematical reasoning. We compare five configurations: (1) Base, the original pretrained model evaluated in a zero-shot setting; (2) GRPO (LoRA), a standard LoRA baseline fine-tuned with GRPO representing a typical PEFT approach; (3) GRPO (LoRA-MoE), a LoRA-MoE model trained with GRPO using only the task reward to isolate the effect of unguided routing; (4) RO-GRPO (Smooth), our method with the curriculum-based reward scheduling strategy (Section 3.3); and (5) RO-GRPO (Relative), our method with the relative improvement gating strategy for the routing reward (Section 3.3).

Evaluation Metrics. We evaluate both task performance and internal mechanism efficiency. Task Performance is measured by accuracy (%) on the respective benchmarks. Routing Performance is quantified by two metrics: (1) **routing entropy**, the average per-token Shannon entropy as formulated in Eq. (2), indicating decision confidence; and (2) **load balancing MSE**, the mean squared error between the expert utilization distribution and a uniform one as formulated in Eq. (3), indicating load balance. We report these raw, un-normalized values for direct interpretability.

Implementation Details. To ensure a fair comparison, we control for trainable-parameter budget: the LoRA baseline uses a rank of 16, while LoRA-MoE models use E experts ($E \in \{2,4,8\}$), each with rank of 8. For LoRA-MoE, modules are inserted into the gate, up, and down projections of the Feed-Forward Network (FFN) in each transformer block During training, only the PEFT parameters are updated while the base model weights remain frozen. Across all experiments, we use a consistent system prompt for both training and evaluation to encourage step-by-step reasoning. The overall routing reward $R_{\rm route}$ is integrated into the total reward using a global scaling coefficient of $w_{\rm route} = 0.2$. Key hyperparameters for our routing strategies were determined through validation. For the Smooth strategy, we set the initial entropy weight $\lambda_H^{\rm start} = 0.5$ and the final balance weight $\lambda_B^{\rm end} = 2.0$. For the Relative strategy, the performance baseline was calculated over a moving window of the 1000 most recent samples. Further details are in the Appendix.

4.2 MAIN RESULTS

As shown in Tables 1 and 2, RO-GRPO yields consistent performance gains over vanilla GRPO with LoRA-MoE across all expert counts ($E \in \{2,4,8\}$) in both unimodal and multimodal settings. The two reward variants show complementary strengths: the Smooth strategy performs best on GSM8K and SVAMP, while the Relative strategy excels on Geometry3k, MathVista, and WeMath.

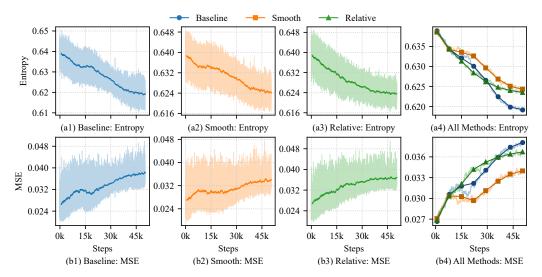


Figure 3: Training dynamics of routing metrics on the unimodal mathematical reasoning task. (Top) Average routing entropy over the course of training. (Bottom) Load balancing MSE over the course of training.

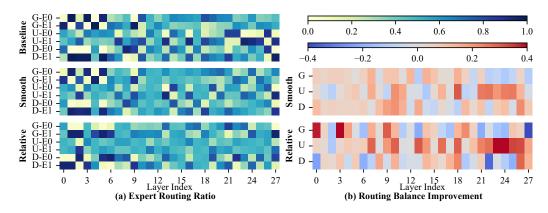


Figure 4: Visual analysis of routing behavior improvements with RO-GRPO on the MathVista benchmark. (a) Left Panel: Heatmaps show the routing ratio of the most frequently selected expert for the baseline and our two RO-GRPO methods. Darker colors represent a higher selection ratio for the dominant expert. (b) Right Panel: Heatmaps quantify the improvement in routing balance relative to the vanilla LoRA-MoE baseline. Positive values (warmer colors) indicate a reduced routing ratio for the dominant expert, signifying a shift toward the desired 1/E equilibrium.

Unimodal. On GSM8K, RO-GRPO (Smooth, E=2) achieves the top score of 91.51%, an improvement of +3.03 pp over GRPO (LoRA) and +1.29 pp over the best-performing vanilla LoRA-MoE model (at E=8). On SVAMP, RO-GRPO (Relative, E=4) obtains 93.33%, surpassing GRPO (LoRA) by +2.66 pp and vanilla LoRA-MoE (at E=4) by +2.03 pp. While improvements on MATH are modest, they are consistent at matched expert counts (e.g., 70.64% vs. 70.36% for E=2 Smooth vs. vanilla). The largest gain on MGSM is at E=4, where RO-GRPO (Relative) achieves 54.58%, an +8.43 pp increase over the vanilla model.

Multimodal. RO-GRPO (Relative, E=2) attains the highest Geometry3k score of 41.93%, outperforming vanilla LoRA-MoE (E=2) by +3.66 pp and GRPO (LoRA) by +3.49 pp. On MathVista and WeMath, the best results are with RO-GRPO (Relative, E=4), which reaches 60.20% and 66.26%. These scores represent gains of +1.60 pp and +2.29 pp over GRPO (LoRA), and +3.80 pp and +2.81 pp over the best vanilla LoRA-MoE model (at E=2), respectively. On MathVerse, the top performance of 33.30% is shared by GRPO (LoRA) and RO-GRPO (Relative, E=2).

Table 3: Ablation and causal analysis on GSM8K (E=2). Our full RO-GRPO model significantly outperforms the vanilla baseline. Subsequent experiments demonstrate that both reward components (R_H,R_B) are necessary for optimal performance, and control experiments validate that the gains are causally driven by our targeted reward signal.

Configuration	GSM8K	Entropy	MSE
RO-GRPO (Smooth)	91.51	0.639	0.016
RO-GRPO (Relative)	90.22	0.639	0.017
GRPO (LoRA-MoE)	89.39	0.640	0.020
Ablations on Reward Cor	nponents:		
w/o R_B (Smooth)	90.75	0.639	0.018
w/o R_H (Smooth)	89.92	0.639	0.019
w/o R_B (Relative)	89.01	0.638	0.019
w/o R_H (Relative)	90.14	0.637	0.019
Causal & Sanity Controls	s:		
Shuffled (Smooth)	89.23	0.640	0.018
Shuffled (Relative)	89.28	0.641	0.020

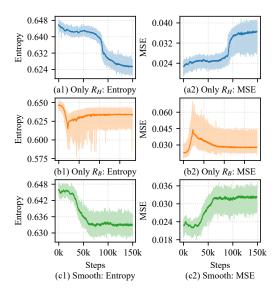


Figure 5: Training dynamics of routing metrics when the R_{task} is set to zero.

Overall, across expert sizes, at matched E our routing-aware training either matches or surpasses vanilla LoRA-MoE on nearly every benchmark, with the largest margins on Geometry3k, SVAMP, and WeMath. These results reaffirm that aligning the router with mechanism-aware rewards translates into stronger task performance.

4.3 Analysis of Routing Mechanism

Unguided routing (vanilla GRPO) is brittle. At $E{=}2$ in the multimodal setting, vanilla LoRA-MoE appears confident (Entropy 0.619) but exhibits routing collapse (MSE 0.038). At larger E, the raw MSE decreases (e.g., multimodal 0.019 and 0.012 for $E{=}4$ and $E{=}8$), yet accuracy does not reliably improve and can even drop (Geometry3k score of 28.95% at $E{=}4$), revealing instability in the absence of mechanism-aware feedback.

RO-GRPO restores balance at matched E. Across all experimental configurations, RO-GRPO reduces or matches the MSE of the vanilla baseline at the same E (unimodal: $0.020 \rightarrow 0.016/0.017$, $0.009 \rightarrow 0.009/0.008$, $0.008 \rightarrow 0.006/0.007$; multimodal: $0.038 \rightarrow 0.033/0.036$, $0.019 \rightarrow 0.014/0.013$, $0.012 \rightarrow 0.011/0.010$), while maintaining comparable entropy. These improvements in routing correspond to the largest accuracy gains on Geometry3k, WeMath, and SVAMP.

Routing rewards mitigate text degeneration. On Geometry3k with E=4, the vanilla GRPO (LoRA-MoE) model exhibits repetitive-loop failures in 7.5% of generations and scores 28.95%. In contrast, RO-GRPO (Smooth) reduces these failures to 0.17% and RO-GRPO (Relative) eliminates them entirely, achieving accuracies of 40.10% and 40.27%, respectively.

4.4 ABLATION AND CAUSAL VERIFICATION

Ablation experiments on GSM8K confirm our approach's integrity (Table 3).

Contribution of Reward Components. We first investigate the individual contributions of the confidence (R_H) and balancing (R_B) rewards. As shown in Table 3, removing either component from our best-performing model, RO-GRPO (Smooth), reduces accuracy. Specifically, removing the balancing reward $(w/\circ R_B)$ or the confidence reward $(w/\circ R_H)$ reduces the GSM8K score from 91.51% to 90.75% and 89.92%, respectively. This dependency is more pronounced for the RO-GRPO (Relative) variant: removing its balancing reward $(w/\circ R_B)$ causes performance to drop to

89.01%, below the vanilla baseline. These results demonstrate that the two reward components are synergistic and critical for optimal performance.

Contribution of the Reward Signal. To ensure performance gains are driven by meaningful feedback, we performed a control experiment. In the Shuffled Control, we randomly permuted routing rewards within each batch, breaking the causal link between an action and its reward. As shown in Table 3, performance under this condition dropped significantly for both the Smooth (89.23%) and Relative (89.28%) variants, falling to the level of the vanilla GRPO (LoRA-MoE) baseline (89.39%). This result strongly suggests the gains from RO-GRPO are causally driven by the targeted feedback from our reward signal, not by an artifact of the reward structure.

5 DISCUSSION

Our experiments demonstrate the empirical success of RO-GRPO across E=2,4,8. This section analyzes why a unified scalar reward can supervise a model's internal router and why this becomes more important as E grows. A more detailed derivation is available in Appendix E.

The Rationale for a Curriculum-Based Reward. The Smooth curriculum is effective because single-objective optimization is suboptimal: rewarding only low entropy degrades balance (MSE rises), whereas rewarding only balance is initially too weak to shape specialization (Figure 5). By first encouraging confident routing and then increasing pressure on balance, the curriculum builds specialized experts and subsequently organizes them. This dynamic mirrors our empirical trends at $E{=}4,8$, where mechanism-aware supervision not only improves accuracy but also suppresses degeneration on Geometry3k.

A Unified View via Constrained Optimization. RO-GRPO can be read as a penalty method for constrained policy optimization: maximize task reward subject to implicit constraints on routing entropy and load balance. The total reward $R_{\text{total}} = R_{\text{task}} + R_{\text{route}}$ functions as a simplified Lagrangian, allowing us to enforce mechanistic properties with standard policy gradients.

Grounding the Reward Components. The confidence reward, R_H , which promotes low-entropy routing, can be understood through the Information Bottleneck (IB) principle (Tishby et al., 2000). The IB principle states that an optimal representation should compress an input while preserving task-relevant information. In our framework, the router's decision acts as this bottleneck. By rewarding low-entropy (confident) decisions, RO-GRPO incentivizes the router to learn a minimal sufficient representation of its input. It is encouraged to discard noisy features and focus on information predictive of task success, a process that naturally fosters expert specialization.

The balancing reward R_B directly optimizes parameter utilization. Maximizing our balancing reward, which is formulated using MSE, is formally equivalent to minimizing the variance of the expert load distribution. This ensures the reward signal provides a direct and efficient gradient for combating routing collapse and ensuring the model leverages its full parametric capacity, a principle established in supervised MoE training (Shazeer et al., 2017).

6 CONCLUSION

We addressed a core limitation of applying LoRA-MoE to GRPO: the task reward is blind to routing. RO-GRPO remedies this by transforming routing statistics into a mechanism-aware reward that plugs into GRPO without architecture changes or extra stages. Across E=2, 4, 8 and both unimodal and multimodal math reasoning, RO-GRPO improves load balance at matched E, boosts accuracy, and reduces text degeneration. These results indicate that reinforcement learning can align not only external behavior but also internal mechanisms, suggesting a path toward principled alignment for complex modular architectures.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our findings, we provide a comprehensive set of resources. The complete source code for all experiments, including model implementation, training and evaluation scripts, is available in the supplementary material. Further details on the experimental environment, including the specific hardware and software configurations used, are documented in Appendix A. A complete list of hyperparameters for all model configurations and our routing reward strategies is provided in Appendix B, alongside pseudocode for our reward calculation algorithms.

REFERENCES

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL https://arxiv.org/abs/2502.13923.
- Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding mixture of experts in deep learning, 2022. URL https://arxiv.org/abs/2208.02813.
- Kyunghyun Cho and Yoshua Bengio. Exponentially increasing the capacity-to-computation ratio for conditional computation in deep learning. *arXiv* preprint arXiv:1406.7362, 2014.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
- Peizhuang Cong, Aomufei Yuan, Shimao Chen, Yuxuan Tian, Bowen Ye, and Tong Yang. Prediction is all moe needs: Expert load distribution goes from fluctuating to stabilizing. *arXiv preprint arXiv:2404.16914*, 2024a.
- Peizhuang Cong, Aomufei Yuan, Shimao Chen, Yuxuan Tian, Bowen Ye, and Tong Yang. Prediction is all moe needs: Expert load distribution goes from fluctuating to stabilizing, 2024b. URL https://arxiv.org/abs/2404.16914.
- Damai Dai, Li Dong, Shuming Ma, Bo Zheng, Zhifang Sui, Baobao Chang, and Furu Wei. Stablemoe: Stable routing strategy for mixture of experts. *arXiv preprint arXiv:2204.08396*, 2022.
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. Loramoe: Alleviate world knowledge forgetting in large language models via moe-style plugin, 2024. URL https://arxiv.org/abs/2312.09979.
- David Eigen, Marc'Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2022. URL https://arxiv.org/abs/2101.03961.
- Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James T. Kwok, and Yu Zhang. Mixture of cluster-conditional lora experts for vision-language instruction tuning, 2024. URL https://arxiv.org/abs/2312.12379.
- Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning for large models: A comprehensive survey, 2024. URL https://arxiv.org/abs/2403.14608.
- Daniel Fidel Harvey, George Weale, and Berk Yilmaz. Optimizing moe routers: Design, implementation, and evaluation in transformer models, 2025. URL https://arxiv.org/abs/2506.16419.

- 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, 2021. URL https://arxiv.org/abs/2103.03874.
 - Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL https://arxiv.org/abs/2106.09685.
 - Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback, 2024. URL https://arxiv.org/abs/2312.14925.
 - Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding, 2020. URL https://arxiv.org/abs/2006.16668.
 - Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers: Simplifying training of large, sparse models. In *International Conference on Machine Learning*, pp. 6265–6274. PMLR, 2021.
 - Dengchun Li, Yingzi Ma, Naizheng Wang, Zhengmao Ye, Zhiyuan Cheng, Yinghao Tang, Yan Zhang, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. Mixlora: Enhancing large language models fine-tuning with lora-based mixture of experts, 2024a. URL https://arxiv.org/abs/2404.15159.
 - Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*, 13(9):9, 2024b.
 - Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning, 2021. URL https://arxiv.org/abs/2105.04165.
 - Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024. URL https://arxiv.org/abs/2310.02255.
 - Siyuan Mu and Sen Lin. A comprehensive survey of mixture-of-experts: Algorithms, theory, and applications, 2025. URL https://arxiv.org/abs/2503.07137.
 - Nabil Omi, Siddhartha Sen, and Ali Farhadi. Load balancing mixture of experts with similarity preserving routers, 2025. URL https://arxiv.org/abs/2506.14038.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL https://arxiv.org/abs/2203.02155.
 - Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems?, 2021. URL https://arxiv.org/abs/2103.07191.
 - Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma GongQue, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, Runfeng Qiao, Yifan Zhang, Xiao Zong, Yida Xu, Muxi Diao, Zhimin Bao, Chen Li, and Honggang Zhang. We-math: Does your large multimodal model achieve human-like mathematical reasoning?, 2024. URL https://arxiv.org/abs/2407.01284.

- Qwen,:, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL https://arxiv.org/abs/2305.18290.
- Stephen Roller, Sainbayar Sukhbaatar, Jason Weston, et al. Hash layers for large sparse models. *advances in neural information processing systems*, 34:17555–17566, 2021.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners, 2022. URL https://arxiv.org/abs/2210.03057.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method, 2000. URL https://arxiv.org/abs/physics/0004057.
- Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free load balancing strategy for mixture-of-experts, 2024. URL https://arxiv.org/abs/2408.15664.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems?, 2024. URL https://arxiv.org/abs/2403.14624.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv* preprint *arXiv*:2202.08906, 2022.

A COMPUTING INFRASTRUCTURE AND SOFTWARE

All experiments were conducted on a high-performance computing cluster. The specific hardware and software configurations are provided to ensure full reproducibility.

Hardware: Each experiment was run on a single node equipped with 8x NVIDIA A800 (80GB VRAM) GPUs. Each node was powered by an Intel(R) Xeon(R) Platinum 8336C CPU with 1875 GB of system RAM.

Software: The operating system was Ubuntu 20.04.6 LTS. The core software stack included:

Python 3.10.18

- PyTorch 2.5.1 (built with CUDA 12.1)
- CUDA Toolkit 12.2
- Hugging Face Transformers 4.51.0
- Hugging Face PEFT 0.14.0
- The ms-swift framework, version 3.7.0.dev0, was used for all training scripts.

B HYPERPARAMETER AND IMPLEMENTATION DETAILS

Our approach to hyperparameter selection was designed to be both systematic and efficient. For established components of the training pipeline, such as the GRPO algorithm and the LoRA architecture, we adopted values from seminal works and common practices to establish strong, competitive baselines. Our primary tuning efforts were concentrated on the novel parameters introduced by the RO-GRPO framework, ensuring a rigorous evaluation of our core contributions.

Core Training and Architecture Parameters.

For all experiments, we used a learning rate of 1×10^{-5} and a batch size of 64. The GRPO configuration included a KL coefficient (β) of 0.1 and sampling 8 responses per prompt (k=8) for advantage estimation. For the base LoRA architecture, we set the rank to r=16 and alpha to $\alpha=32$. For our LoRA-MoE models, we used $E\in\{2,4,8\}$ experts, each with a rank of r=8 and an alpha of $\alpha=32$, maintaining a similar parameter budget. The training duration was set to 1 epoch for the unimodal Numina-Math dataset and 3 epochs for the more complex multimodal Geometry3k dataset to ensure convergence. The external task reward weight was consistently set to 1.0.

RO-GRPO Routing Reward Parameters.

The most critical hyperparameters are those governing the mechanism-aware routing reward, $R_{\rm route}$. We conducted a grid search to determine the optimal settings for both our adaptive strategies, using a held-out validation set.

Parameter	Value
GRPO Configuration	
Learning Rate	1×10^{-5}
KL Coefficient (β)	0.1
Batch Size	64
Generations per Prompt (k)	8
Epochs (Unimodal)	1
Epochs (Multimodal)	3

LoRA / LoRA-MoE Configuration

LoRA Rank (r)	16
LoRA-MoE Rank (r)	8 (per expert)
Number of Experts (E)	{2,4,8}
LoRA Alpha (α)	32
LoRA Dropout	0.05

RO-GRPO Specific (Optimal Values)

Global Routing Weight (Wroute)	0.2
Smooth Strategy	
λ_H^{start} (Entropy Weight Start)	0.5
λ_B^{end} (Balance Weight End)	2.0
Sigmoid Steepness (k)	20.0
Sigmoid Center (c)	0.5
Relative Strategy	
History Window Size (S_{hist})	1000

History Window Size (Signature Constant (C)

1.0 controls the tran

Note: Sigmoid steepness (k) controls the transition speed of the curriculum, and the center (c) defines the transition point in terms of training progress.

Table 4: Hyperparameters for all experiments.

For the Curriculum-Based Reward Scheduling (Smooth) strategy, we explored the key parameters controlling the curriculum's shape and intensity. The search space included the final load bal-

ancing weight $\lambda_B^{\text{end}} \in \{1.0, 2.0, 5.0\}$, the initial entropy weight $\lambda_H^{\text{start}} \in \{0.5, 1.0\}$, and the sigmoid steepness $k \in \{15, 20, 25\}$. The pseudocode for this strategy is detailed in Algorithm 2.

For the **Relative Improvement Gating (Relative)** strategy, the key parameter is the history_size, which defines the window for the moving average baseline. We searched over values in $\{100, 500, 1000\}$. The logic for this strategy is presented in Algorithm 1.

A separate grid search was performed for the global routing reward weight, which scales the entire R_{route} term, across the range $\{0.1, 0.2, 0.5\}$. Our experiments indicated that a weight of **0.2** provided the best trade-off between improving task accuracy and optimizing routing efficiency (i.e., minimizing load balancing MSE and routing entropy). This value was used for all reported RO-GRPO results.

The final, optimal hyperparameters selected through this process are summarized in Table 4.

Algorithm 1 RO-GRPO Reward Calculation (Relative Strategy)

```
Input: Routing statistics stats, history buffer B_{\text{hist}}
```

Parameters: Reward constant C, history buffer size S_{hist}

{Require sufficient history to establish a baseline}

```
\begin{aligned} & \textbf{if } |B_{\text{hist}}| < S_{\text{hist}} \textbf{ then} \\ & \textbf{return } 0 \\ & \textbf{end if} \end{aligned}
```

```
{Compute metrics for the current sample} (\mathcal{M}_{curr}, \bar{\mathcal{H}}_{curr}) \leftarrow ComputeMetrics(stats)
```

```
{Compute historical average baseline} \mathcal{M}_{hist} \leftarrow \text{Average}(B_{hist}.\text{mse}) \mathcal{H}_{hist} \leftarrow \text{Average}(B_{hist}.\text{entropy})
```

```
 \begin{aligned} & \{ \text{Grant reward only if both metrics improve} \} \\ & \text{if } \mathcal{M}_{\text{curr}} < \mathcal{M}_{\text{hist}} \text{ and } \bar{\mathcal{H}}_{\text{curr}} < \bar{\mathcal{H}}_{\text{hist}} \text{ then} \\ & R_{\text{route}} \leftarrow C \\ & \text{else} \\ & R_{\text{route}} \leftarrow 0 \end{aligned}
```

{Update the history buffer with current metrics} Update(B_{hist} , (\mathcal{M}_{curr} , $\bar{\mathcal{H}}_{curr}$))

```
return R_{\text{route}} = 0
```

end if

Algorithm 2 RO-GRPO Reward Calculation (Smooth Strategy)

```
Input: Routing statistics stats, current step t_{
m curr}, max steps t_{
m max}
```

Parameters: Global weight w_{route} , entropy start weight λ_B^{start} , balance end weight λ_B^{end} , sigmoid center c, sigmoid steepness k

{Calculate curriculum progress and sigmoid value}

$$p \leftarrow t_{\text{curr}}/t_{\text{max}}$$
$$\sigma \leftarrow (1 + e^{-k(p-c)})^{-1}$$

```
{Schedule the weights for entropy and balance} w_H \leftarrow \lambda_H^{\text{start}} \cdot (1 - \sigma) w_B \leftarrow \lambda_B^{\text{end}} \cdot \sigma
```

{Compute average normalized metrics from stats}

```
\bar{\mathcal{H}}_{norm} \leftarrow AverageNormalizedEntropy(stats)
\mathcal{M}_{norm} \leftarrow AverageNormalizedMSE(stats)
```

```
{Calculate final reward (a negative penalty)} R_{\text{route}} \leftarrow -w_{\text{route}} \cdot (w_H \cdot \bar{\mathcal{H}}_{\text{norm}} + w_B \cdot \mathcal{M}_{\text{norm}})
```

return $R_{\text{route}} = 0$

C System Prompt

We used a consistent system prompt for all models during both training and evaluation to encourage step-by-step reasoning. The prompt is provided below:

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <think> </think> and<answer> </answer>tags, respectively, i.e., <think> reasoning process here

Aspect	Content				
Problem	Mary is an avid gardener. Yesterday, she received 18 new potted plants from her favorite plant nursery. She already has potted plants on each of the 40 window ledges of her large country home. Feeling generous, she has decided that she wi give 1 potted plant from each ledge to friends and family tomorrow. How many potted plants will Mary remain with? Answer:				
Solution	Yesterday, before receiving the plants, Mary had $2\times 40=80$ potted plants. After receiving an additional 18 plants, she therefore had a total of $80+18=98$ potted plants. Tomorrow, Mary's plant giveaway will be $40\times 1=40$ potted plants. She will therefore remain with $98-40=58$ potted plants. Final Answer: 58				
Base Output	<think> First, we calculate the total number of potted plants Mary already has. Since she has 2 potted plants on each of the 40 window ledges, the total is $2 \times 40 = 80$ potted plants. Next, we determine how many potted plants she will give away. Since she plans to give 1 potted plant from each of th 40 window ledges, she will give away $1 \times 40 = 40$ potted plants. Finally, we subtract the number of potted plants she gives away from the total number she originally had. So $80 - 40 = 40$. Therefore, after giving away the potted plants, Mary will remain with 40 potted plants. <answer>40 <answ< td=""></answ<></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></think>				
GRPO (LoRA- MoE) Out- put					
RO-GRPO (Smooth) Output	<think> First, we calculate the total number of potted plants Mary already has. Since she has 2 potted plants on eac of the 40 window ledges, the total is $2 \times 40 = 80$ potted plants. Next, we consider the 18 new potted plants she received yesterday. Adding these to her existing plants gives $80+18 = 9$ potted plants in total. Mary plans to give away 1 potted plant from each of the 40 window ledges. This means she will give away $40 \times 1 = 4$ potted plants. To find out how many potted plants she will have left after giving some away, we subtract the number of plants give away from the total number of plants she has: $98 - 40 = 58$. <th< td=""></th<></think>				
RO-GRPO (Relative) Output	<think> Mary initially has 2 potted plants on each of the 40 window ledges. So, the total number of potted plants shas before receiving the new ones is: 2 plants/ledge × 40 ledges = 80 plants After receiving 18 new potted plants, the total number of potted plants becomes:</think>				
	80 plants + 18 plants = 98 plants				
	Mary plans to give away 1 potted plant from each of the 40 window ledges. Therefore, the number of plants she will give away is:				
	1 plant/ledge \times 40 ledges = 40 plants				
	Subtracting the number of plants given away from the total number of plants gives us the number of plants Mary wi remain with:				
	98 plants -40 plants $=58$ plants				
	<pre> <answer> Mary will remain with 58 potted plants. </answer></pre>				
Analysis	This case study reveals a critical reasoning failure in the baselines. Both the Base Model and GRPO (LoRA-Mofignore the 18 new potted plants, leading to an incorrect answer. In contrast, our RO-GRPO models successfull integrate all information to derive the correct solution. This demonstrates that our mechanism-aware rewards foster robureasoning and address the core logical deficiencies observed in the baselines, moving beyond superficial format adherence				

Table 5: Case study comparing model outputs on a sample from the GSM8K benchmark.

D QUALITATIVE CASE STUDY

To illustrate the practical impact of our method on reasoning quality, we present a case study on a multi-step arithmetic problem from the GSM8K dataset. Table 5 compares the outputs generated by the baseline models against those from our RO-GRPO framework.

E DETAILED THEORETICAL ANALYSIS

This appendix provides the detailed mathematical derivations and expanded interpretations for the theoretical analysis.

E.1 CONSTRAINED OPTIMIZATION INTERPRETATION

The RO-GRPO framework can be viewed as a practical, penalty-based approach to solving a constrained policy optimization problem. The objective is to maximize the expected task reward, subject to constraints on the policy's internal routing behavior:

$$\max_{\theta} \qquad \mathbb{E}_{y \sim \pi_{\theta}}[R_{\text{task}}(y)]$$
subject to
$$\mathbb{E}_{y \sim \pi_{\theta}}[\bar{\mathcal{H}}_{\text{norm}}(y)] \leq \varepsilon_{H},$$

$$\mathbb{E}_{y \sim \pi_{\theta}}[\mathcal{M}_{\text{norm}}(y)] \leq \varepsilon_{M},$$
(9)

where ε_H and ε_M are desired thresholds for the average normalized routing entropy and load balancing MSE, respectively.

The standard method for solving such a problem is via its Lagrangian relaxation. The Lagrangian $L(\theta, \lambda_H, \lambda_M)$ is:

$$L = \mathbb{E}[R_{\text{task}}] - \lambda_H(\mathbb{E}[\bar{\mathcal{H}}_{\text{norm}}] - \varepsilon_H) - \lambda_M(\mathbb{E}[\mathcal{M}_{\text{norm}}] - \varepsilon_M),$$
(10)

where $\lambda_H, \lambda_M \geq 0$ are the Lagrange multipliers. Our RO-GRPO objective can be expressed as:

$$\max_{\theta} \mathbb{E}\left[R_{\text{task}} - w_{\text{route}}\left(w_H(t) \cdot \bar{\mathcal{H}}_{\text{norm}} + w_B(t) \cdot \mathcal{M}_{\text{norm}}\right)\right]. \tag{11}$$

Comparing our objective in Eq. equation 11 with the Lagrangian in Eq. equation 10 reveals that RO-GRPO maximizes a simplified Lagrangian. The weights $w_H(t)$ and $w_B(t)$ function as fixed (or scheduled) Lagrange multipliers, and the constraint thresholds $\varepsilon_H, \varepsilon_M$ are implicitly absorbed into the objective. This formulation positions RO-GRPO as a **fixed-penalty method**, which gains significant simplicity by integrating the constraints directly into the scalar reward signal.

E.2 Information Bottleneck Principle

The efficacy of the entropy reward, R_H , can be formally interpreted through the lens of the Information Bottleneck (IB) principle. The IB principle seeks a compressed representation \tilde{X} of an input X that minimizes the mutual information $I(X;\tilde{X})$ while maximizing the mutual information $I(\tilde{X};Y)$ with a relevance variable Y.

We map our routing mechanism to this framework: the input token representation is the source variable $(X = \mathbf{h})$; the router's output probability distribution is the compressed representation $(\tilde{X} = \mathbf{p})$; and the final sequence quality (proxied by R_{task}) is the relevance variable (Y). The mutual information $I(\mathbf{h}; \mathbf{p}) = H(\mathbf{p}) - H(\mathbf{p}|\mathbf{h})$ quantifies the routing policy's complexity. Our confidence reward R_H directly encourages the minimization of the conditional entropy $H(\mathbf{p}|\mathbf{h})$. The complete RO-GRPO objective, which maximizes $\mathbb{E}[R_{\text{task}} + w_H R_H]$, is thus analogous to the IB objective. It balances maximizing task-relevant information (via R_{task}) with learning a compressed, minimal sufficient representation of the input (via R_H).

E.3 PARAMETER UTILIZATION AND VARIANCE MINIMIZATION

The balancing reward, R_B , is grounded in a direct mathematical relationship with load distribution variance. We show that minimizing our MSE-based metric is equivalent to minimizing the variance of expert utilization.

Let $\bar{\mathbf{p}}$ be the empirical utilization vector over E experts. The variance of this distribution is:

$$\operatorname{Var}(\bar{\mathbf{p}}) = \frac{1}{E} \sum_{e=1}^{E} (\bar{p}_e - \mathbb{E}[\bar{\mathbf{p}}])^2. \tag{12}$$

Since $\sum \bar{p}_e = 1$, the mean utilization $\mathbb{E}[\bar{\mathbf{p}}] = 1/E$. Substituting this gives:

$$\operatorname{Var}(\bar{\mathbf{p}}) = \frac{1}{E} \sum_{e=1}^{E} \left(\bar{p}_e - \frac{1}{E} \right)^2. \tag{13}$$

This expression is precisely the Mean Squared Error (MSE) between the empirical distribution $\bar{\mathbf{p}}$ and a uniform distribution $\mathbf{u} = (1/E, \dots, 1/E)$. Therefore, $Var(\bar{\mathbf{p}}) = MSE(\bar{\mathbf{p}}, \mathbf{u})$.

This equivalence establishes that maximizing our reward $R_B \propto -\mathrm{MSE}(\bar{\mathbf{p}}, \mathbf{u})$ is directly proportional to minimizing the variance of the expert load. This provides a principled and efficient mechanism to promote balanced parameter usage within the RL framework.