

Beyond Length Scaling: Synergizing Breadth and Depth for Generative Reward Models

Anonymous ACL submission

Abstract

Recent advancements in Generative Reward Models (GRMs) have demonstrated that scaling the length of Chain-of-Thought (CoT) reasoning considerably enhances the reliability of evaluation. However, current works predominantly rely on unstructured length scaling, ignoring the divergent efficacy of different reasoning mechanisms: Breadth-CoT (multi-dimensional principle coverage) and Depth-CoT (substantive judgment soundness). To address this, we introduce **Mix-GRM**, a framework that reconfigures raw rationales into structured Breadth-CoT and Depth-CoT through a modular synthesis pipeline, subsequently employing Supervised Fine-Tuning (SFT) and Reinforcement Learning with Verifiable Rewards (RLVR) to internalize and optimize these mechanisms. Comprehensive experiments demonstrate that Mix-GRM establishes a new state-of-the-art across five benchmarks, surpassing leading open-source RMs by an average of 8.2%. Our results reveal a clear divergence in reasoning: Breadth-CoT benefits subjective preference tasks, whereas Depth-CoT excels in objective correctness tasks. Consequently, misaligning the reasoning mechanism with the task directly degrades performance. Furthermore, we demonstrate that RLVR acts as a switching amplifier, inducing an emergent polarization where the model spontaneously allocates its reasoning style to match task demands.

1 Introduction

Reinforcement learning (RL) has proven to be the critical post-training mechanism for eliciting capabilities in Large Language Models (LLMs) (Ouyang et al., 2022; Team, 2025a,b). However, as the ambition of RL expands from single-domain optimization (e.g., math) (Le et al., 2022; Shao et al., 2024; Wang et al., 2025a) to general-purpose alignment (Lee et al., 2024; Shen et al., 2025), the Reward Model (RM) faces the

challenge of providing reliable feedback for increasingly complex queries from diverse, real-world scenarios (Liu et al., 2025d; Li et al., 2025a). Addressing this challenge requires a shift in RM design. Inspired by how CoT (Wei et al., 2023; Yeo et al., 2025) trades inference-time compute for enhanced generalization performance, the community has increasingly adopted Generative Reward Models (GRMs) (Zheng et al., 2023; Yuan et al., 2024; Zhang et al., 2025a). By prompting an explicit evaluation rationale prior to conclusion, GRMs aim to transfer the robust generalization observed in CoT generation to the task of reward modeling.

Building on these successes, existing GRM methods predominantly leverage CoT by simply scaling its length (Chen et al., 2025b,a; Zhang et al., 2025c), feeding it with massive evaluation signals, such as fine-grained features (Kim et al., 2024) or multi-perspective critiques (Ankner et al., 2024). However, prior CoT studies (Sprague et al., 2025; Besta et al., 2025; Wang et al., 2024b; Kambhampati et al., 2024) have established that longer CoTs do not universally guarantee performance gains; rather, the optimal structural bias diverges significantly across domains. Crucially, recent insights from test-time scaling (Li et al., 2025b; Zhang et al., 2025b) provide a theory for this divergence, identifying *parallel thinking* and *sequential thinking* as two fundamental, orthogonal mechanisms for amplifying intelligence. Conceptually, reasoning-heavy tasks (e.g., math, code) necessitate sequential verification to ensure deductive rigor (Wang et al., 2024a; Liu et al., 2025a; Lightman et al., 2023), whereas semantic-heavy tasks (e.g., open-ended generation) benefit from parallel exploration to ensure comprehensive coverage of diverse possibilities (Zheng et al., 2025; Pan et al., 2025).

Drawing on this distinction, we argue that advancing RM requires shifting focus from merely scaling CoT length to aligning its reasoning mechanisms with task demands. Specifically, this ne-

084	cessitates a transition from static, one-size-fits-all	scaling how long it writes.	136
085	CoT templates toward a <i>mix reasoning mechanism</i> .		
086	Thus, we propose Mix-GRM , which implements a	2 Related Work	137
087	dynamic mix reasoning mechanism within a unified	2.1 Generative Reward Model	138
088	reward modeling framework. Specifically, we intro-	Generative Reward Models represent a paradigm	139
089	duce a synthesis framework that reconfigures raw,	shift from scalar regression to explicit reasoning.	140
090	unstructured rationales into two distinct long CoTs:	Developing alongside the prompting-based “LLM-	141
091	Breadth-CoT (B-CoT) and Depth-CoT (D-CoT).	as-a-Judge” paradigm (Zheng et al., 2023), GRMs	142
092	To achieve this, we first decouple unstructured ra-	are explicitly trained to generate natural language	143
093	tionales into atomic “Principle–Judgment–Verdict”	rationales alongside preference decisions (Yuan	144
094	units. This modularity allows us to reassemble the	et al., 2024). Driven by the transformative success	145
095	units into syntactically unified but structurally di-	of long CoT, the research trajectory in this field has	146
096	verse paths. To illustrate, a B-CoT is synthesized	pivoted toward continuously extending the length	147
097	by the parallel aggregation of units across diverse	of these rationales. To achieve this, many work	148
098	principles (e.g., combining an ‘Accuracy’ unit with	leverages RL to explicitly elicit and stabilize longer	149
099	a ‘Clarity’ unit) to ensure coverage. Conversely,	CoT traces (Chen et al., 2025b,a; Whitehouse et al.,	150
100	D-CoT extends the CoT by first performing a di-	2025), while complementary efforts utilize detailed	151
101	rect reasoning pass to solve the instruction, thereby	rubrics/checklists to synthetically expand evalua-	152
102	enabling a re-evaluated judgment grounded in the	tion coverage (Kim et al., 2024; Liu et al., 2025b;	153
103	generated reasoning pass to ensure soundness. To	Gunjal et al., 2025; Viswanathan et al., 2025). How-	154
104	cultivate mechanism-adaptive alignment, we con-	ever, while these strategies successfully scale the	155
105	struct a synergistic mixture dataset by pairing B-	quantity of reasoning, they typically rely on static,	156
106	CoT with subjective preference tasks and D-CoT	task-agnostic structures, overlooking the critical	157
107	with objective correctness tasks. We first initial-	nuance that the optimal reasoning mechanism is	158
108	ize the model via SFT on this mixture and subse-	intrinsically task-dependent.	159
109	quently optimize it through RLVR using normal		
110	RM datasets, where only final labels are available.	2.2 Breadth and Depth in Chain-of-Thought	160
111	Comprehensive experiments across five stan-	The evolution of CoT is fundamentally charac-	161
112	dard benchmarks yield three critical conclusions:	terized by the continuous exploration of diverse	162
113	(1) Universal SOTA Performance and Down-	structures (Shinn et al., 2023; Team, 2025a). Be-	163
114	stream Utility: <i>Mix-GRM</i> establishes a new state-	yond simple linear chains, frameworks such as	164
115	of-the-art, consistently surpassing strong baselines	Tree of Thoughts (Yao et al., 2023) and Graph of	165
116	like <i>Skywork-Reward</i> and <i>FARE-8B</i> on general re-	Thoughts (Besta et al., 2025) introduce branching	166
117	ward benchmarks. Crucially, this superiority ex-	and recurrent topologies, framing reasoning as a	167
118	tends to practical downstream tasks: <i>Mix-GRM</i>	structured search over partial thoughts. Comple-	168
119	demonstrates best-in-class utility in both Offline	menting these complex structures, approaches like	169
120	RL (DPO) and Test-time Scaling (Best-of-N). (2)	Skeleton-of-Thought (Ning et al., 2024) and Self-	170
121	Divergent Roles of Reasoning Mechanisms: Our	Consistency (Wang et al., 2023) demonstrate the	171
122	analysis reveals that B-CoT predominantly ben-	efficacy of parallel exploration, leveraging lateral	172
123	efits subjective preference but degrades objective	breadth to enhance robustness and coverage. Col-	173
124	correctness, while D-CoT excels in correctness	lectively, these studies establish that reasoning is	174
125	at the cost of preference. This confirms that	not structure-agnostic; rather, specific topological	175
126	the efficacy of a reasoning mechanism is task-	priors—ranging from deep sequential trees to broad	176
127	dependent. (3) RLVR as a Switching Amplifier:	parallel ensembles—are required to unlock opti-	177
128	Mixed mechanisms provide a superior base for	mal performance across distinct domains (Sprague	178
129	RL. RLVR boosts <i>Mix-GRM</i> by a larger margin	et al., 2025), a distinction that our work formally	179
130	than the <i>Base-GRM</i> . Our analysis demonstrates	adapts to reward modeling.	180
131	that RLVR automatically sharpens the mechanism		
132	allocation—spontaneously converging on B-CoT	3 Methodology	181
133	for preference and D-CoT for correctness. This	We propose the <i>Mix-GRM</i> , a framework designed	182
134	confirms that optimizing how a model thinks is	to dynamically align the reasoning mechanism	183
135	more critical for post-training efficacy than simply		

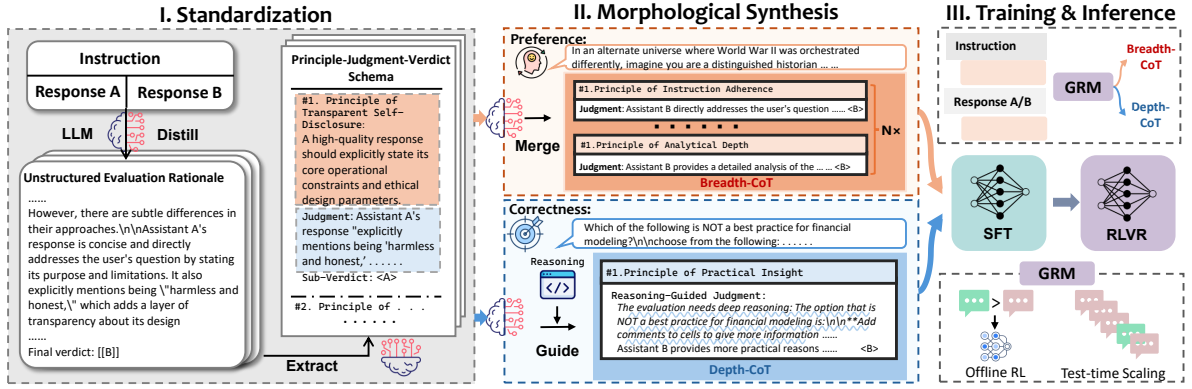


Figure 1: The pipeline of the Mix-GRM. (i) **Standardization**: We extract raw rationales into modular *Principle–Judgment–Verdict* units. (II) **Mechanism Synthesis**: We reconstruct modules into *Breadth-CoT* for preference or *Depth-CoT* for correctness. (III) **Training & Inference**: Following SFT and RLVR training, the model achieves mechanism-adaptive alignment, automatically deploying the optimal mechanism for inference and providing reliable signals for downstream tasks like Offline RL and test-time scaling.

with intrinsic task demands. Moving beyond static, unstructured rationale sequences, our approach formalizes evaluation into two orthogonal CoTs: **Breadth-CoT**, which enforces the lateral aggregation of diverse principles, and **Depth-CoT**, which necessitates the expansion of judgment. As illustrated in Figure 1, our methodology comprises three key phases: modular schema standardization (§3.2), mechanism synthesis (§3.3), and mechanism-adaptive alignment (§3.4).

3.1 Problem Formulation

Supposing $\{y_A, y_B\}$ denote two candidate responses generated by two assistants A and B for a given task instruction x , a normal GRM \mathcal{M} produces an output sequence consisting of an explicit evaluation rationale c followed by a preference verdict v , comparing the quality of y_A and y_B .

$$(c, v) = \mathcal{M}(y_A, y_B | x).$$

The objective is to ensure that the v aligns with human preference. In our framework, we denote the full input triplet as $I = (x, y_A, y_B)$.

3.2 Modular Schema Standardization

Conventional GRMs typically produce the rationale c as an unstructured, free-form sequence. Inspired by recent checklist-based evaluation (Viswanathan et al., 2025), which advocates for the atomization of the complex evaluation process into checklist-driven points, we propose to reconfigure these raw rationales into a structured “Principle–Judgment–Verdict” Schema (Figure 1, Stage I). By transforming tangled rationales into atomic units, we en-

sure that the RM’s reasoning process is both interpretable and granularly verifiable. Formally, we utilize a LLM to parse the raw c into structured atomic units \mathcal{S} :

$$\mathcal{S} = \{(p_k, j_k, v_k)\}_{k=1}^K,$$

where p_k denotes a discrete evaluation **Principle** (e.g., “Instruction Adherence”), j_k represents the specific **Judgment** (e.g., “Response B directly addresses...”) analyzing that principle, and v_k is the following **Sub-Verdict** (e.g., “ is Better”). Here, K typically ranges from 3 to 5.

This atomic decomposition yields cleaner learning signals and ensures syntactic uniformity (Li et al., 2025c), ensuring that performance gains are driven by thinking mechanisms (i.e., Breadth vs. Depth) rather than superficial stylistic patterns.

3.3 Mechanism Synthesis

Building on the \mathcal{S} , we introduce a dual-track synthesis pipeline (Figure 1, Stage II) to synthesize **B-** and **D-CoT** as follows:

B-CoT Synthesis. We define B-CoT as the parallel aggregation of distinct principles, designed to overcome the narrow focus of single-pass rationale. In subjective preference tasks, where a “good” response is defined by the simultaneous satisfaction of multi-dimensional factors (e.g., tone, helpfulness, and creativity), single-track reasoning often fixates on dominant traits while overlooking subtle, fine-grained details. By exploring diverse reasoning paths concurrently, parallel thinking provides a deliberative breadth that aligns with the

multifaceted nature of human preference. To simulate parallel thinking, we treat independent sampling as a stochastic exploration of the instruction’s evaluative manifold. By sampling N independent rationales $\{c_n\}_{n=1}^N$ from multiple cognitive trajectories, we elicit a diverse set of hidden principles that might otherwise remain dormant. These rationales are parsed into structured schemas $\{\mathcal{S}_n\}$ and subsequently unified via an LLM-based **Merge & Deduplicate** transformation $\mathcal{T}_{\text{merge}}$:

$$C_{\text{breadth}} = \mathcal{T}_{\text{merge}} \left(\bigcup_{n=1}^N (p, j, v) \in \mathcal{S}_n \right).$$

Here, we filter out lowest-frequency principles. This synthesis yields a comprehensive, non-redundant spectrum of principles, effectively expanding the model’s horizontal evaluative scope.

D-CoT Synthesis. We define D-CoT as the expansion of judgment to ensure substantive reasoning soundness by mitigating superficial shortcuts. In contrast to subjective preferences, a “good” response in objective correctness tasks depends on rigorous logical constraints (*e.g.*, mathematical proofs or functional code). Normal rationales often fixate on surface-level fluency (*e.g.*, professional tone or formatting) while failing to verify the underlying logical validity. By enforcing the sequential verification of logical dependencies, sequential thinking provides a deductive rigor that naturally aligns with the strict requirements of objective correctness. To simulate sequential thinking, we first elicit a Reasoning Trace z —a self-solving pass derived from x that explicitly outlines the optimal solution paths required for a correct response. Recognizing that depth-oriented reasoning demands higher cognitive load per unit, we intentionally trade off horizontal coverage for deductive rigor by sampling a focused subset $\mathcal{S}_{\text{sub}} \subset \mathcal{S}$ (typically $|K| \leq 3$). In this stage of **Reasoning-Guided Judgment**, each unit’s judgment is re-generated as a derivative of the trace z :

$$\tilde{j}_k = \mathcal{T}_{\text{refine}}(p_k | z)$$

To ensure the evaluative process is transparent and explicitly grounded in the model’s own logic, we inject z directly into the lead unit \tilde{j}_1 . The final C_{depth} is constructed by serializing these refined units, transforming the verdict into a substantive analytical process anchored by the trace z .

3.4 Mechanism-Adaptive Alignment

Training proceeds in two stages (Figure 1, Panel III): SFT on mixture CoT datasets, followed by GRPO (Shao et al., 2024) to align verdicts with human labels.

SFT. Following Frick et al. (2025), we categorize general RM training data into two domains: **Preference** (subjective) and **Correctness** (objective). We construct the mixture dataset \mathcal{D}_{mix} by assigning C_{breadth} to preference instances and C_{depth} to correctness instances. We first initialize the policy π_θ via SFT on \mathcal{D}_{mix} . Given the I , the model is trained to generate the corresponding CoT $c \in \{c_{\text{breadth}}, c_{\text{depth}}\}$ alongside the verdict v .

RLVR via GRPO. To optimize verdict accuracy, we employ RLVR via GRPO (Shao et al., 2024), rewarding the model solely for consistency with ground-truth labels:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\substack{I \sim \mathcal{D} \\ \{o_i\} \sim \pi_{\theta_{\text{old}}}}} \left[\frac{1}{G} \sum_{i=1}^G \left(\frac{\pi_\theta(o_i | I)}{\pi_{\theta_{\text{old}}}(o_i | I)} \hat{A}_i - \beta \mathbb{D}_{\text{KL}}(\pi_\theta || \pi_{\text{ref}}) \right) \right]$$

This process acts as a **switching amplifier**, inducing an emergent polarization: the model spontaneously learns to couple B-CoT with preference tasks and D-CoT with correctness tasks to maximize rewards, as empirically verified in §5. This confirms that the model autonomously converges on the optimal thinking style for each domain.

4 Experiment

We evaluate Mix-GRM across three objectives: (1) **Overall Performance** against SoTA baselines; (2) **Mechanism Efficiency** to quantify the domain-specific benefits of B- and D-CoT; and (3) **Downstream Utility** in Offline RL and Test-time Scaling.

4.1 Experimental Setup

General Reward Benchmarks. We employ five widely recognized benchmarks tailored for general-purpose reward modeling: RewardBench (Lambert et al., 2024), RewardBench-v2 (Malik et al., 2025), RMB (Zhou et al., 2025a), RM-Bench (Liu et al., 2025c), and PPE (Frick et al., 2025). These benchmarks encompass a broad spectrum of tasks, ranging from common tasks like math, coding, and

Models	Stage	Data	Benchmarks					
			RB-v1 [†]	RB-v2 [†]	RM-BENCH	RMB	PPE	Avg.
<i>Reference: Proprietary Models</i>								
DeepSeek-V3.2	–	–	95.5	92.1	91.4	83.9	69.0	86.4
Gemini-3-Flash	–	–	95.3	91.1	93.8	79.2	76.4	87.2
<i>Open-Source Baselines</i>								
Skywork-Reward-8B	BT	44K	93.9	79.7	72.4	74.4	61.7	76.5
JudgeLRM-7B	RL	100K	79.0	55.6	78.5	73.1	57.9	68.8
RM-R1-7B (Distill)	SFT, RL	9K, 64K	83.5	48.7	76.6	65.1	62.0	67.2
RM-R1-7B (Instruct)	SFT, RL	9K, 64K	82.3	61.4	75.1	69.9	62.0	70.1
FARE-8B	SFT	2.5M	86.3	73.4	74.1	83.2	62.5	75.9
RubricRM-8B	SFT	36K	86.7	71.9	74.0	78.5	62.5	74.7
DeepSeek-GRM-16B	SFT, RL	1.2M, 237K	76.8	56.0	63.5	70.8	59.1	65.2
<i>Ours: Mix-GRMs</i>								
<i>Stage I: SFT-trained</i>								
Base-GRM	SFT	9K	84.5	64.7	77.0	79.2	61.1	73.3
Mix-GRM (Ours)	SFT	9K	87.2	67.8	<u>79.2</u>	78.9	62.1	75.1
<i>Stage II: RLVR-trained</i>								
Base-GRM	SFT, RL	9K, 21K	89.0	74.0	78.8	78.5	<u>64.0</u>	<u>76.9</u>
Mix-GRM (Ours)	SFT, RL	9K, 21K	<u>91.8</u>	<u>77.5</u>	82.7	<u>80.1</u>	64.8	79.4

Table 1: Performance of RMs on reward benchmarks. Among open-source models, the highest score per column is **bolded**, and the second-highest is underlined. “Overall” denotes the average score within each benchmark. Proprietary LLMs (gray rows) are included for reference. [†]RB-v1/v2 refers to RewardBench v1 and v2.

open-ended chat, to specialized capabilities including factuality and instruction-following. For Overall Performance, we report standard benchmark-level pairwise comparison accuracy to assess the rewarding capability. For granular Mechanism Efficiency analysis, we aggregate instances from these benchmarks and re-categorize them into two fundamental domains, Correctness and Preference, based on their original task metadata. Detailed statistics and specific domain mappings for these benchmarks are provided in Appendix B.3.

Base Model and Training Data Source. We employ Qwen3-8B-Base (Team, 2025c) trained on a composite corpus 30,000 samples (9K SFT, 21K RLVR) spanning five datasets: HelpSteer3 (Wang et al., 2025b) (chat, stem & multilingual), Code-Preference (coding), Math-DPO (math), WildGuard (Han et al., 2024) (safety), and OffsetBias (Park et al., 2024) (instruction following). Detailed sampling protocols and statistical distributions are provided in Appendix B.2. Other Training Implementation Setting is in Sec. B.

Baselines. We compare our proposed RM with 7 top-tier RMs across two paradigms: (1) *Discriminative*: represented by **Skywork-Reward-v0.2-Llama-3.1-8B** (Liu et al., 2024), a leading scalar model trained via Bradley-Terry modeling; and

(2) *Generative*: encompassing RL-driven reasoning models (**JudgeLRM-7B** (Chen et al., 2025a), **RM-R1-Instruct** (Chen et al., 2025b), **RM-R1-Distill**, **DeepSeek-GRM-16B**) (Liu et al., 2025d), synthetic scaling methods (**FARE-8B** (Xu et al., 2025)), and rubric-based approaches **RubricRM-8B** (Liu et al., 2025b). Notably, RubricRM-8B incorporates two-stage LLMs consisting of a rubric generator and a rubric-based judge.

4.2 Overall Performance in Benchmarks

Table 1 validates the effectiveness of our Mix-GRM through three dimensions.

Effectiveness of Mixture SFT : Via mixture SFT alone, *Mix-GRM* achieves a remarkable average score of 75.1. This performance surpasses GRMs requiring computationally expensive RL to elicit long-CoT capabilities—outperforming *RM-R1-Instruct* by 5.0 and *DeepSeek-GRM-16B* by 9.9. Furthermore, it beats *RubricRM-8B* (+0.4), which relies on a complex but static rubric-template CoT. This confirms that aligning reasoning mechanisms serves as a potent alternative strategy, alongside approaches focused on RL exploration or static template engineering.

Superiority of Data Efficiency : *Mix-GRM* achieves these gains with substantially less data. While *FARE-8B* relies on massive scaling ($\approx 2.5M$

Models	Preference Domain						Correctness Domain						Overall
	RB-v1	RB-v2	RM-B [†]	RMB	PPE	Avg.	RB-v1	RB-v2	RM-B [†]	RMB	PPE	Avg.	
<i>Baselines</i>													
FARE-8B	85.0	57.3	66.9	82.9	59.6	70.4	85.2	67.3	63.0	88.1	63.3	73.3	71.9
RubricRM-8B	82.4	56.0	62.2	77.5	64.9	68.6	87.6	64.2	57.6	86.5	60.4	71.3	70.0
DeepSeek-GRM	80.6	<u>59.6</u>	64.0	76.8	59.8	68.2	76.6	55.8	56.6	86.8	56.8	66.5	67.4
<i>Ours: Mix-GRMs</i>													
<i>Stage I: SFT-trained</i>													
Base-GRM	81.6	55.5	63.3	<u>80.5</u>	60.1	68.2	84.1	63.7	67.7	86.4	59.1	72.2	70.2
Mix-GRM (Breadth)	83.7	59.1	65.9	77.9	59.5	69.3 ^{↑1.1}	81.1	60.2	64.1	86.8	58.7	70.2 ^{↓2.0}	69.8
Mix-GRM (Depth)	80.3	50.2	70.6	70.1	58.6	65.9 ^{↓2.3}	88.0	63.7	66.7	81.1	64.7	72.8 ^{↑0.6}	69.4
Mix-GRM	84.9	55.7	71.2	78.7	59.2	70.0 ^{↑1.8}	88.4	65.8	67.7	81.9	63.7	73.5 ^{↑1.3}	71.8
<i>Stage II: RLVR-trained</i>													
Base-GRM	83.0	58.0	68.5	73.8	61.4	68.9 ^{↑0.7}	89.8	69.5	69.9	89.5	63.4	76.4 ^{↑4.2}	72.7
Mix-GRM (Breadth)	86.2	58.8	70.1	79.2	60.7	71.0 ^{↑2.8}	82.8	63.4	64.3	86.5	60.7	71.5 ^{↓0.7}	71.3
Mix-GRM (Depth)	85.2	57.8	75.6	75.4	61.2	71.0 ^{↑2.8}	91.8	70.3	72.9	87.4	66.2	77.7 ^{↑5.5}	74.4
Mix-GRM	86.2	64.4	<u>72.7</u>	78.1	<u>61.7</u>	72.6 ^{↑3.7}	92.2	72.5	74.5	<u>88.9</u>	<u>65.4</u>	78.7 ^{↑6.5}	75.7

Table 2: Performance of RMs grouped by domain. ‘‘Avg.’’ denotes the domain average. We annotate the performance gap relative to the *Base-GRM in SFT* baseline within the same stage using colored subscripts (\uparrow for gain, \downarrow for drop). Highest score per column is **bolded**, second-highest is underlined. [†]RM-B refers to RM-Bench.

Models	Instruction-Following			Mathematical Reasoning				
	ALPACA-V2	ARENA-HARD	Avg.	GSM8K	MATH	STEM	TABMWP	Avg.
SFT	6.4	4.2	5.3	75.1	25.2	38.6	40.9	45.0
<i>DPO Training (Different RMs)</i>								
\hookrightarrow RubricRM-8B	8.5	12.5	10.5	76.0	<u>26.9</u>	41.4	38.8	<u>45.9</u>
\hookrightarrow FARE-8B	<u>8.9</u>	15.1	<u>12.0</u>	75.7	<u>26.9</u>	39.0	41.4	45.8
\hookrightarrow RM-R1-Instruct	7.9	14.3	11.1	<u>76.3</u>	26.5	38.5	41.7	45.8
\hookrightarrow DeepSeek-GRM-16B	8.0	14.1	11.1	75.6	26.6	38.7	41.6	45.6
\hookrightarrow Ours (Mix-GRM)	9.2	<u>15.0</u>	12.1	77.6	27.1	<u>39.0</u>	41.9	46.4

Table 3: Performance of DPO-trained policy models using different reward models on instruction-following and math-reasoning benchmarks. ‘‘Avg.’’ is the average score of all benchmarks in each domain. In each column, the highest score is **bolded** and the second-highest is underlined.

387 samples) to reach 75.9, *Mix-GRM* attains a compa- 406
388 rable 75.1 in the SFT stage using merely 9K sam-
389 ples. This finding highlights that optimizing CoT
390 mechanisms yields a substantially higher training
391 signal density, enabling data efficiency compared
392 to brute-force dataset expansion.

393 **Switching Amplification via RLVR** : Mix CoT
394 maximizes the efficacy of the RLVR stage, unlock-
395 ing greater performance gains than unstructured
396 CoT. RLVR boosts *Mix-GRM* by 4.3 (75.1 \rightarrow
397 79.4), compared to a 3.6 gain for *Base-GRM*
398 (73.3 \rightarrow 76.9). Consequently, the performance
399 gap over the *Base-GRM* widens from 1.8 (SFT)
400 to 2.5 (RLVR), confirming that the aligned me-
401 chanism offers a more exploitable base for the RL. Fur-
402 thermore, our subsequent analysis (Sec 5) reveals
403 that these gains are fundamentally underpinned by
404 an emergent polarization in mechanism allocation,
405 where RLVR sharpens the model’s reasoning style

to match task-specific demands.

4.3 Mechanism Efficiency 407

408 Table 2 reveals that mechanism efficacy is strictly
409 task-dependent. In the SFT stage, we observe
410 a **distinct performance trade-off**: B-CoT im-
411 proves Preference via lateral coverage but degrades
412 Correctness (72.2 \rightarrow 70.2), whereas D-CoT en-
413 hances deductive soundness but fails in Preference
414 (68.2 \rightarrow 65.9). These results indicate that simply
415 extending CoT length does not guarantee universal
416 gains; while principle expansion facilitates multi-
417 dimensional evaluation, it offers no inherent ad-
418 vantage for deep reasoning. However, *Mix-GRM*
419 overcomes these limitations through a **synergistic**
420 **mutual enhancement**. By integrating orthogo-
421 nal strengths, it not only surpasses the *Base-GRM*
422 (70.2 \rightarrow 71.8) but surprisingly outperforms spe-
423 cialized single-mode models on their respective
424 strongholds (*e.g.*, exceeding Depth-only on Cor-

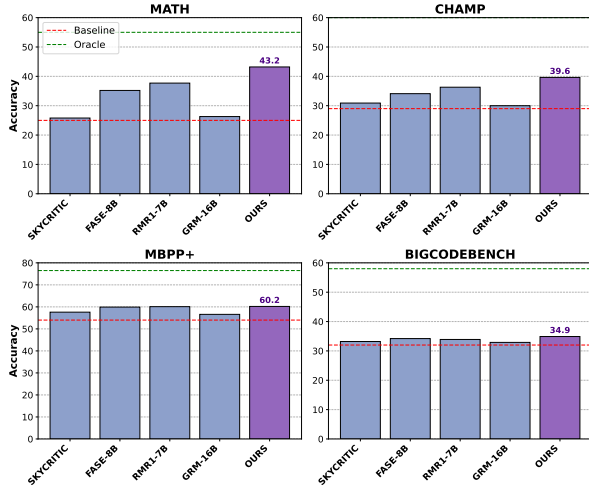


Figure 2: Best-of-10 performance across four challenging reasoning and coding benchmarks. Mix-GRM (ours) consistently achieves the highest accuracy across all tasks, effectively identifying solutions in both mathematical and code generation scenarios. Red and green lines denote random and oracle selection baselines.

rectness). This synergy becomes critical during the RLVR stage, where single-mode mechanisms encounter hard performance ceilings—most notably, *Mix-GRM (Breadth)* plateaus on correctness tasks. In contrast, the *Mix-GRM* enables RL optimization to reach a superior ceiling (78.7). This confirms that the **CoT structure itself acts as a bottleneck for RL optimization**; the mix structures does not merely inherit component strengths but constructs a robust reasoning framework that transcends the inherent limitations of isolated mechanisms.

4.4 Downstream Utility

To validate the practical utility of *Mix-GRM*, we apply it to two downstream applications: (i) serving as a reward signal for **Offline Reinforcement Learning**, and (ii) acting as a verifier for **Test-time Scaling**. We provide detailed descriptions of these application settings in Appendix C.

Reward Model for Offline Reinforcement Learning. In Offline RL via Direct Preference Optimization (DPO) (Rafailov et al., 2023), RM constructs high-quality preference pairs (y_w, y_l) to supervise policy alignment. Table 3 shows that models trained on these signals achieve a peak win rate of **12.1** in instruction-following, surpassing *FARE-8B* (12.0) and *RubricRM* (10.5). Crucially, this alignment gain does not compromise reasoning capabilities; in the math domain, *Mix-GRM* maintains a SOTA accuracy of **46.4**, edging out *RubricRM* (45.9) and *RM-R1-Instruct* (45.8). Specif-

ically, *Mix-GRM* achieves 77.6% on GSM8K, demonstrating a clear lead over the SFT baseline (75.1%). These results confirm that *Mix-GRM* provides reliable supervision, enabling policies to internalize both helpfulness and correctness.

Reward Model for Test-time Scaling. For test-time scaling, leveraging increased inference-time compute to enhance generalization, *Mix-GRM* functions as a robust verifier to re-rank candidates to identify the optimal solution via Best-of- N selection. Following the JETTS protocol (Zhou et al., 2025b), we evaluate $N = 10$ samples from a Llama-3.1-8B generator across 4 diverse benchmarks: MATH and CHAMP (math), as well as MBPP+ and BigCodeBench (coding). As shown in Figure 2, our method consistently secures the highest accuracy, setting a new SOTA for 8B-scale rerankers. The performance advantage is particularly pronounced in reasoning-heavy tasks; for instance, on MATH, our model achieves an accuracy of 43.2%, outperforming the RL-driven *RM-R1* (37.7%) and the data-intensive *FARE-8B* (35.2%). This confirms that ours provides a more discriminative signal for logical verification than methods relying on massive data scaling or generic RL.

5 Analysis

Switching CoT Mechanism Analysis. Visualizing structural transformations (Figure 3) reveals how our pipeline reshapes reasoning mechanisms. First, the polarization of Breadth and Depth baselines confirms the rigidity of static templates, which create capability blind spots by sacrificing either reasoning depth or semantic coverage. Second, the balanced profile of *Mix-GRM (SFT)* indicates successful internalization of different distinct mechanisms. Most pivotally, the global expansion during **RLVR** validates our hypothesis of mechanism polarization. By optimizing for verdict accuracy, the model spontaneously converges on domain-specific mechanism biases—amplifying D-CoT for correctness while reinforcing B-CoT for preference. This emergent specialization confirms that our proposed alignment is not a handcrafted heuristic, but an inherent structural necessity discovered by the model to maximize evaluation efficacy.

Scaling & Selection Analysis. To understand the mechanics of B-CoT, we decouple the impact of quantity (aggregation scale) from quality (principle selection) as shown in Figure 4.1)

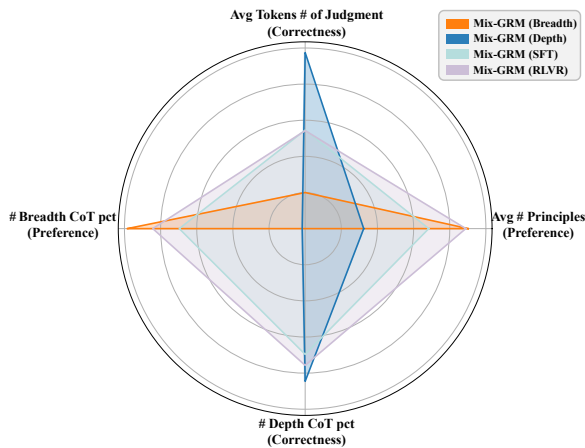


Figure 3: **Structural evolution of CoT mechanisms.** The chart tracks 4 indicators: the average token length per judgment, average principle count, and the percentage of CoT classified as having Breadth or Depth characteristics. Single-mode strategies show extreme trade-offs: **Mix-GRM (Breadth)** expands horizontally (high principle count), while **Mix-GRM (Depth)** extends vertically (long judgments). In contrast, **Mix-GRM (SFT)** achieves a robust union of both, which is further expanded into a broader reasoning manifold by **RLVR**.

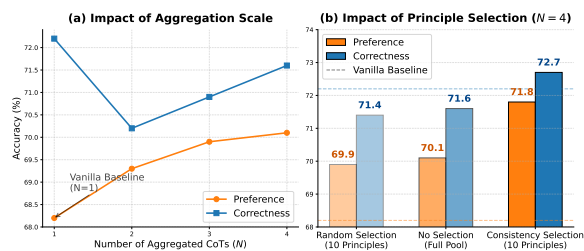


Figure 4: **Ablation of B-CoT synthesis.** (a) Aggregation Scale: Performance as aggregated rationales (N) increases from 1 (Vanilla) to 4. (b) Principle Selection: Comparison of Random, Full, and Consistency (Top-10) selection from the $N = 4$ pool. Orange/blue lines denote Preference/Correctness; dashed lines indicate the Vanilla baseline.

Quantity (Scaling): Figure 4(a) demonstrates that performance improves monotonically as the number of parallel CoTs (N) increases from 1 to 4. This confirms that “breadth” functions by expanding coverage; by aggregating diverse perspectives, the model minimizes the risk of overlooking critical error patterns.2) Quality (Selection): However, more is not always better. Figure 4(b) compares three strategies within the $N = 4$ pool: $Breadth_{Rand}$, $Breadth_{Full}$, and $Breadth_{Top10}$, where Top-10 means ten most frequent principles appearing across 4 CoTs. We observe a clear hierarchy: $Breadth_{Top10} > Breadth_{Full} > Breadth_{Rand}$. While the full pool improves over random sampling, it is the Top-10 consensus that achieves the highest

CASE 1: PREFERENCE DOMAIN (B-CoT WINS)	
<i>Inst (JP):</i>	アフガニスタン... (Is Afghanistan a puppet of Pakistan?)
<i>Resp A (EN):</i>	Detailed history... <i>Resp B (JP):</i> アフガニスタンがパキス...
✗ Vanilla	“A provides comprehensive history.” (Ignored language). <i>Verdict:</i> $[[A]]$.
✗ D-CoT	“Deep analysis of history...” (Tunnel Vision). ... <i>Verdict:</i> $[[A]]$.
✓ B-CoT	Multi-dim Scan: “1. Principle of Linguistic Alignment: “Assistant B’s response is in Japanese... Sub-Verdict: «B»” ... <i>Verdict:</i> $[[B]]$.
CASE 2: CORRECTNESS DOMAIN (D-CoT WINS)	
<i>Inst:</i>	On the basis of oxidation-reduction potential, which of the following is most likely to occur?
<i>Resp A:</i>	The order is Alkali > ... > Zn > ... > Ag, ... D. Zn + ... (Correct)
<i>Resp B:</i>	Alkali > Alkaline earth ... B. Mg+K... (Logic Error: Mg>K)
✗ Vanilla	Surface Length: “B analyzes more options and is longer.” <i>Verdict:</i> $[[B]]$.
✗ B-CoT	Superficial Heuristic: “B covers options A-H comprehensively.” (No verification). <i>Verdict:</i> $[[B]]$.
✓ D-CoT	Rigorous Analysis: “Check B: Claims Mg displaces K (False, K>Mg). Check D: Valid. Pick A.” <i>Verdict:</i> $[[A]]$.

Table 4: **Simplified Case Study.** **Case 1:** B-CoT catches language mismatch. **Case 2:** D-CoT verifies logical steps. The detailed Case is shown in Table. 5

gains (71.8/72.7). This suggests a denoising effect where low-frequency principles introduce noise, while high-frequency ones form a more robust “reasoning consensus.” Thus, representativeness—not just volume—is vital for robust breadth.

Case Study. Table 4 elucidates the structural drivers of the observed trade-off. In preference, B-CoT acts as a multi-dimensional scanner, identifying lateral mismatches (e.g., tone) that D-CoT misses due to attentional tunneling. Conversely, for correctness, D-CoT functions as a probe, exposing factual hallucinations (e.g., $K > Mg$) that B-CoT overlooks by mistaking superficial formatting for logical validity. This confirms that while Breadth ensures multi-faceted alignment, Depth remains the non-negotiable driver for rigorous evaluation.

6 Conclusion

This work demonstrates that beyond mere length scaling, the reliability of GRMs is fundamentally driven by the integration of different reasoning mechanisms. By introducing Mix-GRM, we prove that the frontier of reward modeling lies in synergizing two orthogonal reasoning mechanisms: B-CoT for multi-dimensional coverage and D-CoT for judgment soundness. Through mechanism-adaptive alignment, Mix-GRM ensures that the RM’s reasoning mechanism is precisely calibrated to the nature of the task. Ultimately, these findings shift the focus of GRM development from brute-force expansion to structural optimization.

549 Limitations

550 While Mix-GRM significantly enhances evaluation
551 reliability through mechanism alignment, we identify
552 two primary limitations that warrant further
553 investigation:

554 **Granularity of the Reasoning Manifold.** Our
555 framework successfully captures the double disso-
556 ciation between Subjective Preference and Objec-
557 tive Correctness, which we identify as the dominant
558 axes of the reasoning manifold. However, this di-
559 chotomy represents a coarse-grained mapping of
560 the diverse alignment landscape. Real-world tasks
561 often exist on a continuous spectrum or involve
562 hybrid demands that intricately blend deductive
563 rigor with multi-dimensional nuances. While we
564 prove that the model’s reasoning structure sponta-
565 neously converges toward these two primary poles,
566 our current categorization may act as a low-rank
567 approximation of a higher-dimensional space of
568 mechanisms. Future work could explore more gran-
569 ular taxonomies to achieve even more precise task-
570 mechanism calibration.

571 **Rigidity in Ambiguous Task Boundaries.** Our
572 analysis demonstrates that RLVR induces an intrin-
573 sic convergence toward specialized reasoning poles.
574 However, this emergent polarization may introduce
575 a degree of structural rigidity when encountering
576 hybrid tasks that do not fit neatly into the “Subjec-
577 tive vs. Objective” dichotomy. For instance, tasks
578 that require both factual precision and sophisticated
579 stylistic nuance may demand a dynamic fusion of
580 B-CoT and D-CoT. While our current framework
581 focuses on aligning specialized mechanisms with
582 their respective domains, the spontaneous sharp-
583 ening of reasoning styles might come at the cost
584 of generalist flexibility in highly nuanced, cross-
585 domain scenarios. Future research could explore
586 adaptive, soft-routing mechanisms that allow for a
587 more fluid transition across the reasoning manifold.

588 References

589 Zachary Ankner, Mansheej Paul, Brandon Cui,
590 Jonathan D. Chang, and Prithviraj Ammanabrolu.
591 2024. [Critique-out-loud reward models](#). *Preprint*,
592 arXiv:2408.11791.

593 Maciej Besta, Florim Memedi, Zhenyu Zhang, Robert
594 Gerstenberger, Guangyuan Piao, Nils Blach, Piotr
595 Nyczyk, Marcin Copik, Grzegorz Kwaśniewski, Jür-
596 gen Müller, Lukas Gianinazzi, Ales Kubicek, Hu-
597 bert Niewiadomski, Aidan O’Mahony, Onur Mutlu,

and Torsten Hoefler. 2025. Demystifying chains,
trees, and graphs of thoughts. *Transactions on
Pattern Analysis and Machine Intelligence*, page
10967–10989.

Nuo Chen, Zhiyuan Hu, Qingyun Zou, Jiaying Wu,
Qian Wang, Bryan Hooi, and Bingsheng He. 2025a.
[Judgelrm: Large reasoning models as a judge](#).
Preprint, arXiv:2504.00050.

Xiusi Chen, Gaotang Li, Ziqi Wang, Bowen Jin, Cheng
Qian, Yu Wang, Hongru Wang, Yu Zhang, Denghui
Zhang, Tong Zhang, Hanghang Tong, and Heng Ji.
2025b. [RM-R1: Reward modeling as reasoning](#).
Preprint, arXiv:2505.02387.

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.
2021. [Training verifiers to solve math word prob-
lems](#). *Preprint*, arXiv:2110.14168.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming
Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong
Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and
Maosong Sun. 2024. [UltraFeedback: Boosting lan-
guage models with scaled ai feedback](#). *Preprint*,
arXiv:2310.01377.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin,
Shengding Hu, Zhiyuan Liu, Maosong Sun, and
Bowen Zhou. 2023. Enhancing chat language models
by scaling high-quality instructional conversations.
In *Conference on Empirical Methods in Natural Lan-
guage Processing*.

Yann Dubois, Percy Liang, and Tatsunori Hashimoto.
2024. Length-controlled alpaca-eval: A simple de-
biasing of automatic evaluators. In *Conference on
Language Modeling*.

Evan Frick, Tianle Li, Connor Chen, Wei-Lin Chi-
ang, Anastasios Nikolas Angelopoulos, Jiantao Jiao,
Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica.
2025. How to evaluate reward models for RLHF. In
*The Thirteenth International Conference on Learning
Representations*.

Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar
Nath, Yunzhong He, Bing Liu, and Sean Hendryx.
2025. [Rubrics as rewards: Reinforcement
learning beyond verifiable domains](#). *Preprint*,
arXiv:2507.17746.

Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang,
Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and
Nouha Dziri. 2024. WILDGUARD: open one-stop
moderation tools for safety risks, jailbreaks, and re-
fusals of llms. In *International Conference on Neural
Information Processing Systems*.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021. Measuring mathematical

653	problem solving with the MATH dataset. In <i>Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Zhuang Li, Yuncheng Hua, Thuy-Trang Vu, Haolan Zhan, Lizhen Qu, and Gholamreza Haffari. 2025c. SCAR: Data selection via style consistency-aware response ranking for efficient instruction-tuning of large language models . <i>Preprint</i> , arXiv:2406.10882.	709
654			710
655			711
656	Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bham-bri, Lucas Saldyt, and Anil Murthy. 2024. Llms can't plan, but can help planning in llm-modulo frame-works . <i>Preprint</i> , arXiv:2402.01817.	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's verify step by step . <i>Preprint</i> , arXiv:2305.20050.	712
657			713
658			714
659			715
660			716
661	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing fine-grained evaluation capability in language models . In <i>International Conference on Learning Representations</i> .	Chengwu Liu, Ye Yuan, Yichun Yin, Yan Xu, Xin Xu, Zaoyu Chen, Yasheng Wang, Lifeng Shang, Qun Liu, and Ming Zhang. 2025a. Safe: Enhancing mathematical reasoning in large language models via retrospective step-aware formal verification . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	717
662			718
663			719
664			720
665			721
666			722
667			723
668	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. MAWPS: A math word problem repository . In <i>Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 1152–1157.	Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Ju-jie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024. Skywork-reward: Bag of tricks for reward modeling in llms . <i>Preprint</i> , arXiv:2410.18451.	724
669			725
670			726
671			727
672			728
673	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. RewardBench: Evaluating reward models for language modeling . <i>Preprint</i> , arXiv:2403.13787.	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and LINGMING ZHANG. 2023. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation . In <i>Conference on Neural Information Processing Systems</i> .	729
674			730
675			731
676			732
677			733
678			734
679	Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven C.H. Hoi. 2022. CodeRL: mastering code generation through pretrained models and deep reinforcement learning . In <i>International Conference on Neural Information Processing Systems</i> .	Tianci Liu, Ran Xu, Tony Yu, Ilgee Hong, Carl Yang, Tuo Zhao, and Haoyu Wang. 2025b. OpenRubrics: Towards scalable synthetic rubric generation for reward modeling and llm alignment . <i>Preprint</i> , arXiv:2510.07743.	735
680			736
681			737
682			738
683			739
684			740
685	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. 2024. Rlaif vs. rlhf: scaling reinforcement learning from human feedback with ai feedback . In <i>International Conference on Machine Learning</i> .	Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. 2025c. RM-bench: Benchmarking reward models of language models with subtlety and style . In <i>International Conference on Learning Representations</i> .	741
686			742
687			743
688			744
689			745
690			746
691			747
692	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. 2024. From live data to high-quality benchmarks: The arena-hard pipeline .	Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu. 2025d. Inference-time scaling for generalist reward modeling . <i>Preprint</i> , arXiv:2504.02495.	748
693			749
694			750
695			751
696	Yi-Chen Li, Tian Xu, Yang Yu, Xuqin Zhang, Xiong-Hui Chen, Zhongxiang Ling, Ningjing Chao, Lei Yuan, and Zhi-Hua Zhou. 2025a. Generalist reward models: Found inside large language models . <i>Preprint</i> , arXiv:2506.23235.	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning . <i>Preprint</i> , arXiv:2209.14610.	752
697			753
698			754
699			755
700			756
701	Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, Yingying Zhang, Fei Yin, Jiahua Dong, Zhiwei Li, Bao-Long Bi, Ling-Rui Mei, Junfeng Fang, Xiao Liang, Zhi-jiang Guo, and 2 others. 2025b. From system 1 to system 2: A survey of reasoning large language models . <i>Preprint</i> , arXiv:2502.17419.	Saumya Malik, Valentina Pyatkin, Sander Land, Jacob Morrison, Noah A. Smith, Hannaneh Hajishirzi, and Nathan Lambert. 2025. RewardBench 2: Advancing reward model evaluation . <i>Preprint</i> , arXiv:2506.01937.	757
702			758
703			759
704			760
705			761
706			762
707			763
708			764
		Yujun Mao, Yoon Kim, and Yilun Zhou. 2024. CHAMP: A competition-level dataset for fine-grained analyses of llms' mathematical reasoning capabilities . <i>Preprint</i> , arXiv:2401.06961.	764

765	Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang, Huazhong Yang, and Yu Wang. 2024. Skeleton-of-thought: Prompting LLMs for efficient parallel generation. In <i>International Conference on Learning Representations</i> .	821
766		822
767		
768		
769		
770	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. 2022. Training language models to follow instructions with human feedback. In <i>International Conference on Neural Information Processing Systems</i> .	
771		
772		
773		
774		
775		
776		
777		
778		
779	Jiayi Pan, Xiuyu Li, Long Lian, Charlie Snell, Yifei Zhou, Adam Yala, Trevor Darrell, Kurt Keutzer, and Alane Suhr. 2025. Learning adaptive parallel reasoning with language models . <i>Preprint</i> , arXiv:2504.15466.	
780		
781		
782		
783		
784	Junsoo Park, Seungyeon Jwa, Ren Meiyong, Daeyoung Kim, and Sanghyuk Choi. 2024. OffsetBias: Leveraging debiased data for tuning evaluators. In <i>Findings of the Association for Computational Linguistics: EMNLP</i> , pages 1043–1067.	
785		
786		
787		
788		
789	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model . In <i>Conference on Neural Information Processing Systems</i> .	
790		
791		
792		
793		
794	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models . <i>Preprint</i> , arXiv:2402.03300.	
795		
796		
797		
798		
799		
800	Wei Shen, Guanlin Liu, Yu Yue, Ruofei Zhu, Qingping Yang, Chao Xin, and Lin Yan. 2025. Exploring data scaling trends and effects in reinforcement learning from human feedback. In <i>Annual Conference on Neural Information Processing Systems</i> .	
801		
802		
803		
804		
805	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning . <i>Preprint</i> , arXiv:2303.11366.	
806		
807		
808		
809	Zayne Rea Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2025. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. In <i>International Conference on Learning Representations</i> .	
810		
811		
812		
813		
814		
815		
816	DeepSeek-AI Team. 2025a. DeepSeek-R1: Incentivizing reasoning capability in llms via reinforcement learning . <i>Preprint</i> , arXiv:2501.12948.	
817		
818		
819	Kimi Team. 2025b. Kimi k1.5: Scaling reinforcement learning with llms . <i>Preprint</i> , arXiv:2501.12599.	
820		
	Qwen Team. 2025c. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.	
	Vijay Viswanathan, Yanchao Sun, Shuang Ma, Xiang Kong, Meng Cao, Graham Neubig, and Tongshuang Wu. 2025. Checklists are better than reward models for aligning language models . <i>Preprint</i> , arXiv:2507.18624.	823
		824
		825
		826
		827
	Junqiao Wang, Zeng Zhang, Yangfan He, Zihao Zhang, Xinyuan Song, Yuyang Song, Tianyu Shi, Yuchen Li, Hengyuan Xu, Kunyu Wu, Xin Yi, Zhongwei Wan, Xinhang Yuan, Zijun Wang, Kuan Lu, Menghao Huo, Tang Jingqun, Guangwu Qian, Keqin Li, and 2 others. 2025a. Enhancing code llms with reinforcement learning in code generation: A survey . <i>Preprint</i> , arXiv:2412.20367.	828
		829
		830
		831
		832
		833
		834
		835
	Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhihang Sui. 2024a. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	836
		837
		838
		839
		840
		841
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . <i>Preprint</i> , arXiv:2203.11171.	842
		843
		844
		845
		846
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024b. Mmlu-pro: a more robust and challenging multi-task language understanding benchmark. In <i>International Conference on Neural Information Processing Systems</i> .	847
		848
		849
		850
		851
		852
		853
		854
	Zhilin Wang, Jiaqi Zeng, Olivier Delalleau, Hoo-Chang Shin, Felipe Soares, Alexander Bukharin, Ellie Evans, Yi Dong, and Oleksii Kuchaiev. 2025b. Helpsteer3-preference: Open human-annotated preference data across diverse tasks and languages . <i>Preprint</i> , arXiv:2505.11475.	855
		856
		857
		858
		859
		860
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903.	861
		862
		863
		864
		865
	Chenxi Whitehouse, Tianlu Wang, Ping Yu, Xian Li, Jason Weston, Iliia Kulikov, and Swarnadeep Saha. 2025. J1: Incentivizing thinking in llms-as-a-judge via reinforcement learning . <i>Preprint</i> , arXiv:2505.10320.	866
		867
		868
		869
		870
	Austin Xu, Xuan-Phi Nguyen, Yilun Zhou, Chien-Sheng Wu, Caiming Xiong, and Shafiq Joty. 2025. Foundational automatic evaluators: Scaling multi-task generative evaluator training for reasoning-centric domains . <i>Preprint</i> , arXiv:2510.17793.	871
		872
		873
		874
		875

876	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	Yilun Zhou, Austin Xu, Peifeng Wang, Caiming Xiong,	933
877	Thomas L. Griffiths, Yuan Cao, and Karthik	and Shafiq Joty. 2025b. Evaluating judges as evalu-	934
878	Narasimhan. 2023. Tree of thoughts: deliberate prob-	ators: The jetts benchmark of llm-as-judges as test-	935
879	lem solving with large language models. In <i>Internat-</i>	time scaling evaluators . <i>Preprint</i> , arXiv:2504.15253.	936
880	<i>ional Conference on Neural Information Processing</i>		
881	<i>Systems</i> .		
882	Edward Yeo, Yuxuan Tong, Morry Niu, Graham	Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu,	937
883	Neubig, and Xiang Yue. 2025. Demystifying	Wenhao Yu, Ratnadira Widayarsi, Imam Nur Bani	938
884	long chain-of-thought reasoning in llms . <i>Preprint</i> ,	Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon	939
885	arXiv:2502.03373.	Brunner, Chen GONG, James Hoang, Armel Randy	940
886		Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kad-	941
887	Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU,	dour, Ming Xu, Zhihan Zhang, and 14 others. 2025.	942
888	Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li,	Bigcodebench: Benchmarking code generation with	943
889	Adrian Weller, and Weiyang Liu. 2024. Metamath:	diverse function calls and complex instructions. In	944
890	Bootstrap your own mathematical questions for large	<i>International Conference on Learning Representa-</i>	945
891	language models. In <i>International Conference on</i>	<i>tions</i> .	946
892	<i>Learning Representations</i> .		
893	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho,		
894	Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason		
895	Weston. 2024. Self-rewarding language models. In		
896	<i>International Conference on Machine Learning</i> .		
897	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran		
898	Kazemi, Aviral Kumar, and Rishabh Agarwal. 2025a.		
899	Generative verifiers: Reward modeling as next-token		
900	prediction. In <i>International Conference on Learning</i>		
901	<i>Representations</i> .		
902	Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang,		
903	Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo,		
904	Yufei Wang, Niklas Muennighoff, Irwin King, Xue		
905	Liu, and Chen Ma. 2025b. A survey on test-time		
906	scaling in large language models: What, how, where,		
907	and how well? <i>Preprint</i> , arXiv:2503.24235.		
908	Qiyuan Zhang, Yufei Wang, Yuxin Jiang, Liangyou		
909	Li, Chuhan Wu, Yasheng Wang, Xin Jiang, Lifeng		
910	Shang, Ruiming Tang, Fuyuan Lyu, and Chen Ma.		
911	2025c. Crowd comparative reasoning: Unlocking		
912	comprehensive evaluations for LLM-as-a-judge. In		
913	<i>Annual Meeting of the Association for Computational</i>		
914	<i>Linguistics</i> .		
915	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan		
916	Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin,		
917	Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang,		
918	Joseph E. Gonzalez, and Ion Stoica. 2023. Judging		
919	LLM-as-a-judge with MT-bench and chatbot arena.		
920	In <i>Conference on Neural Information Processing Sys-</i>		
921	<i>tems Datasets and Benchmarks Track</i> .		
922	Tong Zheng, Hongming Zhang, Wenhao Yu, Xiaoyang		
923	Wang, Rungeng Dai, Rui Liu, Huiwen Bao, Cheng-		
924	song Huang, Heng Huang, and Dong Yu. 2025.		
925	Parallel-R1: Towards parallel thinking via reinforce-		
926	ment learning . <i>Preprint</i> , arXiv:2509.07980.		
927	Enyu Zhou, Guodong Zheng, Binghai Wang, Zhiheng		
928	Xi, Shihan Dou, Rong Bao, Wei Shen, Limao Xiong,		
929	Jessica Fan, Yurong Mou, Rui Zheng, Tao Gui,		
930	Qi Zhang, and Xuanjing Huang. 2025a. RMB: Com-		
931	prehensively benchmarking reward models in LLM		
932	alignment. In <i>International Conference on Learning</i>		
	<i>Representations</i> .		

Table 5: Case Study. **Case 1** shows how Breadth-CoT aggregates diverse principles to identify subtle preference nuances. **Case 2** shows how Depth-CoT performs step-by-step verification to catch logical errors.

CASE 1: PREFERENCE DOMAIN (BREADTH-CoT WINS)	
Instruction: アフガニスタンがパキスタンの傀儡というのは本当ですか？ (Is it true that Afghanistan is a puppet of Pakistan?)	
Response A (Rejected): [Language: English] "A sensitive topic! ...While Pakistan has historically exerted significant influence..."	Response B (Chosen): [Language: Japanese] 「アフガニスタンとパキスタンの関係については...簡単に『傀儡』と断言するのは適切ではなく...」
Reasoning Comparison	
✗ Vanilla-CoT	Assistant A offers a comprehensive breakdown of historical context... Verdict: [[A]] (Fail: Ignored language mismatch)
✓ Breadth-CoT	<ol style="list-style-type: none"> Principle of Linguistic Alignment: "Assistant B's response is in Japanese... Sub-Verdict: «B»" Principle of Contextual Nuance: "Assistant A provides a detailed explanation... Sub-Verdict: «A»" Principle of Cultural Sensitivity: "..." Final Verdict: [[B]] (Success)
✗ Depth-CoT	<ol style="list-style-type: none"> Principle of Comprehensive Contextualization: The evaluation needs deep reasoning, and let's solve the instruction: アフガニスタンが... Based on this, Assistant A excels... Sub-Verdict: «A» Principle of Balanced Presentation: ... Final Verdict: [[A]] (Fail: Tunnel vision on content depth)
CASE 2: CORRECTNESS DOMAIN (DEPTH-CoT WINS)	
Instruction: On the basis of oxidation-reduction potential, which of the following is most likely to occur? ...	
Response A (Chosen): "The reactivity of metals... The order is: Alkali > ... > Zn > ... > Ag. Based on this, the most likely reaction is: $D. Zn + 2Ag(NO_3)_2$..." (Correct)	Response B (Rejected): "The reactivity... Alkali > Alkaline earth... Analysis: $B. Mg + 2KNO_3 \rightarrow 2K + Mg(NO_3)_2$." (Error: $K > Mg$.)
Reasoning Comparison	
✗ Vanilla-CoT	To evaluate the responses, let's consider factors: 1. Helpfulness; 2. Relevance... Assistant B goes further by analyzing multiple options... Final Verdict: [[B]] (Fail: Fooled by length/detail)
✗ Breadth-CoT	<ol style="list-style-type: none"> Principle of Comprehensive Option Analysis: "Assistant B analyzes all provided options (A-H)... Sub-Verdict: «B»" Principle of Informative Detail: "Assistant B includes more elements... Sub-Verdict: «B»" Final Verdict: [[B]] (Fail: Superficial heuristic)
✓ Depth-CoT	<ol style="list-style-type: none"> Principle of Comprehensive Analysis: The evaluation needs to deep reasoning, and let's solve the instruction: To determine the most likely reaction... The correct order is Alkali > ... Given the options, the most likely reaction is Option D. ... Assistant B correctly identifies the importance but incorrectly identifies Option B... Sub-Verdict: «A» Principle of Direct Relevance: Assistant A directly addresses the question... Sub-Verdict: «A» Final Verdict: [[A]] (Success)

A Case Study

Table 5 elucidates the structural mechanisms behind the observed double dissociation. Case 1 demonstrates why Breadth-CoT dominates preference tasks: acting as a multi-dimensional scanner, it successfully penalizes a detailed but language-mismatched response by validating lateral constraints (e.g., *Linguistic Alignment*), whereas Depth-CoT exhibits attentional tunneling, fixating on verifying historical facts while missing the high-level language mismatch. Conversely, Case 2 reveals why *Depth-CoT* is essential for correctness: its step-by-step derivation acts as a logic probe, allowing it to spot subtle factual hallucinations (e.g., $K > Mg$) hidden within a lengthy explanation. Here, *Breadth-CoT* actively fails due to feature interference, where it mistakes superficial comprehensiveness (length and formatting) for logical validity. This confirms that while Breadth is necessary for satisfying diverse user preferences, Depth is the non-negotiable driver for rigorous verification.

B Training Implementation

B.1 Hyperparameters Setting

We provide the detailed hyperparameter settings in the Table 6 and Table 7.

Hyperparameters	Values
Epochs	2
Learning rate	$2e-5$
Batch Size	128 (gradient accumulation steps = 16)
Seq Length	12, 288
Weight Decay	0.
Warmup	5% linear warmup

Table 6: Hyperparameter settings for SFT.

Hyperparameters	Values
Training Steps	100
Learning Rate	$1e-6$
Batch Size	128
KL Loss Coefficient	0.001
KL Coefficient	0.001
Rollouts	$n = 8$ using vLLM with temperature 0.8

Table 7: Hyperparameter settings for RL.

B.2 Training Data Source Details

To cultivate general rewarding capabilities, it is essential to curate a training corpus that encompasses diversified real-world scenarios. We construct our dataset by performing stratified random sampling from representative data sources, ensuring balanced coverage across distinct alignment domains, including general chat, STEM, coding, math, safety, multilingual, and instruction following. The specific source datasets, their corresponding domains, and the sampling statistics are detailed in Table 8.

Source Dataset	Domain	Samples
HeIpSteer-3 (Single-Turn)	General Chat	4,973
	STEM	2,321
	Code	4,322
	Multilingual	3,260
Code-Preference	Code	4,000
Math-DPO	Math	4,000
WildGuard	Safety	4,000
OffsetBias	Instruction Following	4,000
Total	–	30,876

Table 8: Composition and statistics of the training data sampled from domain-specific sources.

B.3 Training Data Synthesis Details

To synthesize the CoT data for SFT, we utilized DeepSeek-v3 (0324 snapshot) as the backbone generator. The generation process was configured with a sampling temperature of $T = 0.8$ to promote diversity in the trajectories while maintaining logical coherence. Notably, we abstain from consistency filtering: Contrary to common practices that discard samples where the synthesized verdict diverges from the ground-truth human label, our empirical verification reveals that training on the full synthesized CoTs yields superior performance compared to aggressive filtering, regardless of verdict consistency.

B.4 Training Offline Reinforcement Learning Details

To strictly control for temporal data leakage and ensure a fair comparison with the release dates of our evaluation benchmarks, we select **Llama-3-8B** as our base foundation model. The offline reinforcement learning pipeline consists of two phases: SFT initialization and DPO.

Policy Initialization (SFT). We first derive a supervised policy model by fine-tuning Llama-3-8B

Table 9: Task coverage of the evaluated general reward benchmarks.

Benchmark	Tasks	Samples
REWARDBENCH	Chat, Math, Code, Safety	2,985
REWARDBENCH-v2	Focus, IF, Factuality, Math, Safety, Ties	1,865
RM-BENCH	Chat, Math, Code, Safety	11,943
RMB	Harmfulness, Helpfulness (General, Code)	14,725
PPE (Exclude Tie)	Chat, MMLU-Pro, GPQA, IFEval, MBPP	22,991

on a composite dataset. This dataset ensures basic instruction-following and reasoning capabilities, consisting of the **UltraChat** dataset (Ding et al., 2023) and a random subset of 40K samples from **MetaMathQA** (Yu et al., 2024). We train the model for 2 epochs using a learning rate of $2e-5$ and a maximum sequence length of 2,048 tokens. This SFT model serves as the initial policy π_{ref} for the subsequent DPO stage.

DPO Data Construction via RM Labeling. To evaluate the practical utility of different RMs, we employ them to annotate preferences on a unified source dataset. The prompt source comprises 10K instructions randomly sampled from **UltraFeedback** (Cui et al., 2024) and 40K instructions from **MetaMathQA**. For each instruction x , we generate $N = 5$ diverse candidate responses using gpt-4o-mini with a temperature of 0.8.

We adopt a **Pairwise Scoring Aggregation** strategy to construct the final preference pairs (x, y_w, y_l) . Specifically, for the set of 5 responses, we generate all possible combinations of pairs ($\binom{5}{2} = 10$ pairs). The target RM evaluates each pair, assigning +1 point to the preferred response (chosen) and 0 to the non-preferred one (rejected). After traversing all pairs, we calculate the cumulative score for each response. The response with the highest total score is selected as the positive sample (y_w), and the response with the lowest total score is selected as the negative sample (y_l). These labeled pairs are then used to train the policy via DPO.

C Evaluation Implementation

C.1 Core Benchmarks

There is a list of benchmarks and corresponding task coverage.

C.2 Benchmarks for Offline Reinforcement Learning Evaluation

To comprehensively assess the policy derived from DPO, we conduct evaluations across two distinct

1048	domains: mathematical reasoning and open-ended	derived from at least two Vanilla-CoT responses,	1096
1049	instruction following.	followed by a deduplication step to ensure diverse	1097
1050	Mathematical Reasoning. We employ a suite of	In contrast, the synthesis of Depth-CoT	1098
1051	four challenging datasets to evaluate the model’s	relies on a reasoning-guided evaluation mechanism.	1099
1052	deductive logic and problem-solving capabilities:	We initially prompt the model to reason the instruc-	1100
1053	GSM8k (Cobbe et al., 2021), MATH (Hendrycks	tion deeply. We then use this generated reasoning	1101
1054	et al., 2021), MAWPS (Koncel-Kedziorski et al.,	to ground the re-assessment of selected principles	1102
1055	2016), and TabMWP (Lu et al., 2023). These	extracted from the parsed schemas, discarding their	1103
1056	benchmarks cover a wide spectrum of diffi-	previous rationales to ensure the new judgments	1104
1057	culty, ranging from grade-school arithmetic to	are purely driven by rigorous reasoning.	1105
1058	competition-level mathematics and tabular process-		
1059	ing.	E Reward Model Performance Across	1106
		Preference and Correctness	1107
1060	Instruction Following. For general alignment	To provide a more granular view of our model’s	1108
1061	and conversational versatility, we utilize two widely	efficacy, we report detailed performance across spe-	1109
1062	adopted benchmarks: AlpacaEval-2 (Dubois et al.,	cific tasks based on the meta-data provided by each	1110
1063	2024) and Arena-Hard v0.1 (Li et al., 2024). Eval-	benchmark. We categorize these tasks into two	1111
1064	uation is performed using an auto-evaluator in a	distinct tables: Table 10 for subjective preference	1112
1065	head-to-head setting, where the model’s responses	tasks and Table 11 for objective correctness tasks.	1113
1066	are compared against a baseline reference to deter-	This fine-grained reporting serves as a detailed de-	1114
1067	mine win rates. We strictly adhere to the officially	composition of the mechanism-level performance	1115
1068	recommended configurations for reproducibility.	discussed in the main text, offering deeper empir-	1116
		ical evidence for the mechanism-task synergy be-	1117
1069	C.3 Benchmarks for Test-time Scaling	tween B-CoT and D-CoT.	1118
1070	Evaluation		
1071	Following the JETTS setup (Zhou et al., 2025b),		
1072	we perform Best-of-10 reranking evaluations where		
1073	the model selects the optimal solution from a mixed		
1074	pool of candidate responses. We report results on		
1075	the four most challenging subsets of the bench-		
1076	mark: MATH (Hendrycks et al., 2021) for math-		
1077	ematical reasoning, CHAMP (Mao et al., 2024)		
1078	for competition-level math, along with MBPP+		
1079	(Liu et al., 2023) and BigCodeBench (Zhuo et al.,		
1080	2025) for code generation. This selection tests the		
1081	model’s ability to identify correct reasoning paths		
1082	in complex scenarios.		
1083	D Prompts Template		
1084	To align with established community standards, our		
1085	Vanilla-CoT generation employs the representative		
1086	prompts originally introduced in MT-Bench (Zheng		
1087	et al., 2023) and RewardBench (Lambert et al.,		
1088	2024). Upon generating the raw Vanilla-CoT us-		
1089	ing the standard prompts, we employ a specialized		
1090	extraction prompt to parse the unstructured text		
1091	into the modular “Principle–Judgment–Verdict”		
1092	schema. Leveraging the parsed schemas, we in-		
1093	troduce specialized prompts to synthesize the two		
1094	target morphologies. For Breadth-CoT, the synthe-		
1095	sis process entails merging modular components		

Prompt for Vanilla-CoT Generation

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user’s instructions and answers the user’s question better. Your evaluation should consider as many factors as possible. Begin your evaluation by comparing the two responses and provide a thorough reasoning. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your reasoning, output your final verdict by strictly following this format: `[[A]]`if assistant A is better, `[[B]]`if assistant B is better.

[Instruction]

instruction

[The Start of Assistant A’s Answer]

{response_a}

[The End of Assistant A’s Answer]

[The Start of Assistant B’s Answer]

{response_b}

[The End of Assistant B’s Answer]

Table 10: Performance of RMs on preference-related sub-tasks. “Avg.” is the average score among sub-tasks. Best per column is **bolded**; second-best is underlined.

Models	Reward Bench		Reward Bench-v2		RM Bench	RMB		PPE		Avg.
	CHAT	FOCUS	IF	CHAT	HELPFULNESS	HUMAN	IF			
<i>Open-sourced Reward Models</i>										
JudgeLRM-7B	75.8	46.8	31.3	68.4	<u>79.3</u>	60.2	54.3	59.4		
RM-R1-7B (Distill)	75.3	58.7	24.4	58.7	63.2	57.1	53.0	55.8		
RM-R1-7B (Instruct)	80.5	85.3	28.8	64.7	65.1	57.1	53.0	62.1		
FARE-8B	85.0	78.4	<u>36.3</u>	66.9	82.9	63.4	55.7	67.0		
RubricRM-8B	82.4	78.2	33.8	62.2	77.5	63.8	66.0	66.3		
DeepSeek-GRM-16B	80.6	79.2	40.0	64.0	76.8	61.7	57.9	65.7		
<i>Our Proposed Reward Models</i>										
SFT-trained										
Base-GRM	81.6	76.6	34.4	63.3	80.5	64.3	55.9	65.2		
Mix-GRM (Breadth)	83.7	84.5	33.8	65.9	77.9	63.5	55.6	66.4		
Mix-GRM (Depth)	80.3	72.2	28.1	70.6	70.1	63.1	54.1	62.6		
Mix-GRM	84.9	79.6	31.9	71.2	78.7	62.0	56.3	66.4		
RLVR-trained										
Base-GRM	83.0	86.7	29.4	68.5	73.8	65.8	57.1	66.3		
Mix-GRM (Breadth)	86.2	88.9	28.8	70.1	79.2	66.0	55.3	<u>67.8</u>		
Mix-GRM (Depth)	<u>85.3</u>	<u>89.3</u>	26.3	75.6	75.4	66.0	56.4	<u>67.8</u>		
Mix-GRM	86.2	91.3	37.5	<u>72.7</u>	78.1	<u>65.9</u>	<u>57.4</u>	69.9		

Prompt for Schema Extraction

PRIMARY TASK:

Your mission is to analyze a given reasoning Chain-of-Thought from a generative reward model. From this CoT, you will extract, define, and refine the specific, detailed principles (or rubrics, criteria) it uses to judge the quality of AI-generated responses. For each principle, you must provide a corresponding analysis that traces it directly back to the original text.

INSTRUCTIONS:

You will be given a CoT text below. Please follow these four steps precisely:

1. Deconstruct the CoT: First, perform a close reading of the entire CoT. Identify all explicit evaluation criteria mentioned as well as any implicit judgments or preferences revealed in the model's comparative language.

2. Extract the Core Idea of Each Criterion: For each criterion, do not simply use the high-level category name. Your goal is to uncover the specific description of that criterion as used by the model. Ask yourself: What specific actions, qualities, or content does the model praise or criticize? What makes one response "more accurate" or "clearer" according to this specific CoT?

3. Formulate and Refine the Principle: Convert each core idea you extracted into a formal, normative, and reusable principle.

3.1 Name It: Give the principle a clear and descriptive name that captures its essence (e.g., "Principle of Factual Precision," "Principle of Structural Clarity").

3.2 Define It: Write the principle as a concise, actionable, and universal rule. It should be an instructive statement about what constitutes a high-quality response.

3.3 Be Specific: Avoid vague terms. Instead of "The response should be relevant," specify how it should be relevant based on the CoT's logic, such as "A relevant response must directly and unambiguously address the user's primary question."

3.4 Be Normative: Phrase it as a standard to be met (e.g., "A high-quality response must...").

4. Provide Corresponding Judgment: For each principle you formulate, you must write a brief "CoT Judgment Extraction." To do this, quote or closely paraphrase specific phrases from the CoT that support your formulation.

5. Conclude the sub-verdict in this Judgment: For the principle and corresponding judgment, you should conclude this verdict in this segment.

OUTPUT FORMAT:

You must follow this format strictly for your entire response.

```
. . .
### 1. Principle of [Descriptive Name]: [Your refined, normative principle statement.] Judgment:
[In this principle, what judgment on Response A and Response B quotes or paraphrases from the
source CoT.] Sub-Verdict: «A/B», In this principle, the judgment judge which assistant Better]
### 2. Principle of [Descriptive Name]: [Your refined, normative principle statement.] Judgment:
[In this principle, what judgment on Response A and Response B quotes or paraphrases from the
source CoT.]*** Sub-Verdict: «A/B», In this principle, the judgment judge which assistant Better]
(Continue this structure for all principles identified in the CoT)
```

```
. . .
Extract and Analyse the following CoT Text
{Vanilla-CoT}
```

Prompt for Breadth-CoT Generation

PRIMARY TASK:

You are provided with a series of lists, each containing Principles, Judgments, and Sub-Verdicts derived from an independent analysis of a Chain-of-Thought (CoT). Your mission is to merge these lists into a single, master list of unique evaluation principles.

INSTRUCTIONS FOR MERGING AND SYNTHESIS:

1. Deduplication and Semantic Grouping:

* Compare all Principles with the corresponding Judgments across all provided lists. * Identify and group principles that are **semantically similar**, even if they use different wording (e.g., “Principle of Precision” and “Principle of Correctness” are likely the same concept).

2. Principle Refinement:

* For each semantic group, synthesize the most concise, actionable, and specifically-detailed statement for the **Principle Description**. * Select the most descriptive and formal **Name** for the refined principle.

3. Judgment Synthesis:

* For the refined principle, create a new, synthesized **Judgment** block. This block should consist of a curated selection of the most illustrative quotes and paraphrases from the original Judgments across all source lists that led to the consolidated principle. This new Judgment serves as the combined evidence for the principle.

4. Merge Count:

* **COUNT THE SOURCES:** For each synthesized principle, you must count the total number of original, distinct principles/judgments from the source lists that were merged to create it. This number is the **Merge Count**.

5. Sub-Verdict Aggregation:

* The final **Sub-Verdict** for the synthesized principle must reflect the aggregated trend. Since the judgments are now synthesized, simply use the majority verdict (e.g., if a principle appeared 4 times with [[B]] and 1 time with [[A]], conclude [[B]]). If the verdicts are balanced (e.g., 2 [[A]] and 2 [[B]]), state [[MIXED]].

6. Strict Output Adherence:

* Maintain the exact four-part format for every final entry. The output must be one continuous list of unique, synthesized principles.

SOURCE LISTS TO MERGE:

[Insert List 1 Here]

[Insert List 2 Here]

[Insert List 3 Here]

(Continue for all lists)

OUTPUT FORMAT:

You must follow this exact format for your final, merged response.

1. Principle of [Refined, Descriptive Name]: [The synthesized, normative principle statement.] **Judgment:** [A synthesis of the most relevant quotes/paraphrases from the source Judgments that supports this consolidated principle.]*** **Merge Count:** [The total number of original source principles/judgments that were merged to form this entry.] **Sub-Verdict:** «A/B/MIXED», The aggregate verdict for this principle across all CoTs.]

2. Principle of [Refined, Descriptive Name]: [The synthesized, normative principle statement.] **Judgment:** [A synthesis of the most relevant quotes/paraphrases from the source Judgments that supports this consolidated principle.]*** **Merge Count:** [The total number of original source principles/judgments that were merged to form this entry.] **Sub-Verdict:** «A/B/MIXED», The aggregate verdict for this principle across all CoTs.]
(Continue this structure for all unique, synthesized principles)

Q Prompt for Depth-CoT Verification

PRIMARY TASK: Your role is to critically assess the quality of two competing responses (Assistant A and Assistant B) against the user’s question, leveraging the expert reasoning as the ultimate ground truth.

MANDATORY NON-BIAS RULES: Avoid all position biases (do not favor the first response presented). Do not allow the length or formatting of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective, clinical, and data-driven as possible.

PRINCIPLE-BASED EVALUATION

Given some potential principles, {principles}, you should choose the most critical principle (Preferably one principle, with a maximum of three) from them and then evaluate the two Chatbot responses (A and B) based on the choosed principle. This evaluation must directly reference the deep reasoning to instruction and ****must strictly adhere to the following output format for each principle:****

EXPERT REASONING: {reasoning}

Principle of [Critical Principle Name]:

Judgment: [Give your specific and detailed evaluation in this principle, and if you are referring the this reasoning, you ****MUST**** quote using ‘<Answer>’]

Sub-Verdict: «A/B/MIXED», «A» if assistant A is better in this principle, «B» if assistant B is better in this principle, «MIXED» if assistant A and B is Tie.

After providing your complete principle-based evaluation, output your final verdict by strictly following this format: [[A]]if assistant A is better, [[B]]if assistant B is better.

Table 11: Performance of RMs on correctness-related sub-tasks. “Avg.” is the average within this block. Best per column is **bolded**; second-best is underlined.

Models	RewardBench		RewardBench-v2		RM-Bench		RMB	PPE				Avg.
	CODE	MATH	FACTUALITY	MATH	CODE	MATH	CODE	MMLU-Pro	MATH	GPQA	MBPP	
<i>Open-sourced Reward Models</i>												
JudgeLRM-7B	81.6	77.2	53.8	76.5	51.0	86.7	82.1	57.2	65.5	51.3	52.3	66.8
RM-R1-7B (Distill)	91.9	93.7	28.3	73.2	53.3	<u>85.8</u>	74.8	66.7	89.4	56.3	<u>64.4</u>	70.7
RM-R1-7B (Instruct)	81.7	84.1	42.6	67.8	56.7	72.7	74.7	67.0	<u>89.1</u>	<u>55.9</u>	64.8	68.8
FARE-8B	88.1	82.3	65.8	68.9	57.0	69.1	88.1	63.2	79.3	55.2	55.4	70.2
RubricRM-8B	93.6	81.7	50.8	77.6	55.4	59.8	86.5	60.9	75.5	52.8	52.4	67.9
DeepSeek-GRM-16B	84.0	69.1	49.4	62.3	51.5	61.7	86.8	55.2	64.3	54.1	53.7	62.9
<i>Our Proposed Reward Models</i>												
SFT-trained												
Base-GRM	91.0	77.2	46.0	81.4	57.1	78.4	86.4	59.9	71.8	52.2	52.5	68.5
Mix-GRM (Breadth)	90.1	72.0	51.1	69.4	54.0	74.1	86.8	60.0	70.3	51.9	52.4	66.6
Mix-GRM (Depth)	89.8	86.1	45.1	75.4	56.2	77.1	81.1	<u>66.8</u>	83.6	54.7	53.7	70.0
Mix-GRM	88.9	87.9	55.7	76.0	55.8	79.6	81.9	63.9	82.1	54.8	54.0	71.0
RLVR-trained												
Base-GRM	93.2	86.4	<u>62.0</u>	77.0	61.7	78.1	89.5	64.6	84.1	54.3	50.5	72.9
Mix-GRM (Breadth)	96.3	69.3	55.7	71.0	58.0	70.6	86.5	61.7	75.2	53.0	52.8	68.2
Mix-GRM (Depth)	94.6	<u>89.0</u>	61.8	<u>78.7</u>	<u>64.4</u>	81.4	87.4	67.0	86.5	55.6	55.5	<u>74.7</u>
Mix-GRM	<u>95.4</u>	<u>89.0</u>	65.8	79.2	66.6	82.5	<u>88.9</u>	65.0	86.7	54.8	55.2	75.4