

INFERENCE-TIME SCALING FOR GENERALIST REWARD MODELING

Anonymous authors
 Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) has been widely adopted in post-training for large language models (LLMs) at scale. Recently, the incentivization of reasoning capabilities in LLMs from RL indicates that proper learning methods could enable effective inference-time scalability. A key challenge of RL is to obtain accurate reward signals for LLMs in various domains beyond verifiable questions or artificial rules. In this work, we investigate how to improve reward modeling (RM) with more inference compute for general queries, i.e. the **inference-time scalability of generalist RM**. For the RM approach, we adopt pointwise generative reward modeling (GRM) to enable flexibility for different input types and the potential for inference-time scaling. For the learning method, we propose **Self-Principled Critique Tuning (SPCT)** to foster scalable reward generation behaviors in GRMs through online RL, to generate principles adaptively and critiques accurately, resulting in **SPCT-GRM** models. Furthermore, for effective inference-time scaling, we use parallel sampling to expand compute usage, and introduce a meta RM to guide the voting process for better scaling performance. Empirically, we show that SPCT significantly improves the quality and scalability of GRMs, outperforming existing methods and models in various RM benchmarks without severe biases, and it can achieve better performance compared to training-time scaling. SPCT-GRM still meets challenges in some tasks, which we believe can be addressed by future efforts in generalist reward systems. The models will be released.

1 INTRODUCTION

The remarkable advancements in large language models (LLMs) (DeepSeek-AI, 2024b; OpenAI, 2025b) have catalyzed significant shifts in artificial intelligence research, enabling models to perform tasks that require understanding, generation, and nuanced decision-making capabilities. Recently, reinforcement learning (RL) as a post-training method for LLMs has been widely adopted at scale, and resulting in remarkable improvements in human value alignment (Ouyang et al., 2022; Bai et al., 2022a), long-term reasoning (DeepSeek-AI, 2025; OpenAI, 2025c), and environment adaptation (OpenAI, 2025a) for LLMs. Reward modeling (RM) (Gao et al., 2023), as a crucial component in RL, is essential for generating accurate reward signals for LLM responses. Current studies (Lightman et al., 2024; DeepSeek-AI, 2025) also show that, with high-quality and robust rewards in either training or inference time, LLMs can achieve strong performance in specific domains.

However, such high-quality rewards in specific domains are mainly obtained from human-designed environments with clear conditions (Yao et al., 2022; Xie et al., 2024) or from hand-crafted rules for verifiable questions, e.g., mathematical problems (Hendrycks et al., 2021; Veeraboina, 2023) and

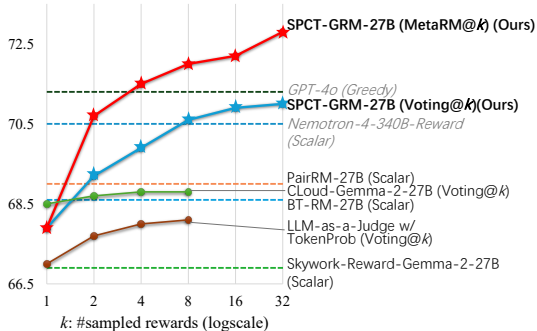


Figure 1: Inference-time scaling performance with different RMs on all tested RM benchmarks. Results are shown with up to 8 samples for each method, and are further scaled to 32 samples for ours. Non-italic font indicates models based on Gemma-2-27B.

coding tasks (Jimenez et al., 2024; Zhuo et al., 2025). In general domains, reward generation is more challenging, as the criteria for rewards are more diverse and complex, and there are often no explicit reference or ground truth. Generalist reward modeling is thus crucial for improving the performance of LLMs in broader applications, either from post-training perspectives, e.g., RL at scale, or from inference perspectives, e.g., RM-guided search. Furthermore, RM performance should be improved by increasing both the training compute (Gao et al., 2023) and the inference compute.

In practice, challenges arise in making RMs both general and effectively scalable in inference time. The former demands (1) flexibility for different input types and (2) accurate reward generation in various domains. We refer to this paradigm as **generalist reward modeling**. Moreover, effective **inference-time scalability** requires the RM (3) to generate higher-quality reward signals with increased inference compute, and (4) to learn scalable behaviors for better performance-compute scaling. Existing research on reward modeling demonstrates several paradigms for reward generation, including scalar (Cobbe et al., 2021; Wang et al., 2024e; Liu et al., 2024), semi-scalar (Ye et al., 2025a; Yu et al., 2025b; Zhang et al., 2025a), and generative (Li et al., 2024a; Kim et al., 2024; Vu et al., 2024; Cao et al., 2024; Arabzadeh et al., 2024; Ye et al., 2025b; Alexandru et al., 2025; Yu et al., 2025a) approaches, and various scoring patterns, such as pointwise (Kendall & Smith, 1940; Gao et al., 2023; Yuan et al., 2024; Winata et al., 2025; Guo et al., 2025) and pairwise (Park et al., 2024; Zheng et al., 2023; Jiang et al., 2023; Wang et al., 2024c; Liu et al., 2025). These approaches inherently determine the input flexibility and the inference-time scalability of RMs ((1)&(3)), as shown in Figure 2. For instance, pairwise RMs only consider the relative preference of paired responses, lacking flexibility to accept single or multiple responses as input; scalar RMs could hardly generate diverse reward signals for the same response, which obstructs getting better rewards through sampling-based inference-time scaling methods (Snell et al., 2025). Also, different learning methods (Wang et al., 2024a; Ankner et al., 2024; Wang et al., 2024c; Mahan et al., 2024) have been proposed to improve the quality of rewards, but few of them focus on inference-time scalability and study the interconnection between the learned reward generation behaviors and the effectiveness of inference-time scaling of RMs, resulting in marginal performance improvement ((2)&(4)). Current research (DeepSeek-AI, 2025) indicates that effective inference-time scalability could be enabled by proper learning methods, which raises the question: *Can we design a learning method aiming to enable effective inference-time scaling for generalist reward modeling?*

In this work, we investigated different approaches for RM, and found that pointwise generative reward modeling (GRM) could unify the scoring of single, paired, and multiple responses within pure language representation, overcoming challenge (1). We explored that certain principles could guide reward generation within proper criteria for GRMs, improving the quality of rewards, which suggested that *inference-time scalability of RM might be achieved by scaling the generation of high-quality principles and accurate critiques*. Based on this preliminary, we propose a novel learning method, **Self-Principled Critique Tuning (SPCT)**, to foster effective inference-time scalable behaviors in GRMs. By leveraging rule-based online RL, SPCT enables GRMs to learn to adaptively posit principles and critiques based on the input query and responses, leading to better outcome rewards in general domains (challenge (2)). We then come up with **SPCT-GRM-27B**, which is post-trained with SPCT based on Gemma-2-27B (Team, 2024). For inference-time scaling, we expand compute usage by sampling multiple times. By sampling in parallel, SPCT-GRM could generate different sets of principles and corresponding critiques, and then vote for the final reward. **With larger-scale sampling, SPCT-GRM could judge more accurately based on more diverse principles, and output rewards with finer granularity**, which resolves challenge (3)&(4). Furthermore, we train a meta RM in addition to voting for better scaling performance. Empirically, we show that SPCT significantly improves the quality and scalability of GRMs, outperforming existing methods and models in multiple comprehensive RM benchmarks without severe domain biases. We also compared the inference-time scaling performance of SPCT-GRM-27B with larger models up to 671B parameters, and found it could achieve better performance compared to training-time scaling on model sizes. Though the current method meets challenges in efficiency and specific tasks, with efforts beyond SPCT, we believe GRMs with enhanced scalability and efficiency could serve as a versatile interface for generalist reward systems, advancing the frontiers of LLM post-training and inference.

In general, our main contributions are as follows.

1. We propose a novel approach, **Self-Principled Critique Tuning (SPCT)**, to foster effective inference-time scalability for generalist reward modeling, resulting in **SPCT-GRM** models.

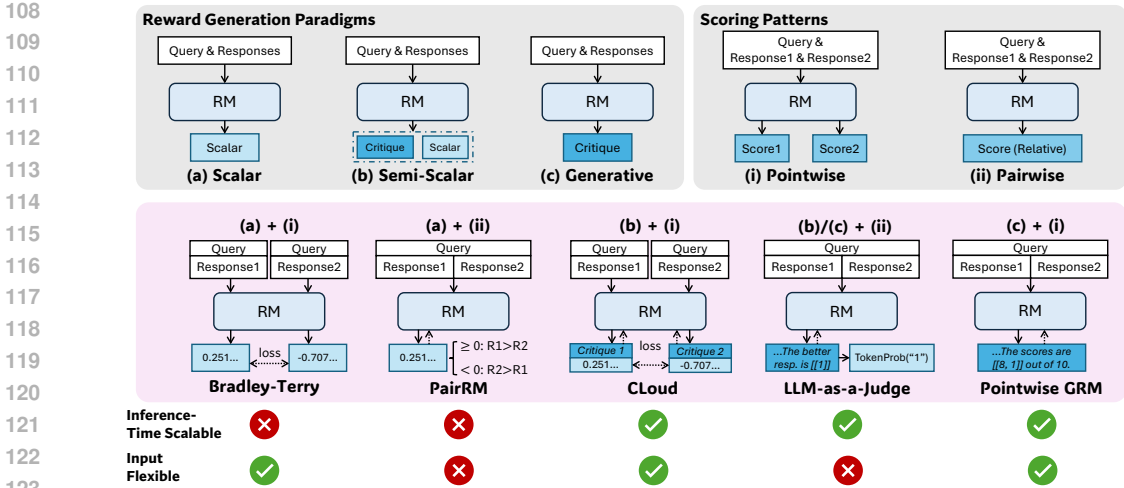


Figure 2: Different paradigms for reward generation, including (a) scalar, (b) semi-scalar, and (c) generative approaches, and different scoring patterns, including (i) pointwise and (ii) pairwise approaches. We list the representative methods for each approach, and corresponding inference-time scalability (whether better rewards could be obtained from multiple sampling) and input flexibility (whether supports rating single and multiple responses).

And we further introduce a meta RM to effectively improve the inference-time scaling performance of SPCT-GRM beyond voting.

2. We empirically show SPCT significantly improves the quality and inference-time scalability of GRMs over existing methods and several strong public models.
3. We also applied the SPCT training schedule to LLMs with larger sizes and found that inference-time scaling could outperform model size scaling in training time.

2 PRELIMINARIES

2.1 COMPARISONS OF DIFFERENT RM APPROACHES

As shown in Figure 2, RM approaches are mainly determined by reward generation paradigms and scoring patterns, which inherently affect the inference-time scalability and the input flexibility of the RM. For **reward generation paradigms**, we distinguish three main approaches: scalar, semi-scalar, and generative. The scalar approach assigns scalar values to the given query and responses, while the semi-scalar approach generates textual judgement, termed “critique”, and the scalar reward value as well. The generative approach only generates critiques as the textual reward, from which the reward value could be extracted. For **scoring patterns**, we distinguish two main approaches: pointwise and pairwise. The pointwise approach assigns an individual score to each response, while the pairwise approach selects a single best response from all candidates.

To expand inference compute usage, we focus on sampling-based methods, which generate multiple sets of rewards for the same query and responses, and then aggregate the final reward. Thus, the *inference-time scalability* of RMs is determined by whether different rewards could be obtained from multiple sampling, where scalar RMs would fail in most cases due to the invariant generation of rewards; and the *input flexibility* is defined by whether the RM supports rating single, paired, and multiple responses, where pairwise RMs could hardly rate single responses and usually require extra techniques (Jiang et al., 2023; Liu et al., 2025) to handle multiple responses. The formulation of pointwise GRMs is:

$$\{S_i\}_{i=1}^n = f_{\text{point}}(\mathcal{R}, \{y_i\}_{i=1}^n) = f_{\text{extract}}(\mathcal{C}), \quad \mathcal{R} = \mathcal{C} \sim r_{\theta}(x, \{y_i\}_{i=1}^n), \quad (1)$$

where x is the query, y_i is the i -th response, r_{θ} is the reward function parameterized by θ , \mathcal{R} is the generated reward, \mathcal{C} is the textual judgement or critique, $S_i \in \mathbb{R}$ is the individual score of y_i , f_{point} denotes the pointwise scoring function, and $f_{\text{extract}}(\cdot)$ is the implementation of f_{point} to extract the rewards from textual generation results. Usually, the scores are discrete, and in this work, we assign $S_i \in \mathbb{N}, 1 \leq S_i \leq 10$ by default. Detailed analysis is in Appendix D.1.

2.2 BOOSTING REWARD QUALITY WITH PRINCIPLES

Generalist RM requires generating high-quality rewards beyond specific domains (Hendrycks et al., 2021; Jimenez et al., 2024), where the criteria for rewards are more diverse and complex, and there are often no explicit reference or ground truth. To this end, for general domains, we adopt principles to guide reward generation in place of artificial rules. Principles for LLMs are first introduced in Constitutional AI (Bai et al., 2022b; Sharma et al., 2025), which are hand-crafted criteria that guide the LLMs or curated classifiers to construct safe data pipelines. With principles, the reward generation of GRMs changes to

$$\mathcal{R} = \mathcal{C} \sim r_{\theta}(x, \{y_i\}_{i=1}^n, \{p_i\}_{i=1}^m), \quad (2)$$

where $\{p_i\}_{i=1}^m$ denotes the principles. We conduct a preliminary experiment to examine the influence of proper principles on reward quality, with the Chat Hard subset of Reward Bench (Lambert et al., 2024) and the IFEval subset of the PPE benchmark (Frick et al., 2025).

In the experiment, data samples contain a query and two responses, with the ground-truth label denoting the better response. We use GPT-4o-2024-08-06 to generate the principles and then pointwise rewards four times for each sample. We filter the principles from the correct reward generation process, where the larger reward value is assigned to the labeled better response. We test different LLMs with principles generated by themselves and the filtered principles, and compare them with the default setting with no principle guidance. The results are shown in Table 1. We found that the self-generated principles barely improve performance, but the filtered principles could significantly boost the reward quality. The result is non-trivial and two main conclusions could be drawn: (a) Current LLMs could generate diverse principles, but not all of them are proper for reward generation. (b) A subset of generated principles could better guide reward generation under correct criteria, indicating a potential of self-bootstrapping. These findings are the foundation of the usage of online RL to optimize GRMs, where they could learn from the principles generated by themselves with a clear signal to tell whether the principles are proper or not. Other details are depicted in Appendix E.

Table 1: Preliminary experiments on the influence of principles on reward quality. The default setting of SPCT-GRM-27B includes self-generated principles.

Method	Chat Hard	IFEval
GPT-4o-2024-08-06	76.1	56.0
w/ Self-Gen. Principles	75.9	55.6
w/ Filtered Principles	77.8	57.5
Gemma-2-27B-it	59.1	56.1
w/ Self-Gen. Principles	64.0	55.8
w/ Filtered Principles	68.0	57.3

3 SELF-PRINCIPLED CRITIQUE TUNING (SPCT)

Inspired by the preliminary results, we developed a novel approach for pointwise GRMs to learn generating adaptive and high-quality principles that could effectively guide the generation of critiques, termed **Self-Principled Critique Tuning** (SPCT). As shown in Figure 3, SPCT consists of two phases: rejective fine-tuning, as the cold start, and rule-based online RL, reinforcing generalist reward generation by advancing the generated principles and critiques. SPCT fosters these behaviors in GRMs for inference-time scaling as well.

3.1 UNPINNING PRINCIPLES FROM UNDERSTANDING TO GENERATION

From preliminary experiments in Section 2.2, we found that proper principles could guide reward generation within certain criteria, which is critical for high-quality rewards. However, it remains challenging to generate effective principles for generalist RM at scale. To address this challenge, we propose to unpin principles from understanding to generation, i.e. view principles as a part of reward generation instead of a preprocessing step.

Formally, principles guide the generation of rewards following Equation 2, when principles are pre-defined. GRMs could generate principles themselves, and then generate critiques based on the principles, formalized as

$$\{p_i\}_{i=1}^m \sim p_{\theta}(x, \{y_i\}_{i=1}^n), \quad \mathcal{R} = \mathcal{C} \sim r_{\theta}(x, \{y_i\}_{i=1}^n, \{p_i\}_{i=1}^m), \quad (3)$$

where p_{θ} is the principle generation function parameterized by θ , that shares the same model with reward generation r_{θ} . In practice, they are implemented with the same language head. This shift enables principles to be generated based on the input query and responses, adaptively aligning reward generation process. The quality and granularity of the principles and corresponding critiques could be improved with post-training on GRMs. With **the principles generated at scale**, GRMs could rate rewards with finer granularity and broader criteria to achieve better inference-time scalability.

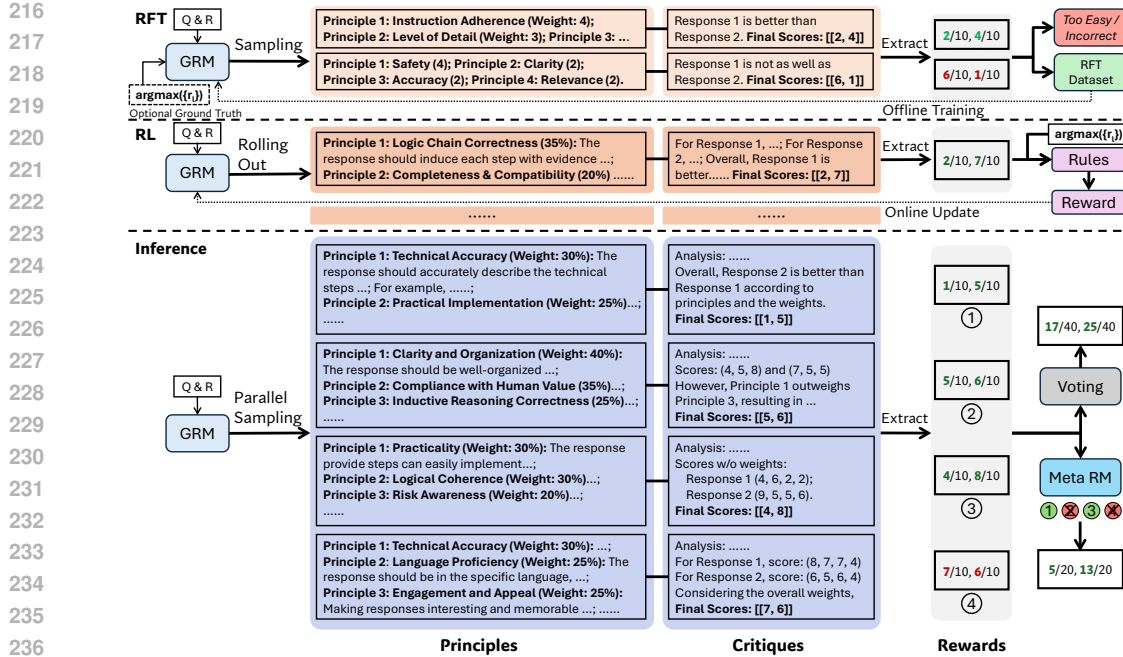


Figure 3: Illustration of SPCT, including rejective fine-tuning, rule-based RL, and scalable behaviors during inference. The inference-time scaling is achieved via voting or meta RM guided voting with principles generated at scale, resulting in finer-grained outcome rewards within a larger value space.

3.2 RULE-BASED REINFORCEMENT LEARNING

To optimize principle and critique generation in GRMs simultaneously, we propose SPCT, which integrates rejective fine-tuning and rule-based RL. The former serves as a cold start.

Rejective Fine-Tuning (Cold Start) The core idea of the rejective fine-tuning stage is to train the GRM to generate principles and critiques in the correct format and for various input types. Unlike previous works (Vu et al., 2024; Cao et al., 2024; Alexandru et al., 2025) that mix RM data for single, paired, and multiple responses in different formats, we adopt pointwise GRM, introduced in Section 2.1, to flexibly generate rewards for any number of responses in the same format. For data construction, besides general instruction data, we sample trajectories with a pretrained GRM by giving the query and corresponding responses. Each RM data point contains a query and one or multiple responses to the query, as well as the ground-truth label denoting the best response. For each RM data point, the sampling of principles and critiques is performed N_{RFT} times. The rejection strategy is also unified, which is to reject trajectories with predicted rewards that are incorrect, and the query and responses with all N_{RFT} trajectories correct (too easy). Formally, let r_i denote the ground-truth reward for the i -th response y_i to the query x , the predicted pointwise rewards $\{S_i\}_{i=1}^n$ are correct if

$$\begin{cases} \forall i \neq j, & S_j > S_i, \quad j = \arg \max_l \{r_l\}_{l=1}^n, & \text{if } n \geq 2, \\ S_1 = r_1, & & \text{if } n = 1. \end{cases} \quad (4)$$

with a guarantee that the ground-truth rewards only contain one maximum. However, similar to previous works (Zhang et al., 2025a), we found pretrained GRMs could hardly generate correct rewards for a portion of data points within limited sampling quota. Thus, we optionally append $\arg \max_l \{r_l\}_{l=1}^n$ to the prompt of the GRM, termed *hinted sampling*, with the expectation that the predicted rewards to align with the ground truth, besides *non-hinted sampling*. Specifically, an additional segment “The best response is: Response $\arg \max_l \{r_l\}_{l=1}^n$ ” will be appended to the input. For hinted sampling, each query and the corresponding responses are sampled once, and trajectories are only rejected when incorrect. Beyond previous studies (Li et al., 2024a; Mahan et al., 2024), we observed that hinted sampled trajectories sometimes take shortcuts in judgement, especially for reasoning tasks, indicating the necessity and potential benefits of non-hinted sampled data.

Rule-Based RL The GRM is further fine-tuned with rule-based online RL. Specifically, we use the original setting of GRPO (Shao et al., 2024) with rule-based outcome rewards. During roll-

out, the GRM generates principles and critiques based on the input query and responses, and then the predicted reward is extracted and compared to the ground truth with accuracy rules. Unlike DeepSeek-AI (2025), no format rewards are used. Instead, a larger coefficient for the KL penalty is applied to ensure the format and avoid severe biases. Formally, the reward for the i -th output o_i to the given query x and responses $\{y_i\}_{i=1}^n$ is

$$\hat{r}_i = \begin{cases} 1, & \text{if } n \geq 2 \text{ and } \forall i' \neq j', S_{j'} > S_{i'}, j' = \arg \max_l \{r_l\}_{l=1}^n, \\ 1, & \text{if } n = 1 \text{ and } S_1 = r_1, \\ -1, & \text{otherwise,} \end{cases} \quad (5)$$

where the pointwise rewards $\{S_i\}_{i=1}^n$ are extracted from o_i . **The reward function encourages GRMs to distinguish the best responses with online optimized principles and critiques, in favor of effective inference-time scaling.** The reward signal could be obtained seamlessly from any preference dataset and labeled LLM responses.

4 INFERENCE-TIME SCALING WITH SPCT

To further improve the performance of SPCT-GRM using more inference compute, we explore sampling-based strategies to achieve effective inference-time scalability.

Voting with Generated Rewards Voting is a widely adopted method for inference-time scaling in RM. Recalling the approaches in Section 2.1, we demonstrate voting results of k samples for semi-scalar and generative RMs. For semi-scalar RMs (Ankner et al., 2024; Zhang et al., 2025a), voting is performed as averaging:

$$S^* = \frac{1}{k} \sum_{i=1}^k S_i, \quad \{\mathcal{R}_i = (S_i, C_i)\}_{i=1}^k \sim r_\theta(x, \{y_i\}_{i=1}^n), \quad (6)$$

where S^* is the final reward. In practice, the scalar value has limited variance which could hinder the scalability. For pairwise GRMs (Mahan et al., 2024; Wang et al., 2024c), voting is performed as selecting the response identified to be the best with the highest frequency, i.e. majority:

$$\hat{y}^* = \arg \max_y \sum_{i=1}^k \mathbb{I}(y = \hat{y}_i), \quad \{\mathcal{R}_i = C_i\}_{i=1}^k \sim r_\theta(x, \{y_i\}_{i=1}^n), \quad (7)$$

where \hat{y}^* is the predicted best response, $f_{\text{pair}}(\cdot, \cdot)$ is a selection function, $\hat{y}_i = f_{\text{pair}}(C_i, \{y_i\}_{i=1}^n)$ is the individually selected best response of each sample, and $\mathbb{I}(\cdot)$ is the indicator function. Though the voting process is scalable, the majority voted result might be biased since ties are not allowed in each sample, and may not be able to tell apart subtle differences between responses due to the lack of quantitative scores. The voting process for pointwise GRMs is defined as summing the rewards:

$$S_i^* = \sum_{j=1}^k S_{i,j}, \quad \{p_{i,j}\}_{i=1}^{m_j} \sim p_\theta(x, \{y_i\}_{i=1}^n), \mathcal{R}_j = C_j \sim r_\theta(x, \{y_i\}_{i=1}^n, \{p_{i,j}\}_{i=1}^{m_j}), j = 1, \dots, k, \quad (8)$$

where S_i^* is the final reward for the i -th response ($i = 1, \dots, n$) and $\{S_{i,j}\}_{i=1}^n = f_{\text{point}}(C_j, \{y_i\}_{i=1}^n)$ is the j -th set of pointwise rewards. **Since $S_{i,j}$ is usually set within a small discrete range, e.g., $\{1, \dots, 10\}$, the voting process actually expands the reward space by k times, and enables the GRM to generate a large number of principles, which benefits the quality and granularity of the final rewards.** An intuitive explanation is that, if each principle could be viewed as a proxy of judgement perspectives, a larger number of principles may reflect the real distribution more accurately, leading to scaling effectiveness. Notably, to avoid positional biases and for diversity, responses are shuffled before sampling.

Meta Reward Modeling Guided Voting The voting process of SPCT-GRM requires multiple sampling and a few generated principles and critiques might be biased or low-quality due to randomness or model limitations. Thus, we train a meta RM to guide the voting process. The meta RM is a pointwise scalar RM, trained to identify the correctness of the principles and critiques generated by SPCT-GRM, with the binary cross-entropy loss, where the label is identified based on Equation 4. The prompt template is in Appendix H, integrates the query, candidate responses, corresponding principles, and the critiques. The dataset comprises trajectories from non-hinted sampling in the RFT stage, and trajectories sampled from the SPCT-GRM, to provide positive and negative rewards and alleviate the gap between training and inference policy (Chow et al., 2025). The guided voting is simple: The meta RM outputs meta rewards for k sampled rewards, and the final outcome is voted by rewards with top $k_{\text{meta}} \leq k$ meta rewards, so that filtering out low-quality samples.

Table 2: Overall results of different methods and models on RM benchmarks. Underlined numbers indicate the best performance, **bold numbers** indicate the best performance among baseline and our methods, and *italicized font* denotes scalar or semi-scalar RMs. For meta RM guided voting (MetaRM), $k_{\text{meta}} = \frac{1}{2}k$. “Avg. Rank (\downarrow)” is calculated based on the rank (from 1 to 6, better ranks denoted by darker shades) of baseline and our methods on each dataset.

Model	RB	PPE Pref.	PPE Correct.	RMB	Overall	Avg. Rank (\downarrow)
<i>Reported Results of Public Models</i>						
<i>Skywork-Reward-Gemma-2-27B</i>	<u>94.1</u>	56.6	56.6	60.2	66.9	-
DeepSeek-V2.5-0905	81.5	62.8	58.5	65.7	67.1	-
Gemini-1.5-Pro	86.8	66.1	59.8	56.5	67.3	-
<i>ArmoRM-8B-v0.1</i>	90.4	60.6	61.2	64.6	69.2	-
<i>InternLM2-20B-Reward</i>	90.2	61.0	63.0	62.9	69.3	-
LLaMA-3.1-70b-Instruct	84.1	65.3	59.2	68.9	69.4	-
Claude-3.5-sonnet	84.2	65.3	58.8	70.6	69.7	-
<i>Nemotron-4-340B-Reward</i>	92.0	59.3	60.8	69.9	70.5	-
GPT-4o	86.7	67.1	57.6	73.8	71.3	-
<i>Reproduced Results of Baseline Methods</i>						
LLM-as-a-Judge	83.4	64.2	58.8	64.8	67.8	4.50
BT-RM-27B	81.7	68.3	66.7	57.9	68.6	3.50
CLOUD-Gemma-2-27B	82.0	67.1	62.4	63.4	68.7	3.50
PairRM-27B	87.1	65.8	64.8	58.2	69.0	2.75
<i>Results of Our Method</i>						
SPCT-GRM-27B-RFT (Ours)	84.5	64.1	59.6	67.0	68.8	4.00
SPCT-GRM-27B (Ours)	86.0	64.7	59.8	69.0	69.9	2.75
<i>Results of Inference-Time Scaling (Voting@32)</i>						
SPCT-GRM-27B (Ours)	88.5	65.3	60.4	69.7	71.0	-
SPCT-GRM-27B (MetaRM) (Ours)	90.4	67.2	63.2	70.3	72.8	-

5 RESULTS ON REWARD MODELING BENCHMARKS

5.1 EXPERIMENT SETTINGS

Benchmarks and Evaluation Metrics We evaluate the performance of different methods on various RM benchmarks of different domains: **Reward Bench (RB)** (Lambert et al., 2024), **PPE** (Preference and Correctness subsets) (Frick et al., 2025), **RMB** (Zhou et al., 2025), **RealMistake** (Kamoi et al., 2024). We use the standard evaluation metrics for each benchmark: accuracy of picking the best response from a set of responses in Reward Bench, PPE, and RMB, and ROC-AUC for RealMistake. To deal with ties of the predicted rewards for multiple responses, we shuffle the responses and determine the best response by $\arg \max_i S_i$, where S_i is the predicted reward for the i -th response after shuffling. Details are in Appendix E.

Method Implementation For the baseline methods, we re-implement **LLM-as-a-Judge** (Zheng et al., 2023), **BT-RM-27B** (Bradley-Terry model) (Kendall & Smith, 1940), **CLOUD-Gemma-2-27B** (Ankner et al., 2024), and **PairRM-27B** (Jiang et al., 2023) based on Gemma-2-27B (Team, 2024) and with all compatible training data and settings as SPCT-GRM. For our methods, we implement **SPCT-GRM-27B-RFT** based on Gemma-2-27B, and **SPCT-GRM** on different sizes of LLMs, including DeepSeek-V2-Lite (16B MoE) (DeepSeek-AI, 2024a), Gemma-2-27B, DeepSeek-V2.5 (236B MoE), and DeepSeek-V3 (671B MoE) (DeepSeek-AI, 2024b). The meta RM is trained on Gemma-2-27B. Default results are reported with greedy decoding, and the inference-time scaling uses temperature = 0.5. Other details are provided in Appendix D.

5.2 RESULTS AND ANALYSIS

Performance on RM Benchmarks The overall results of different methods and models on RM benchmarks are shown in Table 2. We compare the performance of SPCT-GRM-27B with the reported results of public models and the reproduced results of baseline methods. We find that SPCT-GRM-27B outperforms the baseline methods in overall performance, and achieves competitive performance with strong public RMs, such as Nemotron-4-340B-Reward and GPT-4o; with inference-time scaling, SPCT-GRM-27B could further improve and achieve the best overall results. For detailed comparisons, scalar (BT-RM-27B) and semi-scalar (CLOUD-Gemma-2-27B) RMs demonstrate biased results on different benchmarks, with significantly better performance on verifiable tasks (PPE Correctness) than all generative RMs, but fail in different other benchmarks,

Table 3: Inference-time scalability results of different methods on RM benchmarks. Settings are the same as Table 2.

Model	Overall
<i>Reported Results of Public Models</i>	
Nemotron-4-340B-Reward	70.5
GPT-4o	71.3
<i>Results of Inference-Time Scaling (Voting@1)</i>	
LLM-as-a-Judge	67.0
C <small>LOUD</small> -G <small>EMMA</small> -2-27B	68.5
SPCT-GRM-27B-RFT (Ours)	67.8
SPCT-GRM-27B (Ours)	67.9
<i>Results of Inference-Time Scaling (Voting@8)</i>	
LLM-as-a-Judge	67.6 (+0.6)
LLM-as-a-Judge w/ TokenProb	68.1 (+1.1)
C <small>LOUD</small> -G <small>EMMA</small> -2-27B	68.8 (+0.3)
SPCT-GRM-27B-RFT (Ours)	69.3 (+1.5)
SPCT-GRM-27B (Ours)	70.6 (+2.7)
SPCT-GRM-27B (MetaRM) (Ours)	72.0 (+4.1)
<i>Results of Further Inference-Time Scaling (Voting@32)</i>	
SPCT-GRM-27B (Ours)	71.0 (+3.1)
SPCT-GRM-27B (MetaRM) (Ours)	72.8 (+4.9)

Table 4: Ablation studies for different components of the proposed SPCT. **Bold numbers** indicate the best performance.

Method	Overall
<i>Results of Greedy Decoding</i>	
SPCT-GRM-27B	69.9
w/o Principle Generation	67.5
w/o Rejective Sampling	68.7
SPCT-GRM-27B-RFT	68.8
w/o Hinted Sampling (①)	68.0
w/o Non-Hinted Sampling (②)	67.4
w/o Rejective Sampling (①&②)	66.1
w/o General Instruction Data	63.3
<i>Results of Inference-Time Scaling (Voting@8)</i>	
SPCT-GRM-27B	70.6
w/o Principle Generation	68.0
<i>Results of Inference-Time Scaling (Voting@32)</i>	
SPCT-GRM-27B	71.0
SPCT-GRM-27B ($k_{\text{meta}} = 1$)	71.5
SPCT-GRM-27B ($k_{\text{meta}} = 8$)	72.7
SPCT-GRM-27B ($k_{\text{meta}} = 16$)	72.8

respectively. Nonetheless, most public scalar RMs also exhibit severe domain biases. The PairRM approach could alleviate the problem. LLM-as-a-Judge shows similar trends with SPCT-GRM-27B with lower performance, potentially due to the lack of training on rating single responses. In conclusion, **SPCT improves the generalist reward generation capability of GRMs, with significantly fewer biases compared to scalar and semi-scalar RMs.**

Inference-Time Scalability The inference-time scaling results of different methods are shown in Table 3, and the overall trends are demonstrated in Figure 1. Details are in Appendix E.3. With up to 8 samples, we find that SPCT-GRM-27B has the highest performance increase over the greedy decoding and sampling results. SPCT-GRM-27B further shows a strong potential to increase performance with larger inference compute, up to 32 samples. We attribute the effectiveness to the refined principle generation, which expands the output length in a structural way and guides the outcome rewards closer to the ground-truth distribution. The meta RM also reveals its validity in filtering low-quality trajectories for SPCT-GRM on each benchmark. Voted with token probabilities, LLM-as-a-Judge also shows a significant performance increase, indicating that the token probability as quantitative weights could help the reliability over mere majority voting on discrete indices. For CLOUD-GEMMA-2-27B, the performance increase is limited, mainly due to the lack of variance in scalar reward generation, even though the critique has changed a lot. In summary, **SPCT improves the inference-time scalability of GRMs, and the meta RM further boosts the scaling performance in general.**

Ablation Study Table 4 shows the ablation study results of different components of the proposed SPCT, detailed results are listed in Appendix E.3. Surprisingly, without the cold start with rejective sampled critique data, general-instruction-tuned GRMs still improve significantly after undergoing the online RL (66.1 \rightarrow 68.7). Also, the non-hinted sampling seems more important than the hinted sampling, potentially because of the shortcuts appearing in hinted sampled trajectories. These indicate the **importance of online training for GRMs**. Aligned with previous works (Cao et al., 2024), we confirm that the general instruction data is essential for the performance of GRMs. We find that **the principle generation is crucial for the performance of both greedy decoding and inference-time scaling of SPCT-GRM-27B**. For inference-time scaling, the meta RM guided voting shows robustness with different k_{meta} . Further analysis on the generalist RM performance, including input flexibility, domain generalization of training data, etc., is discussed in Appendix F.

Scaling Inference and Training Costs We further investigate the inference-time and training-time scaling performance of SPCT-GRM-27B, by post-training LLMs of different sizes. The models are tested on the Reward Bench, and the results are shown in Figure 4. We find that direct voting with 32 samples of SPCT-GRM-27B could achieve comparable performance to the 671B MoE model, and the meta RM guided voting could achieve the best results with 8 samples, demonstrating the **effectiveness of inference-time scaling of SPCT-GRM-27B compared to scaling model sizes**. Moreover, we test DeepSeek-R1-0120 with a downsampled test set containing 300 samples, and find that

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

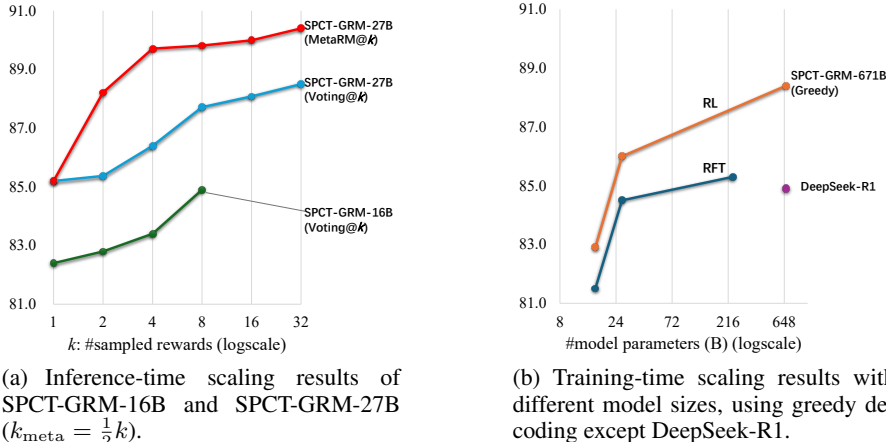


Figure 4: Inference-time scaling v.s. training-time scaling performance on Reward Bench.

its performance is even worse than the 236B MoE RFT model, indicating that expanding long chain-of-thoughts for reasoning tasks could not significantly improve the performance of generalist RM.

6 RELATED WORK

Generative Reward Models GRMs represent a paradigm shift from scalar RMs (Ouyang et al., 2022), modeling reward as textual feedback or scores. (Li et al., 2024a; Kim et al., 2024; Wang et al., 2024c; Cao et al., 2024; Vu et al., 2024; Alexandru et al., 2025), enabling richer reward representations and more flexible for judging single and multiple responses. Previously, the LLM-as-a-judge method (Zheng et al., 2023; Wang et al., 2025a) accommodates reference-based or reference-free pairwise judgement for evaluating LLMs. Recent studies use offline and online RL to train GRMs (Wu et al., 2024; Mahan et al., 2024; Yu et al., 2025a; Ye et al., 2025b; Chen et al., 2025), incorporate tools and external knowledge with GRMs (Li et al., 2024b; Peng et al., 2025), and even train GRMs as an interface to adjust rewards from environments (Baker et al., 2025). Very recently, GRMs are also used in productive systems for RL (OpenAI, 2025d). Though these methods face challenges in efficiency, they demonstrate the potential to improve rewards at scale.

Inference-Time Scaling for LLMs Inference-time scaling for LLMs has been a critical research direction parallel with scaling LLMs in training time. Studies focus on sampling and RM guided aggregation (Lightman et al., 2024; Brown et al., 2024; Snell et al., 2025; Wu et al., 2025). Recently, long-horizon chain-of-thoughts (Wei et al., 2022) incentivized from LLMs significantly improve the reasoning capabilities of the models in both solving (OpenAI, 2024; DeepSeek-AI, 2025; OpenAI, 2025c) and judging (Zheng et al., 2025; Kim et al., 2025) difficult verifiable questions, as another format of inference-time scaling. However, we do not find an effective way to incentivize long-horizon reward generation for generalist reward modeling with general RM data as DeepSeek-AI (2025), and we leave the combination of reasoning and principle-guided reward generation in a single model for future engineering efforts. There is also research using scalable rewards or verifiers to improve the performance of policy models, in domains of coding (Chen et al., 2023), reasoning (Lifshitz et al., 2025), etc. Thus, the development of inference-time scalable generalist RMs in this work might also contribute to the general performance of policy models by inference-time co-scaling.

7 CONCLUSION AND FUTURE WORK

We introduced Self-Principled Critique Tuning (SPCT), a method that enhances the scalability of inference time for generalist reward modeling. With rule-based online RL, SPCT enables adaptive generation of principles and critiques, significantly boosting reward quality and inference-time scalability for GRMs in diverse domains. Empirical results demonstrate that SPCT-GRM surpasses baseline methods and a few strong public RMs, and shows notable improvement through inference-time scaling, particularly with the guidance of the meta RM. Future directions could include integrating GRMs into online RL pipelines as versatile interfaces of reward systems, exploring inference-time co-scaling with policy models, or serving as robust offline evaluators for foundation models.

ETHICS STATEMENT

Our proposed method, Self-Principled Critique Tuning (SPCT), aims to enhance inference-time scalability of generative reward models (GRMs) for general domains. While this advancement promotes accuracy and consistency in reward modeling, several ethical implications might warrant explicit consideration.

Firstly, even though through our empirical analysis that SPCT-GRM shows fewer biases on different domains, the automated generation of principles and critiques can inadvertently perpetuate or amplify biases when the training data is toxic. We argue that further investigation in the meta RM and other bias mitigation strategies should be prioritized to ensure equitable outcomes. Also, our approach does not aim to diminish human oversight. Instead, we advocate maintaining human-in-the-loop frameworks, and developing reliable proxy methods, like SPCT, to scale human oversight more efficiently and effectively.

Secondly, expanded applicability of the inference-time scalable GRMs across diverse domains might raise concerns regarding transparency, accountability, etc. Since the reward generation behavior is largely emerged from self-bootstrapping, the potential for unfaithful principles and critiques is non-negligible. We demonstrate case studies in Appendix G.1 and limitations in Appendix C, and open-source the model under public supervision, which is essential for maintaining trust and ensuring responsible deployment of the artifact.

Finally, robust validation and ongoing vigilance across varied RM benchmarks and practical scenarios remain crucial. Ethical use of SPCT-GRM necessitates proactive management of risks and continuous evaluation against biases, requiring efforts in research about RM evaluation.

REPRODUCIBILITY STATEMENT

We have made several efforts to ensure the reproducibility of our work. First of all, we will open-source the models for research purposes and public use. Due to the size of the models, they could not be directly uploaded to the supplementary materials. Although we could not release the in-house training data due to privacy concerns, we have particularly described the open-sourced data sources and statistics in Section 5.1 and Appendix D. Following the training recipe (Section 3), one is able to reproduce the training of SPCT-GRM with open-source training infrastructure, e.g., verl (Sheng et al., 2025), and open-source LLMs, e.g., DeepSeek-V3-Base (DeepSeek-AI, 2024b). With replicated models, one should see the similar trends of the experimental results, including (1) SPCT-GRM has fewer domain biases in reward generation and outperforms other RM approaches; (2) SPCT-GRM shows significant improvement with inference-time scaling, especially with the meta RM; (3) The improvement from inference-time scaling is comparable or better to scaling model sizes in SPCT. Further improvements over the work could be and have been made in the community.

REFERENCES

- Andrei Alexandru, Antonia Calvi, Henry Broomfield, Jackson Golden, Kyle Dai, Mathias Leys, Maurice Burger, Max Bartolo, Roman Engeler, Sashank Pisupati, Toby Drane, and Young Sun Park. Atla selene mini: A general purpose evaluation model. *Computing Research Repository*, arXiv:2501.17195, 2025. URL <https://arxiv.org/abs/2501.17195>.
- Wei An, Xiao Bi, Guanting Chen, Shanhuang Chen, Chengqi Deng, Honghui Ding, Kai Dong, Qishi Du, Wenjun Gao, Kang Guan, Jianzhong Guo, Yongqiang Guo, Zhe Fu, Ying He, Panpan Huang, Jiashi Li, Wenfeng Liang, Xiaodong Liu, Xin Liu, Yiyuan Liu, Yuxuan Liu, Shanghao Lu, Xuan Lu, Xiaotao Nie, Tian Pei, Junjie Qiu, Hui Qu, Zehui Ren, Zhangli Sha, Xuecheng Su, Xiaowen Sun, Yixuan Tan, Minghui Tang, Shiyu Wang, Yaohui Wang, Yongji Wang, Ziwei Xie, Yiliang Xiong, Yanhong Xu, Shengfeng Ye, Shuiping Yu, Yukun Zha, Liyue Zhang, Haowei Zhang, Mingchuan Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, and Yuheng Zou. Fire-flyer ai-hpc: A cost-effective software-hardware co-design for deep learning. *Computing Research Repository*, arXiv:2408.14158, 2024. URL <https://arxiv.org/abs/2408.14158>.

- 540 Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D. Chang, and Prithviraj Ammanabrolu.
541 Critique-out-loud reward models. *Computing Research Repository*, arXiv:2408.11791, 2024.
542 URL <https://arxiv.org/abs/2408.11791>.
543
- 544 Negar Arabzadeh, Siqing Huo, Nikhil Mehta, Qingyun Wu, Chi Wang, Ahmed Hassan Awadal-
545 lah, Charles L. A. Clarke, and Julia Kiseleva. Assessing and verifying task utility in LLM-
546 powered applications. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Pro-
547 ceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp.
548 21868–21888, Miami, Florida, USA, November 2024. Association for Computational Linguis-
549 tics. doi: 10.18653/v1/2024.emnlp-main.1219. URL [https://aclanthology.org/
550 2024.emnlp-main.1219/](https://aclanthology.org/2024.emnlp-main.1219/).
- 551 Amanda Askeell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy
552 Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny
553 Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown,
554 Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as
555 a laboratory for alignment. *Computing Research Repository*, arXiv:2112.00861, 2021. URL
556 <https://arxiv.org/abs/2112.00861>.
- 557 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askeell, Anna Chen, Nova DasSarma, Dawn
558 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jack-
559 son Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Her-
560 nandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine
561 Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann,
562 and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from
563 human feedback. *Computing Research Repository*, arXiv:2204.05862, 2022a. URL [https://arxiv.org/abs/
564 2204.05862](https://arxiv.org/abs/2204.05862).
- 565 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askeell, Jackson Kernion, Andy Jones,
566 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Ol-
567 sson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-
568 Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse,
569 Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mer-
570 cado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna
571 Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Con-
572 erly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario
573 Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai:
574 Harmlessness from ai feedback. *Computing Research Repository*, arXiv:2212.08073, 2022b.
575 URL <https://arxiv.org/abs/2212.08073>.
- 576 Bowen Baker, Joost Huizinga, Leo Gao, Zehao Dou, Melody Y. Guan, Aleksander Madry, Wojciech
577 Zaremba, Jakub Pachocki, and David Farhi. Monitoring reasoning models for misbehavior and
578 the risks of promoting obfuscation. *OpenAI Publication*, 2025. URL [https://cdn.openai.
579 com/pdf/34f2ada6-870f-4c26-9790-fd8def56387f/CoT_Monitoring.pdf](https://cdn.openai.com/pdf/34f2ada6-870f-4c26-9790-fd8def56387f/CoT_Monitoring.pdf).
- 580 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and
581 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sam-
582 pling. *Computing Research Repository*, arXiv:2407.21787, 2024. URL [https://arxiv.
583 org/abs/2407.21787](https://arxiv.org/abs/2407.21787).
- 584
- 585 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui
586 Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye
587 Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting
588 Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaying Li, Jingwen Li, Linyang Li,
589 Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun
590 Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang
591 Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song,
592 Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang,
593 Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingdong Xiong,
Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia
Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo

- 594 Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang,
595 Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng
596 Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. Internlm2 technical report. *Computing Research
597 Repository*, arXiv:2403.17297, 2024. URL <https://arxiv.org/abs/2403.17297>.
- 598
599 Maosong Cao, Alexander Lam, Haodong Duan, Hongwei Liu, Songyang Zhang, and Kai Chen.
600 Compassjudge-1: All-in-one judge model helps model evaluation and evolution. *Computing
601 Research Repository*, arXiv:2410.16256, 2024. URL <https://arxiv.org/abs/2410.16256>.
- 602
603 Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu
604 Chen. Codet: Code generation with generated tests. In *The Eleventh International Confer-
605 ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=
606 ktrw68Cmu9c](https://openreview.net/forum?id=ktrw68Cmu9c).
- 607
608 Nuo Chen, Zhiyuan Hu, Qingyun Zou, Jiaying Wu, Qian Wang, Bryan Hooi, and Bingsheng
609 He. Judgelrm: Large reasoning models as a judge. *Computing Research Repository*,
arXiv:2504.00050, 2025. URL <https://arxiv.org/abs/2504.00050>.
- 610
611 Yinlam Chow, Guy Tennenholtz, Izzeddin Gur, Vincent Zhuang, Bo Dai, Aviral Kumar, Rishabh
612 Agarwal, Sridhar Thiagarajan, Craig Boutilier, and Aleksandra Faust. Inference-aware fine-
613 tuning for best-of-n sampling in large language models. In *The Thirteenth International Confer-
614 ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=
615 77gQUdQhE7](https://openreview.net/forum?id=77gQUdQhE7).
- 616
617 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
618 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
619 Schulman. Training verifiers to solve math word problems. *Computing Research Repository*,
arXiv:2110.14168, 2021. URL <https://arxiv.org/abs/2110.14168>.
- 620
621 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong
622 Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. ULTRAFEEDBACK: Boost-
623 ing language models with scaled AI feedback. In Ruslan Salakhutdinov, Zico Kolter, Kather-
624 ine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Pro-
625 ceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proced-
626 ings of Machine Learning Research*, pp. 9722–9744. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/cui24f.html>.
- 627
628 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language
629 model. *Computing Research Repository*, arXiv:2405.04434, 2024a. URL [https://arxiv.
630 org/abs/2405.04434](https://arxiv.org/abs/2405.04434).
- 631
632 DeepSeek-AI. Deepseek-v3 technical report. *Computing Research Repository*, arXiv:2412.19437,
633 2024b. URL <https://arxiv.org/abs/2412.19437>.
- 634
635 DeepSeek-AI. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*,
636 645(7952):633–638, September 2025. doi: 10.1038/s41586-025-09422-z.
- 637
638 Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. Understanding dataset difficulty with
639 \mathcal{V} -usable information. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,
640 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine
641 Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 5988–6008. PMLR,
642 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/ethayarajh22a.
643 html](https://proceedings.mlr.press/v162/ethayarajh22a.html).
- 644
645 Jan-Philipp Fränken, Eric Zelikman, Rafael Rafailov, Kanishk Gandhi, Tobias Gerstenberg, and
646 Noah Goodman. Self-supervised alignment with mutual information: Learning to follow princi-
647 ples without preference labels. In *The Thirty-eighth Annual Conference on Neural Information
648 Processing Systems*, 2024. URL <https://openreview.net/forum?id=UvbbpEhGaw>.
- 649
650 Evan Frick, Tianle Li, Connor Chen, Wei-Lin Chiang, Anastasios Nikolas Angelopoulos, Jiantao
651 Jiao, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. How to evaluate reward models for
652 RLHF. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
653 <https://openreview.net/forum?id=cbttLtO94Q>.

- 648 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In An-
649 dreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan
650 Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume
651 202 of *Proceedings of Machine Learning Research*, pp. 10835–10866. PMLR, 23–29 Jul 2023.
652 URL <https://proceedings.mlr.press/v202/gao23h.html>.
- 653 Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Mari-
654 beth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham,
655 Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth
656 Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa
657 Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William
658 Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey
659 Irving. Improving alignment of dialogue agents via targeted human judgements. *Computing*
660 *Research Repository*, arXiv:2209.14375, 2022. URL [https://arxiv.org/abs/2209.](https://arxiv.org/abs/2209.14375)
661 [14375](https://arxiv.org/abs/2209.14375).
- 662 Fang Guo, Wenyu Li, Honglei Zhuang, Yun Luo, Yafu Li, Le Yan, Qi Zhu, and Yue Zhang.
663 Mcranker: Generating diverse criteria on-the-fly to improve pointwise llm rankers. In *Proceed-*
664 *ings of the Eighteenth ACM International Conference on Web Search and Data Mining, WSDM*
665 *'25*, pp. 944–953, New York, NY, USA, 2025. Association for Computing Machinery. ISBN
666 9798400713293. doi: 10.1145/3701551.3703583. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3701551.3703583)
667 [3701551.3703583](https://doi.org/10.1145/3701551.3703583).
- 668 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
669 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.
670 In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks*
671 *Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=7Bywt2mQsCe>.
- 672 Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. LLM-blender: Ensembling large language mod-
673 els with pairwise ranking and generative fusion. In Anna Rogers, Jordan Boyd-Graber, and
674 Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Com-*
675 *putational Linguistics (Volume 1: Long Papers)*, pp. 14165–14178, Toronto, Canada, July
676 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.792. URL
677 <https://aclanthology.org/2023.acl-long.792/>.
- 678 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
679 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth*
680 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=VTF8yNQm66)
681 [net/forum?id=VTF8yNQm66](https://openreview.net/forum?id=VTF8yNQm66).
- 682 Ryo Kamoi, Sarkar Snigdha Sarathi Das, Renze Lou, Jihyun Janice Ahn, Yilun Zhao, Xiaoxin
683 Lu, Nan Zhang, Yusen Zhang, Haoran Ranran Zhang, Sujeeth Reddy Vummanthala, Salika
684 Dave, Shaobo Qin, Arman Cohan, Wenpeng Yin, and Rui Zhang. Evaluating LLMs at de-
685 tecting errors in LLM responses. In *First Conference on Language Modeling*, 2024. URL
686 <https://openreview.net/forum?id=dnwRScljXr>.
- 687 M. G. Kendall and B. Babington Smith. On the method of paired comparisons. *Biometrika*, 31(3/4):
688 324–345, 1940. ISSN 00063444. URL <http://www.jstor.org/stable/2332613>.
- 689 Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Gra-
690 ham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Prometheus 2: An open source
691 language model specialized in evaluating other language models. In Yaser Al-Onaizan, Mohit
692 Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Meth-*
693 *ods in Natural Language Processing*, pp. 4334–4353, Miami, Florida, USA, November 2024.
694 Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.248. URL
695 <https://aclanthology.org/2024.emnlp-main.248/>.
- 696 Seungone Kim, Ian Wu, Jinu Lee, Xiang Yue, Seongyun Lee, Mingyeong Moon, Kiril Gashteovski,
697 Carolin Lawrence, Julia Hockenmaier, Graham Neubig, and Sean Welleck. Scaling evaluation-
698 time compute with reasoning models as process evaluators. *Computing Research Repository*,
699 arXiv:2503.19877, 2025. URL <https://arxiv.org/abs/2503.19877>.
- 700
701

- 702 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
703 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Re-
704 wardbench: Evaluating reward models for language modeling. *Computing Research Repository*,
705 arXiv:2403.13787, 2024. URL <https://arxiv.org/abs/2403.13787>.
- 706 Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, hai zhao, and Pengfei Liu. Generative judge
707 for evaluating alignment. In *The Twelfth International Conference on Learning Representations*,
708 2024a. URL <https://openreview.net/forum?id=gtkFw6sZGS>.
- 709 Lei Li, Yekun Chai, Shuohuan Wang, Yu Sun, Hao Tian, Ningyu Zhang, and Hua Wu. Tool-
710 augmented reward modeling. In *The Twelfth International Conference on Learning Representations*,
711 2024b. URL <https://openreview.net/forum?id=d94x0gWTUX>.
- 712 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
713 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following
714 models. https://github.com/tatsu-lab/alpaca_eval, 5 2023.
- 715 Shalev Lifshitz, Sheila A. McIlraith, and Yilun Du. Multi-agent verification: Scaling test-time
716 compute with multiple verifiers. In *Second Conference on Language Modeling*, 2025. URL
717 <https://openreview.net/forum?id=LriQ3NY9uL>.
- 718 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
719 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth
720 International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=v8L0pN6EOi>.
- 721 Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang
722 Liu, and Yahui Zhou. Skywork-reward: Bag of tricks for reward modeling in llms. *Computing
723 Research Repository*, arXiv:2410.18451, 2024. URL <https://arxiv.org/abs/2410.18451>.
- 724 Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. Pairjudge rm: Perform best-
725 of-n sampling with knockout tournament. *Computing Research Repository*, arXiv:2501.13007,
726 2025. URL <https://arxiv.org/abs/2501.13007>.
- 727 Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato,
728 Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. Generative reward models. *Computing
729 Research Repository*, arXiv:2410.12832, 2024. URL <https://arxiv.org/abs/2410.12832>.
- 730 Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D Kivlichan, Molly
731 Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for language model
732 safety. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
733 URL <https://openreview.net/forum?id=QVtwpT5Dmg>.
- 734 Niklas Muennighoff, Qian Liu, Armel Randy Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo,
735 Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. Octopack: Instruc-
736 tion tuning code large language models. In *The Twelfth International Conference on Learning
737 Representations*, 2024. URL <https://openreview.net/forum?id=mw1PWNSWZP>.
- 738 OpenAI. Openai o1 system card. *Computing Research Repository*, arXiv:2412.16720, 2024. URL
739 <https://arxiv.org/abs/2412.16720>.
- 740 OpenAI. Deep research system card. *OpenAI Publication*, 2025a. URL <https://cdn.openai.com/deep-research-system-card.pdf>.
- 741 OpenAI. Openai gpt-4.5 system card. *OpenAI Publication*, 2025b. URL <https://cdn.openai.com/gpt-4-5-system-card-2272025.pdf>.
- 742 OpenAI. Openai o3-mini system card. *OpenAI Publication*, 2025c. URL <https://cdn.openai.com/o3-mini-system-card-feb10.pdf>.
- 743 OpenAI. Grader guide. <https://platform.openai.com/docs/guides/graders>,
744 2025d. Accessed: 2025-09-25.

- 756 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
757 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
758 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
759 and Ryan Lowe. Training language models to follow instructions with human feedback. In *Pro-
760 ceedings of the 36th International Conference on Neural Information Processing Systems, NIPS
761 '22*, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- 762 Junsoo Park, Seungyeon Jwa, Ren Meiyang, Daeyoung Kim, and Sanghyuk Choi. OffsetBias:
763 Leveraging debiased data for tuning evaluators. In Yaser Al-Onaizan, Mohit Bansal, and Yun-
764 Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*,
765 pp. 1043–1067, Miami, Florida, USA, November 2024. Association for Computational Lin-
766 guistics. doi: 10.18653/v1/2024.findings-emnlp.57. URL [https://aclanthology.org/
767 2024.findings-emnlp.57/](https://aclanthology.org/2024.findings-emnlp.57/).
- 768 Hao Peng, Yunjia Qi, Xiaozhi Wang, Zijun Yao, Bin Xu, Lei Hou, and Juanzi Li. Agentic reward
769 modeling: Integrating human preferences with verifiable correctness signals for reliable reward
770 systems. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar
771 (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics
772 (Volume 1: Long Papers)*, pp. 15934–15949, Vienna, Austria, July 2025. Association for Com-
773 putational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.775. URL
774 <https://aclanthology.org/2025.acl-long.775/>.
- 775 Paul Röttger, Hannah Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy.
776 XSTest: A test suite for identifying exaggerated safety behaviours in large language models.
777 In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Confer-
778 ence of the North American Chapter of the Association for Computational Linguistics: Human
779 Language Technologies (Volume 1: Long Papers)*, pp. 5377–5400, Mexico City, Mexico, June
780 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.301. URL
781 <https://aclanthology.org/2024.naacl-long.301/>.
- 782 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
783 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of math-
784 ematical reasoning in open language models. *Computing Research Repository*, arXiv:2402.0330,
785 2024. URL <https://arxiv.org/abs/2402.03300>.
- 786 Mrinank Sharma, Meg Tong, Jesse Mu, Jerry Wei, Jorrit Kruthoff, Scott Goodfriend, Euan Ong, Al-
787 win Peng, Raj Agarwal, Cem Anil, Amanda Askell, Nathan Bailey, Joe Benton, Emma Blumke,
788 Samuel R. Bowman, Eric Christiansen, Hoagy Cunningham, Andy Dau, Anjali Gopal, Rob
789 Gilson, Logan Graham, Logan Howard, Nimit Kalra, Taesung Lee, Kevin Lin, Peter Lofgren,
790 Francesco Mosconi, Clare O’Hara, Catherine Olsson, Linda Petrini, Samir Rajani, Nikhil Sax-
791 ena, Alex Silverstein, Tanya Singh, Theodore Sumers, Leonard Tang, Kevin K. Troy, Constantin
792 Weisser, Ruiqi Zhong, Giulio Zhou, Jan Leike, Jared Kaplan, and Ethan Perez. Constitutional
793 classifiers: Defending against universal jailbreaks across thousands of hours of red teaming. *Com-
794 puting Research Repository*, arXiv:2501.18837, 2025. URL [https://arxiv.org/abs/
795 2501.18837](https://arxiv.org/abs/2501.18837).
- 796 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
797 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *EuroSys*, pp.
798 1279–1297, 2025. URL <https://doi.org/10.1145/3689031.3696075>.
- 800 Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute
801 optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth Interna-
802 tional Conference on Learning Representations, 2025*. URL [https://openreview.net/
803 forum?id=4FWAwZtd2n](https://openreview.net/forum?id=4FWAwZtd2n).
- 804 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
805 Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *Proced-
806 ings of the 34th International Conference on Neural Information Processing Systems, NIPS '20*,
807 Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 808 Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Daniel Cox,
809 Yiming Yang, and Chuang Gan. Principle-driven self-alignment of language models from scratch

- 810 with minimal human supervision. In *Thirty-seventh Conference on Neural Information Processing*
811 *Systems*, 2023. URL <https://openreview.net/forum?id=p40XRfBX96>.
812
- 813 Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong Zhou, Zhenfang Chen, David Daniel Cox,
814 Yiming Yang, and Chuang Gan. SALMON: Self-alignment with instructable reward models.
815 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=xJbsmB8UMx>.
816
- 817 Gemma Team. Gemma 2: Improving open language models at a practical size. *Computing Research*
818 *Repository*, arXiv:2408.0011, 2024. URL <https://arxiv.org/abs/2408.00118>.
819
- 820 Hemish Veeraboina. Aime problem set 1983-2024, 2023. URL <https://www.kaggle.com/datasets/hemishveeraboina/aime-problem-set-1983-2024>.
821
- 822 Tu Vu, Kalpesh Krishna, Salaheddin Alzubi, Chris Tar, Manaal Faruqui, and Yun-Hsuan Sung.
823 Foundational autoraters: Taming large language models for better automatic evaluation. In Yaser
824 Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on*
825 *Empirical Methods in Natural Language Processing*, pp. 17086–17105, Miami, Florida, USA,
826 November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.
827 949. URL <https://aclanthology.org/2024.emnlp-main.949/>.
828
- 829 Chenglong Wang, Yang Gan, Yifu Huo, Yongyu Mu, Qiaozhi He, MuRun Yang, Bei Li, Tong
830 Xiao, Chunliang Zhang, Tongran Liu, and JingBo Zhu. GRAM: A generative foundation reward
831 model for reward generalization. In *Forty-second International Conference on Machine Learning*,
832 2025a. URL <https://openreview.net/forum?id=rxKC8v2uHc>.
- 833 Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences
834 via multi-objective reward modeling and mixture-of-experts. In Yaser Al-Onaizan, Mohit Bansal,
835 and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP*
836 *2024*, pp. 10582–10592, Miami, Florida, USA, November 2024a. Association for Computational
837 Linguistics. doi: 10.18653/v1/2024.findings-emnlp.620. URL <https://aclanthology.org/2024.findings-emnlp.620/>.
838
- 839 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhi-
840 fang Sui. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In
841 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meet-*
842 *ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9426–9439,
843 Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/
844 v1/2024.acl-long.510. URL <https://aclanthology.org/2024.acl-long.510/>.
845
- 846 Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu,
847 Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught eval-
848 uators. *Computing Research Repository*, arXiv:2408.02666, 2024c. URL <https://arxiv.org/abs/2408.02666>.
849
- 850 Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer:
851 Evaluating safeguards in LLMs. In Yvette Graham and Matthew Purver (eds.), *Findings of the*
852 *Association for Computational Linguistics: EAACL 2024*, pp. 896–911, St. Julian’s, Malta, March
853 2024d. Association for Computational Linguistics. URL <https://aclanthology.org/2024.findings-eaACL.61/>.
854
- 855 Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang,
856 Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. Helpsteer 2: Open-source dataset for training
857 top-performing reward models. In *The Thirty-eight Conference on Neural Information Process-*
858 *ing Systems Datasets and Benchmarks Track*, 2024e. URL <https://openreview.net/forum?id=PvVKUFhaNy>.
859
- 860 Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii
861 Kuchaiev, and Yi Dong. Helpsteer2-preference: Complementing ratings with preferences. In
862 *The Thirteenth International Conference on Learning Representations*, 2025b. URL <https://openreview.net/forum?id=MnfHxPP5gs>.
863

- 864 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,
865 Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language
866 models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Ad-
867 vances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/
868 forum?id=_VjQlMeSB_J](https://openreview.net/forum?id=_VjQlMeSB_J).
- 869
870 Genta Indra Winata, David Anugraha, Lucky Susanto, Garry Kuwanto, and Derry Tanti Wijaya.
871 Metametrics: Calibrating metrics for generation tasks using human preferences. In *The Thirteenth
872 International Conference on Learning Representations*, 2025. URL [https://openreview.
873 net/forum?id=s103xTt4CG](https://openreview.net/forum?id=s103xTt4CG).
- 874 Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason Weston,
875 and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment with
876 llm-as-a-meta-judge. *Computing Research Repository*, arXiv:2407.19594, 2024. URL <https://arxiv.org/abs/2407.19594>.
- 877
878 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws:
879 An empirical analysis of compute-optimal inference for LLM problem-solving. In *The Thirteenth
880 International Conference on Learning Representations*, 2025. URL [https://openreview.
881 net/forum?id=VNckp7JEHn](https://openreview.net/forum?id=VNckp7JEHn).
- 882
883 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing
884 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio
885 Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. OSWorld: Benchmarking multimodal
886 agents for open-ended tasks in real computer environments. In *The Thirty-eight Conference on
887 Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL [https://
888 openreview.net/forum?id=tN61DTr4Ed](https://openreview.net/forum?id=tN61DTr4Ed).
- 889
890 Shunyu Yao, Howard Chen, John Yang, and Karthik R Narasimhan. Webshop: Towards scalable
891 real-world web interaction with grounded language agents. In Alice H. Oh, Alekh Agarwal,
892 Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Sys-
893 tems*, 2022. URL <https://openreview.net/forum?id=R9KnuFlvnU>.
- 894
895 Zihuiwen Ye, Fraser David Greenlee, Max Bartolo, Phil Blunsom, Jon Ander Campos, and Matthias
896 Gallé. Improving reward models with synthetic critiques. In Luis Chiruzzo, Alan Ritter, and
897 Lu Wang (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*, pp.
898 4506–4520, Albuquerque, New Mexico, April 2025a. Association for Computational Linguistics.
899 ISBN 979-8-89176-195-7. doi: 10.18653/v1/2025.findings-naacl.254. URL <https://aclanthology.org/2025.findings-naacl.254/>.
- 900
901 Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun
902 LIU. Learning LLM-as-a-judge for preference alignment. In *The Thirteenth International Con-
903 ference on Learning Representations*, 2025b. URL [https://openreview.net/forum/
904 id=HZVIQE1MsJ](https://openreview.net/forum?id=HZVIQE1MsJ).
- 905
906 Jiachen Yu, Shaoning Sun, Xiaohui Hu, Jiayu Yan, Kaidong Yu, and Xuelong Li. Improve llm-as-
907 a-judge ability as a general ability. *Computing Research Repository*, arXiv:2502.11689, 2025a.
908 URL <https://arxiv.org/abs/2502.11689>.
- 909
910 Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan, Chenguang Zhu, Richard Yuanzhe Pang, Yundi
911 Qian, Xuwei Wang, Suchin Gururangan, Chao Zhang, Melanie Kambadur, Dhruv Mahajan,
912 and Rui Hou. Self-generated critiques boost reward modeling for language models. In Luis
913 Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations
914 of the Americas Chapter of the Association for Computational Linguistics: Human Language
915 Technologies (Volume 1: Long Papers)*, pp. 11499–11514, Albuquerque, New Mexico, April
916 2025b. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/
917 2025.naacl-long.573. URL <https://aclanthology.org/2025.naacl-long.573/>.
- 918
919 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu,
920 and Jason E Weston. Self-rewarding language models. In Ruslan Salakhutdinov, Zico Kolter,
921 Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.),

- 918 *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Pro-*
919 *ceedings of Machine Learning Research*, pp. 57905–57923. PMLR, 21–27 Jul 2024. URL
920 <https://proceedings.mlr.press/v235/yuan24d.html>.
921
- 922 Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. Evaluating
923 large language models at evaluating instruction following. In *The Twelfth International Confer-*
924 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=tr0KidwPLc)
925 [tr0KidwPLc](https://openreview.net/forum?id=tr0KidwPLc).
- 926 Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agar-
927 wal. Generative verifiers: Reward modeling as next-token prediction. In *The Thirteenth Interna-*
928 *tional Conference on Learning Representations*, 2025a. URL [https://openreview.net/](https://openreview.net/forum?id=Ccwp4tFEtE)
929 [forum?id=Ccwp4tFEtE](https://openreview.net/forum?id=Ccwp4tFEtE).
- 930 Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu,
931 Jingren Zhou, and Junyang Lin. The lessons of developing process reward models in mathematical
932 reasoning. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar
933 (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 10495–10516,
934 Vienna, Austria, July 2025b. Association for Computational Linguistics. ISBN 979-8-89176-256-
935 5. doi: 10.18653/v1/2025.findings-acl.547. URL [https://aclanthology.org/2025.](https://aclanthology.org/2025.findings-acl.547/)
936 [findings-acl.547/](https://aclanthology.org/2025.findings-acl.547/).
- 937 Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu,
938 Jingren Zhou, and Junyang Lin. ProcessBench: Identifying process errors in mathematical rea-
939 soning. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar
940 (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*
941 *(Volume 1: Long Papers)*, pp. 1009–1024, Vienna, Austria, July 2025. Association for Com-
942 putational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.50. URL
943 <https://aclanthology.org/2025.acl-long.50/>.
944
- 945 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
946 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
947 Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International*
948 *Conference on Neural Information Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2023.
949 Curran Associates Inc.
- 950 Enyu Zhou, Guodong Zheng, Binghai Wang, Zhiheng Xi, Shihan Dou, Rong Bao, Wei Shen, Limao
951 Xiong, Jessica Fan, Yurong Mou, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. RMB:
952 Comprehensively benchmarking reward models in LLM alignment. In *The Thirteenth Interna-*
953 *tional Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=kmgrlG9TR0)
954 [forum?id=kmgrlG9TR0](https://openreview.net/forum?id=kmgrlG9TR0).
- 955 Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam
956 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen GONG, James
957 Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan
958 Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang,
959 David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Le-
960 andro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls and
961 complex instructions. In *The Thirteenth International Conference on Learning Representations*,
962 2025. URL <https://openreview.net/forum?id=YrycTjlll0>.
963
964
965
966
967
968
969
970
971

972	CONTENTS	
973		
974	1 Introduction	1
975		
976	2 Preliminaries	3
977		
978	2.1 Comparisons of Different RM approaches	3
979	2.2 Boosting Reward Quality with Principles	4
980		
981	3 Self-Principled Critique Tuning (SPCT)	4
982		
983	3.1 Unpinning Principles from Understanding to Generation	4
984	3.2 Rule-Based Reinforcement Learning	5
985		
986	4 Inference-Time Scaling with SPCT	6
987		
988	5 Results on Reward Modeling Benchmarks	7
989		
990	5.1 Experiment Settings	7
991	5.2 Results and Analysis	7
992		
993	6 Related Work	9
994		
995	7 Conclusion and Future Work	9
996		
997	A The Use of Large Language Models (LLMs)	21
998		
999	B Additional Related Work	21
1000		
1001	C Limitations and Future Directions	21
1002		
1003	D Implementation Details	22
1004		
1005	D.1 Comparisons of Different RM Approaches	22
1006	D.2 Model Training	23
1007	D.3 Baseline Implementation	24
1008		
1009	E Experiment Details	24
1010		
1011	E.1 Hyper-Parameters	24
1012	E.2 Benchmarks	24
1013	E.3 Detailed Results	26
1014		
1015	F Additional Experiments	27
1016		
1017	F.1 Input Flexibility of the Pointwise GRM Approach	27
1018	F.1.1 Generating Rewards for Many Responses	27
1019	F.1.2 Generating Rewards for Single Responses	28
1020	F.1.3 Generating Rewards with Reference	28
1021	F.2 Transferability of Generated Principles	28
1022		
1023		
1024		
1025		

1026	F.3 Generalization beyond Training Data	29
1027		
1028	F.4 Response Length Analysis for Rule-Based RL	29
1029		
1030	G Qualitative Analysis	29
1031	G.1 Case Study	29
1032		
1033	G.2 Failure Mode Analysis	39
1034		
1035	H Prompt Templates	40
1036		
1037		
1038		
1039		
1040		
1041		
1042		
1043		
1044		
1045		
1046		
1047		
1048		
1049		
1050		
1051		
1052		
1053		
1054		
1055		
1056		
1057		
1058		
1059		
1060		
1061		
1062		
1063		
1064		
1065		
1066		
1067		
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

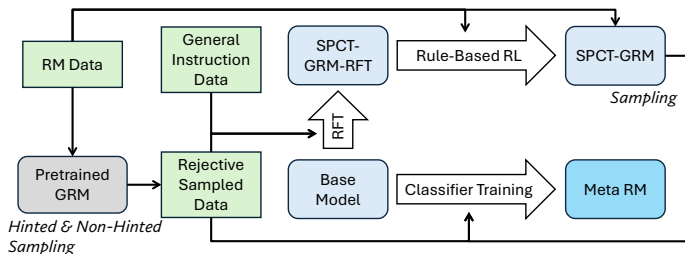


Figure 5: Illustration of the derivation of SPCT-GRM-RFT, SPCT-GRM, and Meta RM in the SPCT pipeline.

A THE USE OF LARGE LANGUAGE MODELS (LLMs)

We promise that all the LLM generated content in this work is solely used as research objectives or tools to polish writings in grammar or sentence levels. We do not use LLMs to generate any part of the scientific content from scratch, including but not limited to the ideas, algorithms, experimental results, and analysis. We human authors will take full responsibility for the presented work.

B ADDITIONAL RELATED WORK

Constitutional AI Constitutional AI has emerged as a scalable alternative to traditional reinforcement learning from human feedback (Ouyang et al., 2022), aiming to align language models with human values through a set of guiding principles or “constitutions” (Bai et al., 2022b; Sun et al., 2023; 2024), replacing human critiques with AI-generated feedback (Fränken et al., 2024) or classifiers (Sharma et al., 2025) based on these hand-crafted principles. Similarly, rule-based approaches like Sparrow (Glaese et al., 2022) and Rule-Based Rewards (RBR) (Mu et al., 2024) incorporate explicit natural language rules into the training loop for specific domains like safety. Although effective, these methods rely on static, manually written constitutions that are limited in scope, potentially biased, and inflexible. This has motivated interests in automating the generation or refinement of principles, which aligns with our target in this work.

Scalar Reward Models Scalar reward modeling for LLMs are proposed the earliest to serve as a proxy model for human feedback (Stiennon et al., 2020; Gao et al., 2023). Recent studies focus on Bradley-Terry modeling (Kendall & Smith, 1940) and other regression approaches for better expressiveness for scalar reward models (Cai et al., 2024; Wang et al., 2024e;a; Liu et al., 2024; Wang et al., 2025b) of general preference. In contrast to these outcome reward models, process reward models are proposed as step verifiers for reasoning problems, e.g., math, etc. (Cobbe et al., 2021; Wang et al., 2024b; Zhang et al., 2025b), demonstrating the feasibility of scalar RMs in a formal domain with extensive reasoning and knowledge. Scalar RM excels in simplicity and is computationally efficient, but suffers from limited expressivity and struggles to generalize across diverse input types or refine reward signals at inference time.

Semi-Scalar Reward Models Semi-scalar reward models aim to enrich scalar reward signals through textual intermediate representations. (Ye et al., 2025a; Ankner et al., 2024) Consequently, works (Yu et al., 2025b) proposed to enhance the quality of generated critiques to eventually improve reward generation. Some studies use the token probability to substitute the scalar head for reward extraction (Mahan et al., 2024; Zhang et al., 2025a). These works show that semi-scalar RMs face challenges in inference-time scaling based on sampling and voting, resulting in limited performance improvement. The semi-scalar approach trades off between scalar RMs and GRMs in terms of both efficiency and effectiveness.

C LIMITATIONS AND FUTURE DIRECTIONS

Limitations Though SPCT significantly leverages the performance and inference-time scalability of GRMs and surpasses (semi-)scalar RMs in general domains, it still faces a few limitations. (1)

The efficiency of the generative RMs is largely lagging behind the scalar RMs at the same scale by nature, which inhibits its large-scale usage in online RL pipelines. However, since we adopt parallel sampling for inference-time scaling, the latency of reward generation with a reasonable amount of, e.g., eight samplings will not increase significantly. Further research around the efficient generation of LLMs and innovations in RM applications could alleviate the problem. (2) In specific domains such as verifiable tasks, SPCT-GRM still lags behind scalar models. This could be because the scalar RMs capture hidden features of reasoning queries and responses, while GRMs need stronger reasoning capabilities to examine responses thoroughly. However, scalar RMs suffer severe biases and scalability issues. For GRMs, we found that both reference-based reward generation (Appendix F.1.3 and long-horizon reasoning (Appendix E.3) could mitigate this limitation. (3) Due to the universality of the pointwise GRM approach, SPCT-GRM could potentially serve as a process RM in addition to the outcome RM. Though we have not explored much in this direction in the paper, the performance in the Reasoning subset of Reward Bench, which mainly comprises of MATH-prm data (Lightman et al., 2024), could partially support the potential of this application.

Future Directions There are also several promising directions for future research based on SPCT or SPCT-GRM models. (1) Tool incorporation of RMs is studied by previous work (Li et al., 2024b), and could also be used for SPCT-GRM augmentation. **With tools such as code interpreters and search engine interfaces**, the generated critiques could be more accurate for tasks that requires strict procedures or extensive knowledge, and the cases in which GRMs fail to follow principles related to numeric calculations, pattern matching, etc. could be avoided. (2) **The generation paradigm for principles and critiques could be decomposed** into separate stages, that is, the principles could be generated ahead of time for each query and the responses to be rated and stored, and then the critiques are generated with GRMs, rules, or other agentic approaches. The principle generation serves as an interface for the following critiques. This might improve the efficiency of current GRMs for the integration of RL pipelines. (3) The SPCT-GRM could be potentially **used in LLM offline evaluation**. Since each principle reflects a criteria, we can get criteria from all data points that a particular LLM is inferior than one another, as a interpretable protocol of the weaknesses of the particular LLM. (4) The SPCT-GRM might **be benefit from long-horizon reasoning**. However, this will further affect its efficiency. These directions should be studied in the future work.

D IMPLEMENTATION DETAILS

D.1 COMPARISONS OF DIFFERENT RM APPROACHES

Reward Generation Paradigms Classic RMs adopt the **(a) scalar** approach to generate rewards (\mathcal{R}), which assigns scalar values to the given query and responses. The scalar approach is further extended to the **(b) semi-scalar** approach, which generates texts besides the scalar value. And the **(c) generative** approach only generates textual rewards.

$$\mathcal{R} = \begin{cases} \mathcal{S} & \text{(Scalar)} \\ (\mathcal{S}, \mathcal{C}) & \text{(Semi-Scalar)} \sim r_\theta(x, \{y_i\}_{i=1}^n), \\ \mathcal{C} & \text{(Generative)} \end{cases}, \quad (9)$$

where x is the query, y_i is the i -th response, r_θ is the reward function parameterized by θ , $\mathcal{S} \in \mathbb{R}^m$, $m \leq n$ is the scalar reward, and \mathcal{C} is the critique.

Scoring Patterns We distinguish two main scoring approaches for rewards: pointwise and pairwise. The **(i) pointwise** approach assigns an individual score to each response:

$$\{S_i\}_{i=1}^n = f_{\text{point}}(\mathcal{R}, \{y_i\}_{i=1}^n), \quad \mathcal{R} \sim r_\theta(x, \{y_i\}_{i=1}^n), S_i \in \mathbb{R}, \quad (10)$$

where $f_{\text{point}}(\cdot, \cdot)$ is a splitting function. In contrast, the **(ii) pairwise** approach can be viewed as a best-of- n method, selecting a single best response from all candidates:

$$\hat{y} = f_{\text{pair}}(\mathcal{R}, \{y_i\}_{i=1}^n), \quad \mathcal{R} \sim r_\theta(x, \{y_i\}_{i=1}^n), \hat{y} \in \{y_i\}_{i=1}^n, \quad (11)$$

where $f_{\text{pair}}(\cdot, \cdot)$ is a selection function and $n = 2$ in most cases. Though the pairwise approach could be extended to $n > 2$, it could not be applied to single response scoring ($n = 1$).

Representative Methods Figure 2 illustrates how the three reward generation paradigms (scalar, semi-scalar, generative) can be combined with the two scoring patterns (pointwise, pairwise). Specifically, Bradley-Terry model (Kendall & Smith, 1940) ((a)+(ii)) is trained with pairwise preference data and outputs scalar rewards pointwisely

$$\{S_i\}_{i=1}^n = f_{\text{point}}(\mathcal{R}, \{y_i\}_{i=1}^n) = \mathcal{S} \in \mathbb{R}^n. \quad (12)$$

PairRM (Jiang et al., 2023) ((a)+(ii)) compares a pair of responses with the sign of the scalar reward

$$\hat{y} = f_{\text{pair}}(\mathcal{R}, \{y_i\}_{i=1}^n) = y_{\lfloor \frac{1}{2}(3 - \text{sgn}(\mathcal{S})) \rfloor}, \quad n = 2, \mathcal{S} \in \mathbb{R}. \quad (13)$$

The scalar methods above could barely perform inference-time scaling due to the lack of diversity in reward generation. CLOUD (Ankner et al., 2024) ((b)+(i)) generates scalar rewards for each response based on pre-generated critiques, similar to Equation 12. LLM-as-a-Judge (Zheng et al., 2023; Wang et al., 2025a) ((c)+(ii)) judges the preference order between paired responses textually,

$$\hat{y} = f_{\text{pair}}(\mathcal{R}, \{y_i\}_{i=1}^n) = y_{f_{\text{extract}}(\mathcal{C})}, \quad n = 2, \quad (14)$$

where $f_{\text{extract}}(\cdot)$ extracts the index of best response from language representations. However, this approach defaults to neglect ties of the paired responses. Following Zhang et al. (2025a), the generation probability of the token that indicates the preference order could be used as the scalar reward ((b)+(ii)): $\mathcal{S} = \text{TokenProb}(\hat{C}) = r_{\theta}(\hat{C}|x, \{y_i\}_{i=1}^n)$, where \hat{C} is a pre-defined token related to the preference order. However, without additional constraints, GRMs are able to generate pointwise rewards for multiple responses within pure language representation ((c)+(i)):

$$\{S_i\}_{i=1}^n = f_{\text{point}}(\mathcal{R}, \{y_i\}_{i=1}^n) = f_{\text{extract}}(\mathcal{C}), \quad (15)$$

where $f_{\text{extract}}(\cdot)$ extracts the rewards assigned to each response from generation results. Usually, the rewards are discrete, and in this work we assign $S_i \in \mathbb{N}, 1 \leq S_i \leq 10$ by default. This approach promisingly allows both inference-time scalability and input flexibility.

D.2 MODEL TRAINING

For the rule-based online RL, we use the standard GRPO setting (Shao et al., 2024). The overall objective is

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)] \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[\frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \text{clip} \left(\frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta} || \pi_{\text{ref}}] \right\}, \quad (16)$$

where $\hat{A}_{i,t} = \frac{\hat{r}_i - \text{mean}(\hat{\mathbf{r}})}{\text{std}(\hat{\mathbf{r}})}$, G is the group size, β is the coefficient of KL penalty, and $q = (x, \{y_i\}_{i=1}^n)$ with prompts. We performed grid search on hyper-parameter $\beta \in \{0.00, 0.01, 0.02, 0.08\}$ and found that $\beta = 0.08$ is the most stable configuration for SPCT-GRM-27B. And with too small KL coefficient, SPCT-GRM-27B tends to collapse on a few subsets in benchmarks, e.g., Chat in the Reward Bench benchmark and Harmlessness in the RMB benchmark, and shows biases towards some other domains. For smaller SPCT-GRM-16B, we use $\beta = 0.002$ because it is less vulnerable to the KL loss coefficient. We set $G = 4$ for a better trade-off between efficiency and performance.

The training set comprises of 1256K RFT data, including 1070K general instruction data and 186K rejective sampled data, and 237K RL data. General instruction data are from in-house datasets. Rejective sampled data and RL data are from the same RM datasets, containing the preference for single, paired, and multiple responses, constructed from internal data and open-source datasets, including the training sets from MATH (Hendrycks et al., 2021), UltraFeedback (Cui et al., 2024), Offset-Bias (Park et al., 2024), Skywork-Reward-Preference-80K-v0.2 (Liu et al., 2024), and HelpSteer2-Preference (Wang et al., 2025b). Specifically, we re-tagged the preference label of a part of UltraFeedback due to its quality issues; we sampled and filtered trajectories on MATH by rule-based ground-truth matching, resulting in pairwise preference data; for rating single responses, we set the ground-truth reward to 1 for correct responses and 0 for incorrect ones, only incorporating verifiable

questions. For rejective sampling, we use DeepSeek-v2.5-0905 to generate the trajectories with principles and critiques. The sampling time N_{RFT} is set to 3. During hinted sampling on HelpSteer2, we add the preference strengths labeled in the original dataset as the hint. We also remove the samples that are viewed too easy for DeepSeek-V2-Lite-Chat, i.e. all generated rewards are correct for three times according to Equation 4, from the RL data.

The derivation of SPCT-GRM models and the meta RM is illustrated in Figure 5. All SPCT-GRM models are trained from the pretrained version of LLMs. For the training of the meta RM, we reuse the rejective sampled data from the RFT stage, and use SPCT-GRM-27B to perform rejective sampling with $N_{\text{RFT}} = 3$, in order to avoid potential bias (Chow et al., 2025) in the meta RM guided voting. The learning rate is 1×10^{-5} and the batch size is 512 for the meta RM training. The training time of RFT and RL for SPCT-GRM-27B is depicted in Table 5, Gemma-2-27B based models are trained with 128 A100 GPUs on the Fire-Flyer platform (An et al., 2024). The learning rate is 5×10^{-6} for the RFT stage and 4×10^{-7} for the RL stage, and the batch size is 1024 for the RFT stage and 512 for the RL stage. Both stages are trained for 900 steps. Due to resource constraints, SPCT-GRM models larger than 27B does not undergo the rule-based RL and only trained with 50K rejective sampled data.

Table 5: Training times of RFT and RL stages for SPCT-GRM-27B in hours.

Stage	Time (h)
RFT	19.2
Rule-Based RL	15.6

D.3 BASELINE IMPLEMENTATION

For the baseline methods, we re-implement **LLM-as-a-Judge** (Zheng et al., 2023), **BT-RM-27B** (Kendall & Smith, 1940), **CLOUD-GEMMA-2-27B** (Ankner et al., 2024), and **PAIR-RM-27B** (Jiang et al., 2023) based on Gemma-2-27B (Team, 2024) and with all compatible training data and settings as SPCT-GRM.

For **LLM-as-a-Judge**, we use exactly the same training configuration as SPCT-GRM-27B, including RFT with rejective sampled data from DeepSeek-v2.5-0905 and rule-based online RL. Due to its scoring pattern, only pairwise data could be used in the RL stage. For **CLOUD-GEMMA-2-27B**, we also generate pointwise critiques from DeepSeek-v2.5-0905 using the same prompt template. However, it is not feasible to perform rejective sampling, since no rewards could be extracted without a trained value head. We fine-tune Gemma-2-27B with the same general instruction data of SPCT-GRM-27B along with the sampled critique, resulting in a critique generation model. Specifically, we fine-tune another Gemma-2-27B model with a value head for reward generation, instead of training value heads post hoc on the critique model. The training of the value head of CLOUD-GEMMA-2-27B, **BT-RM-27B**, and **PAIR-RM-27B** (Jiang et al., 2023) uses the same dataset from the RL stage of SPCT-GRM-27B, except for single response rating data.

E EXPERIMENT DETAILS

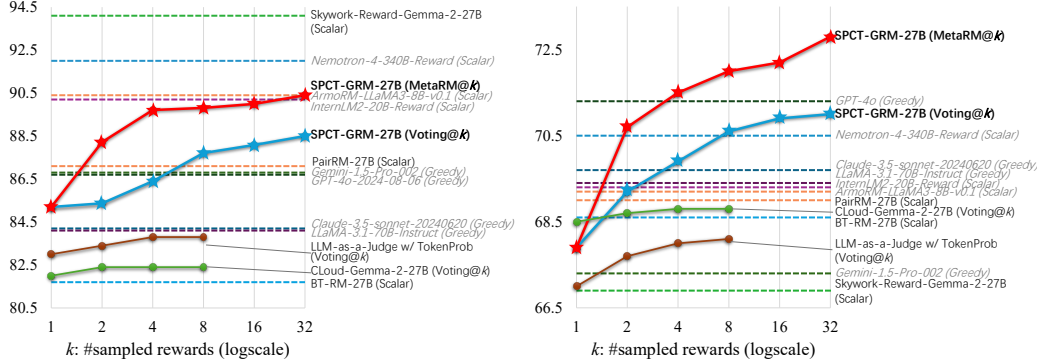
E.1 HYPER-PARAMETERS

For inference-time scaling results of SPCT-GRM-27B, SPCT-GRM-16B, LLM-as-a-Judge, and CLOUD-GEMMA-2-27B, the temperature is set to 0.5 for each model. And for other experiments, temperature is set to 0 for all models. Without specific description, $k_{\text{meta}} = \frac{1}{2}k$ by default in the meta RM guided voting for SPCT-GRM-27B. For inference on DeepSeek-R1-0120, the temperature is set to 0.6. Please note that we let SPCT-GRM to output rewards in the same range for rating single responses in the ReaLMistake benchmark as other benchmarks.

E.2 BENCHMARKS

We evaluate the performance of different methods on various RM benchmarks of different domains: (1) **Reward Bench** (Lambert et al., 2024), a common benchmark for RM evaluation, with semi-automatically collected chat (Li et al., 2023; Zheng et al., 2023; Zeng et al., 2024), reasoning (Lightman et al., 2024; Muennighoff et al., 2024), and safety (Röttger et al., 2024; Wang et al., 2024d) preference data, where two responses require to be ranked for each query; (2) **PPE** (Frick et al., 2025), a large-scale benchmark containing crowdsourced preference data and correctness data for

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349



(a) Results on the Reward Bench benchmark. (b) Results on all tested reward modeling benchmarks.

Figure 6: Inference-time scaling performance with different RMs on different reward modeling benchmarks. Non-italic font indicates models based on Gemma-2-27B.

Table 6: Detailed results of inference-time scalability experiments (Table 3) of different methods and models on RM benchmarks. Underlined numbers indicate the best performance, **bold numbers** indicate the best performance among baseline and our methods, and *italicized font* denotes scalar or semi-scalar RMs. For meta RM guided voting (MetaRM), $k_{meta} = \frac{1}{2}k$. Numbers in the parentheses is the performance change after inference-time scaling.

Model	Reward Bench	PPE Preference	PPE Correctness	RMB	Overall
<i>Reported Results of Public Models</i>					
<i>Nemotron-4-340B-Reward</i>	<u>92.0</u>	59.3	60.8	69.9	70.5
GPT-4o	86.7	67.1	57.6	<u>73.8</u>	71.3
<i>Results of Inference-Time Scaling (Voting@1)</i>					
LLM-as-a-Judge	83.0	63.4	57.4	64.3	67.0
<i>C</i> Cloud-Gemma-2-27B	82.0	67.0	62.0	63.2	68.5
SPCT-GRM-27B-RFT (Ours)	84.0	62.2	59.4	65.8	67.8
SPCT-GRM-27B (Ours)	85.2	62.4	59.5	64.4	67.9
<i>Results of Inference-Time Scaling (Voting@8)</i>					
LLM-as-a-Judge	83.4	63.8	58.2	65.2	67.6 (+0.6)
LLM-as-a-Judge w/ TokenProb	83.8	64.6	58.8	65.2	68.1 (+1.1)
<i>C</i> Cloud-Gemma-2-27B	82.4	<u>67.3</u>	62.4	63.2	68.8 (+0.3)
SPCT-GRM-27B-RFT (Ours)	85.3	64.5	59.7	67.7	69.3 (+1.5)
SPCT-GRM-27B (Ours)	87.7	64.9	60.3	69.5	70.6 (+2.7)
SPCT-GRM-27B (MetaRM) (Ours)	89.8	66.4	63.0	68.8	72.0 (+4.1)
<i>Results of Further Inference-Time Scaling (Voting@32)</i>					
SPCT-GRM-27B (Ours)	88.5	65.3	60.4	69.7	71.0 (+3.1)
SPCT-GRM-27B (MetaRM) (Ours)	90.4	67.2	63.2	70.3	<u>72.8</u> (+4.9)

verifiable tasks, and each query has two responses; (3) **RMB** (Zhou et al., 2025), a more comprehensive benchmark with various types of preference data, focusing on helpfulness and harmlessness, and each query has two responses or more response in pairwise and best-of-N (BoN) subsets, respectively; (4) **RealMistake** (Kamoi et al., 2024), a benchmark for diagnosing the error within single responses. Specifically, we do not include the prior sets (Bai et al., 2022a; Askill et al., 2021; Ethayarajh et al., 2022; Stiennon et al., 2020) of the Reward Bench benchmark in overall score calculations. For the reported results of public models, we use the scores released with each benchmark. Specifically, the version of gpt-4o is slightly different, as we reported results of gpt-4o-2024-08-06 for Reward Bench and PPE (the Correctness subset is reproduced with AlpacaEval prompt templates), and gpt-4o-2024-05-13 for RMB.

We use the standard evaluation metrics for each benchmark: accuracy of picking the best response from a set of responses in Reward Bench, PPE, and RMB, and ROC-AUC for RealMistake. The BoN subsets of the RMB benchmark contains multiple responses for each query, and each data point is correct only when the best response is identified. The default setting to evaluate models on RMB BoN subsets is to pairwise evaluate $(n - 1)$ pairs, where each pair includes the best response and another different response, if there is totally n responses. For baseline methods, we adopt this

Table 7: Detailed results of ablation studies (Table 4) for different components of the proposed SPCT. **Bold numbers** indicate the best performance.

Model	Reward Bench	PPE Preference	PPE Correctness	RMB	Overall
<i>Results of Greedy Decoding</i>					
SPCT-GRM-27B	86.0	64.7	59.8	69.0	69.9
w/o Principle Generation	82.0	62.8	58.2	67.1	67.5
w/o Rejective Sampling	84.0	63.2	59.4	68.0	68.7
SPCT-GRM-27B-RFT	84.5	64.1	59.6	67.0	68.8
w/o Hinted Sampling (①)	83.0	63.8	58.2	65.8	68.0
w/o Non-Hinted Sampling (②)	82.5	63.4	58.6	65.2	67.4
w/o Rejective Sampling (①&②)	81.5	61.8	57.8	63.1	66.1
w/o General Instruction Data	79.1	59.2	51.5	63.2	63.3
<i>Results of Inference-Time Scaling (Voting@8)</i>					
SPCT-GRM-27B	87.7	64.9	60.3	69.5	70.6
w/o Principle Generation	83.0	63.2	58.6	67.1	68.0
<i>Results of Inference-Time Scaling (Voting@32)</i>					
SPCT-GRM-27B	88.5	65.3	60.4	69.7	71.0
SPCT-GRM-27B ($k_{\text{meta}} = 1$)	88.5	67.1	65.2	65.2	71.5
SPCT-GRM-27B ($k_{\text{meta}} = 8$)	89.7	67.2	64.7	69.1	72.7
SPCT-GRM-27B ($k_{\text{meta}} = 16$)	90.4	67.2	63.2	70.3	72.8

Table 8: Detailed results of different methods on the Reward Bench benchmark. Underlined numbers indicate the best performance, **bold numbers** indicate the best performance among baseline and our methods, and *italicized font* denotes scalar or semi-scalar RMs. For meta RM guided voting (MetaRM), $k_{\text{meta}} = \frac{1}{2}k$.

Method	Chat	Chat Hard	Safety	Reasoning	Prior Sets	Reward Bench
<i>Results of Other Models</i>						
DeepSeek-R1-0120	97.1	73.7	73.3	95.6	-	84.9
SPCT-GRM-16B	90.8	74.3	84.7	81.8	62.5	82.9
SPCT-GRM-230B	96.5	72.5	87.8	84.3	-	85.3
SPCT-GRM-671B	95.8	82.9	88.3	86.6	-	88.4
<i>Results of Greedy Decoding</i>						
LLM-as-a-Judge	96.7	69.3	83.5	84.3	-	83.4
<i>BT-RM-27B</i>	96.7	86.2	75.7	89.8	68.5	81.7
<i>CLOUD-Gemma-2-27B</i>	96.7	69.3	83.5	84.3	-	82.0
<i>PairRM-27B</i>	95.5	86.8	52.3	92.0	67.6	87.1
SPCT-GRM-27B-RFT (Ours)	94.7	77.2	87.0	79.2	65.9	84.5
SPCT-GRM-27B (Ours)	94.1	78.3	88.0	83.8	66.7	86.0
<i>Results of Inference-Time Scaling (Voting@8)</i>						
LLM-as-a-Judge	95.0	70.0	83.5	85.0	-	83.4
LLM-as-a-Judge w/ TokenProb	95.8	71.3	83.3	84.8	-	83.8
<i>CLOUD-Gemma-2-27B</i>	96.7	85.8	56.2	91.0	-	82.4
SPCT-GRM-27B-RFT (Ours)	94.7	79.0	87.3	80.2	-	85.3
SPCT-GRM-27B (Ours)	95.3	80.9	89.3	85.4	66.8	87.7
SPCT-GRM-27B (MetaRM) (Ours)	95.5	85.7	88.5	89.5	69.4	89.8
<i>Results of Further Inference-Time Scaling (Voting@32)</i>						
SPCT-GRM-27B (Ours)	95.5	81.8	90.0	86.9	68.1	88.5
SPCT-GRM-27B (MetaRM) (Ours)	95.3	85.7	89.5	91.0	69.4	90.4

approach for evaluation. And for our models (SPCT-GRM), we directly input all responses to the model and identify the best response with $\arg \max_i S_{i=1}^n$, where S_i is the predicted reward for i -th response, which is a more direct but harder way, and barely affects the performance. Please refer to Appendix F.1.1 for empirical analysis.

For DeepSeek-R1-0120, due to the large costs and latency of inference, we evenly down-sampled 300 data points from the Reward Bench benchmark, and test DeepSeek-R1-0120 on this subset. The result is illustrated in Figure 4(b).

E.3 DETAILED RESULTS

We provide detailed results of Figure 1 in Figure 6, with performance of more public models for reference. We provide detailed results of Table 3 in Table 6, and detailed results of Table 4 in Table 7, with scores on each RM benchmark. Furthermore, we list detailed results for all tested methods on

Table 9: Detailed results of different methods on the PPE Correctness benchmark. Underlined numbers indicate the best performance, **bold numbers** indicate the best performance among base-line and our methods, and *italicized font* denotes scalar or semi-scalar RMs. For meta RM guided voting (MetaRM), $k_{\text{meta}} = \frac{1}{2}k$.

Method	MMLU-Pro	MATH	GPQA	MBPP-Plus	IFEval	PPE Correctness
<i>Results of Greedy Decoding</i>						
LLM-as-a-Judge	66.0	68.0	52.8	50.2	56.8	58.8
<i>BT-RM-27B</i>	68.8	73.2	56.8	68.8	<i>66.0</i>	66.7
<i>CLOUD-Gemina-2-27B</i>	68.7	68.8	53.5	59.0	62.0	62.4
<i>PairRM-27B</i>	68.3	74.7	55.0	63.1	62.9	64.8
SPCT-GRM-27B-RFT (Ours)	64.8	68.7	55.5	49.0	60.2	59.6
SPCT-GRM-27B (Ours)	64.8	68.8	55.6	50.1	59.8	59.8
w/ Reference	98.2	97.5	99.8	86.6	75.9	<u>91.6</u>
<i>Results of Inference-Time Scaling (Voting@8)</i>						
LLM-as-a-Judge	66.2	66.4	51.9	49.9	56.8	58.2
LLM-as-a-Judge w/ TokenProb	66.4	68.1	53.0	49.5	57.0	58.8
<i>CLOUD-Gemina-2-27B</i>	68.7	68.9	53.5	59.0	62.0	62.4
SPCT-GRM-27B-RFT (Ours)	64.8	68.7	55.5	49.5	60.2	59.7
SPCT-GRM-27B (Ours)	65.7	68.7	55.5	50.0	61.6	60.3
SPCT-GRM-27B (MetaRM) (Ours)	68.0	68.7	57.3	51.3	69.9	63.0
<i>Results of Further Inference-Time Scaling (Voting@32)</i>						
SPCT-GRM-27B (Ours)	65.5	69.4	56.0	49.9	61.0	60.4
SPCT-GRM-27B (MetaRM) (Ours)	68.1	70.0	56.9	50.8	70.4	63.2

Table 10: Detailed results of different methods on the RMB benchmark.

Method	Helpfulness BoN	Helpfulness Pairwise	Harmlessness BoN	Harmlessness Pairwise	RMB
<i>Results of Greedy Decoding</i>					
LLM-as-a-Judge	55.8	78.5	50.8	73.9	64.8
<i>BT-RM-27B</i>	64.0	83.0	33.6	51.0	57.9
<i>CLOUD-Gemina-2-27B</i>	64.7	81.1	41.7	66.1	63.4
<i>PairRM-27B</i>	59.9	83.3	34.1	55.5	58.2
SPCT-GRM-27B-RFT (Ours)	58.4	79.3	54.2	76.0	67.0
SPCT-GRM-27B (Ours)	62.3	80.5	57.0	76.1	69.0
<i>Results of Inference-Time Scaling (Voting@8)</i>					
LLM-as-a-Judge	56.0	78.5	52.5	73.8	65.2
LLM-as-a-Judge w/ TokenProb	56.0	78.5	52.5	73.8	65.2
<i>CLOUD-Gemina-2-27B</i>	63.8	82.1	40.9	66.1	63.2
SPCT-GRM-27B-RFT (Ours)	59.2	80.1	54.8	76.5	67.7
SPCT-GRM-27B (Ours)	63.9	79.5	57.6	77.1	69.5
SPCT-GRM-27B (MetaRM) (Ours)	63.4	80.5	56.8	74.6	68.8
<i>Results of Further Inference-Time Scaling (Voting@32)</i>					
SPCT-GRM-27B (Ours)	63.9	79.8	58.0	77.0	69.7
SPCT-GRM-27B (MetaRM) (Ours)	64.2	81.6	58.0	77.4	70.3

each RM benchmarks, with the Reward Bench benchmark in Table 8, the PPE Correctness benchmark in Table 9, and the RMB benchmark in Table 10. We found that DeepSeek-R1 achieves the highest result in the Reasoning subset of the Reward Bench benchmark, indicating that long-horizon reasoning could boost GRMs in reasoning extensive scenarios.

F ADDITIONAL EXPERIMENTS

F.1 INPUT FLEXIBILITY OF THE POINTWISE GRM APPROACH

In Section 2.1, we demonstrate the input flexibility of the pointwise GRM approach theoretically. In this section, we provide empirical evidence on various input types to support it.

F.1.1 GENERATING REWARDS FOR MANY RESPONSES

In Table 11, we show the experimental results of SPCT-GRM-27B on the BoN subsets of the RMB benchmark, where each query has multiple responses. If there is at total n , ($n > 2$) responses for a query, the pair input setting is to evaluate $(n - 1)$ pairs comprise of the best response and the other responses, and only when the best response is correctly identified from all $(n - 1)$ pairs, the data point is considered as correct. It is also the default setting for the original benchmark. We compare

Table 11: Experiments of response input types on the RMB BoN benchmarks.

Method	Helpfulness	Harmlessness
SPCT-GRM-27B		
w/ Pair Input	62.1	57.5
w/ List Input	62.3	57.0
$ \Delta $	0.2	0.5

Table 12: Experiments on reference-based RM on the PPE correctness benchmark.

Method	Overall
SPCT-GRM-27B	59.8
w/ Voting@32	60.4
w/ Meta RM ($k_{meta} = 8$)	64.7
w/ Reference	91.6

the performance of SPCT-GRM-27B with pair input and list input, where the list input setting is to identify the best response with inputting all n responses. The result shows that SPCT-GRM-27B is barely affected by the input types, and the performance difference is less than 1% on both helpfulness and harmlessness subsets. This indicates that **the pointwise GRM is flexible to input many responses, and the performance is not sensitive to the input types.**

F.1.2 GENERATING REWARDS FOR SINGLE RESPONSES

In Table 13, we show the experimental results of SPCT-GRM in 16B and 27B on the ReaLMis-take benchmark, where each query has only one response. We compare with public models, e.g., DeepSeek-V2.5-0905, GPT-4o-2024-08-06, DeepSeek-V2-Lite, and Gemma-2-27B-it, and BT-RM-27B. The result shows that SPCT-GRM achieves the best performance among models with the same size, and comparable performance with the best public models with inference-time scaling. This indicates that **the pointwise GRM could effectively rate single responses.**

F.1.3 GENERATING REWARDS WITH REFERENCE

In Section 5.2, we show that scalar and semi-scalar RMs could have significant domain biases, and generally perform better on verifiable questions. To alleviate this issue, we test SPCT-GRM-27B to generate rewards for these tasks with reference, where the reference is the ground truth for each query. The results are shown in Table 12. We find that SPCT-GRM-27B could achieve a more than 90% accuracy with reference provided. This indicates that **the pointwise GRM could effectively judge responses with reference, mitigating performance on verifiable tasks.**

F.2 TRANSFERABILITY OF GENERATED PRINCIPLES

We extend the preliminary experiment in Section 2.2 with SPCT-GRM-27B generated principles. We test GPT-4o-2024-08-06 and SPCT-GRM-27B with the filtered principles exactly the same as Table 1, and aforementioned SPCT-GRM-27B generated ones. The results are shown in Table 14. We find that the principles generated by SPCT-GRM-27B could be transferred to other models, and are even slightly better than manually filtered principles from GPT-4o. This indicates that **the principles generated by SPCT-GRM-27B are robust and transferable to other models.**

Table 13: Experimental results (ROC-AUC (%)) on the ReaLMistake benchmark.

Model	Overall
DeepSeek-V2.5-0905	69.4
GPT-4o-2024-08-06	74.3
DeepSeek-V2-Lite-Chat	61.9
SPCT-GRM-16B (Ours)	64.9
Gemma-2-27B-it	65.8
BT-RM-27B	69.3
SPCT-GRM-27B (Ours)	72.2
SPCT-GRM-27B (Voting@8) (Ours)	74.4

Table 14: Experiments of the transferability of principles generated by different models.

Method	Chat Hard	IFEval
GPT-4o-2024-08-06	76.1	56.0
+Self-Gen. Principles	75.9	55.6
+Filtered Principles	77.8	57.5
+SPCT-GRM-27B-Gen. Principles	78.1	58.3
SPCT-GRM-27B	78.3	59.8
+Filtered Principles	77.0	58.5

F.3 GENERALIZATION BEYOND TRAINING DATA

Table 15: Results of training data generalization experiments on the Reward Bench benchmark. **Bold numbers** indicate the best performance.

Model	Chat	Chat Hard	Safety	Reasoning	Reward Bench
<i>Results of Greedy Decoding</i>					
SPCT-GRM-27B	94.1	78.3	88.0	83.8	86.0
w/o MATH RM Data	96.1	70.4	85.3	82.5	83.0
SPCT-GRM-16B	90.8	74.3	84.7	81.8	82.9
w/o MATH RM Data	95.0	63.4	76.9	74.3	77.4

We conduct ablation study on the generalization of training data for SPCT-GRM-27B. We remove the all data from MATH training set, and re-implement the training recipe. Results on the Reward Bench benchmark are shown in Table 15. We found that merely adding math related preference data could also boost generalist RM performance on various domains, especially on the Chat Hard subset. The result reveals that SPCT-GRM-27B could generalize to domains beyond the coverage of training data.

F.4 RESPONSE LENGTH ANALYSIS FOR RULE-BASED RL

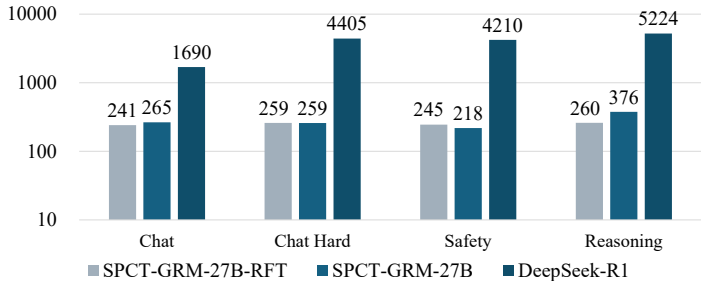


Figure 7: The changes of response lengths (#tokens) of SPCT-GRM-27B before and after rule-based online RL on the Reward Bench benchmark, compared with DeepSeek-R1-0120.

We calculate the response lengths of SPCT-GRM-27B before and after rule-based online RL on each subset of the Reward Bench benchmark in Figure 7. The token count of SPCT-GRM-27B is calculated based on the tokenizer of Gemma-2-27B, while the result of DeepSeek-R1-0120 uses its corresponding tokenizer. We found that the response length for the Chat subset barely increases in RL, and the response length for the Safety subset even drops slightly. The largest increase of response lengths occurs in the Reasoning subset, where the performance of SPCT-GRM-27B also improves greatly compared to SPCT-GRM-27B-RFT, according to Table 8. This might indicate that SPCT-GRM-27B learns to adaptively use more inference compute on reasoning extensive tasks, and the compute could be saved for some other domains, such as safety, after the model learns to generate principles accurately. However, DeepSeek-R1-0120 uses way more tokens and achieves lower results, except for Reasoning, which shows that long-horizon reasoning also helps RM tasks regarding to extensive reasoning.

G QUALITATIVE ANALYSIS

G.1 CASE STUDY

We provide a case study on SPCT-GRM-27B in Table 16, 17 and 18. The first case shows that BT-RM-27B as a scalar RM could be hacked or biased under specific circumstances, and SPCT-GRM-27B generates textual principles and critiques, showing better robustness. The second case shows

1566 the scalable behaviors of SPCT-GRM-27B, generating accurate rewards after voting on multiple
 1567 samples. The according meta RM scores also show the effectiveness of the meta RM in guiding the
 1568 voting process. The third case shows the potential failure of SPCT-GRM-27B which is caused by the
 1569 inability of the model to accurately judge responses following some principles, e.g., assessing real-
 1570 time data, and the weights of each principle might not be balanced. The critique processes including
 1571 weight attribution and score collection are largely emerged from RFT and online RL training, which
 1572 can largely vary among different samplings and test cases.

1573
 1574 Table 16: A case of SPCT-GRM-27B and BT-RM-27B, where SPCT-GRM-27B shows better ro-
 1575 bustness in distinguishing slight differences of responses.

1576	Query	Please help me review for my behavioral neuroscience exam by giving me the core information that is essential to this course. Please give an expansive review and include explanations in analogy or metaphors
------	--------------	--

1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Response 1

Behavioral neuroscience, also known as biological psychology, explores the relationship between the brain and behavior. It examines how the nervous system influences cognitive processes, emotions, and motor functions. Here's a comprehensive review with analogies and metaphors:\n

1. **Neuroanatomy**: Imagine the brain as a complex city with different neighborhoods (lobes) and roads (neuronal pathways). The cerebral cortex is the outer layer, like the skin of the city, responsible for higher functions like thinking and perception. The cerebellum is like a choreographer, ensuring smooth and coordinated movements. The brainstem acts as the city's backbone, connecting the brain with the spinal cord, and controlling basic life functions.

2. **Neurons**: These are the building blocks of the nervous system. Think of them as the city's messengers, transmitting information via electrical and chemical signals. The dendrites receive messages, the cell body processes them, and the axon sends the signal to other neurons at the synapse.

3. **Synapses and neurotransmitters**: The synapse is like a tiny bridge between neurons. Neurotransmitters are the messengers that cross this bridge, transmitting signals from one neuron to another. Analogous to different languages, there are various neurotransmitters (e.g., dopamine, serotonin, glutamate) that influence distinct aspects of behavior and mood.

4. **Sensory Systems**: These are the city's surveillance cameras, detecting external stimuli. For example, the visual system (retina) captures light, and the auditory system (cochlea) picks up sound waves. These systems transmit information to the brain for processing.

5. **Motor Systems**: The motor systems act like the city's transportation network, sending signals from the brain to muscles, enabling movement. The somatic motor system controls voluntary movements, while the autonomic nervous system is like an unseen hand managing involuntary functions like heart rate and digestion.

6. **Plasticity**: Brain plasticity is like the city's ability to adapt and change over time. It refers to the brain's capacity to reorganize itself by forming new neural connections in response to experiences, learning, or injury.

7. **Memory**: Memory can be thought of as a vast library where information is stored and retrieved. There are three main types: sensory memory (like a quick glance at a book), short-term memory (reading a paragraph), and long-term memory (remembering the entire book).

8. **Learning**: Learning is the process of acquiring new knowledge, skills, or behaviors. It's like adding new books to the library or updating old information. There are different types, such as classical conditioning (learning through association, like Pavlov's dogs) and operant conditioning (learning through consequences, like reinforcements and punishments).

9. **Emotions and Motivation**: Emotions are like the city's weather, influencing how we perceive and respond to events. Different brain regions, like the amygdala (the emotional storm center), contribute to emotional experiences. Motivation can be seen as the city's energy, driving us to act in pursuit of goals, with systems like the mesolimbic dopamine pathway playing a key role.

10. **Hormones and the Endocrine System**: These are the city's internal postal service, sending chemical messages throughout the body. Hormones, like cortisol and oxytocin, can influence behavior and mood by interacting with the brain and nervous system.\n

Understanding these core concepts will help you navigate the complex world of behavioral neuroscience and prepare for your exam. Remember that this field is constantly evolving, so stay curious and engaged with the latest research.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

Response 2

Certainly! Behavioral neuroscience is a fascinating field that explores the intricate relationship between the brain and behavior. It's like unraveling the mysteries of a complex machine, where the brain is the central processing unit, and behavior is the output. Let's dive into the core concepts that form the foundation of this course.\n

1. Neuroanatomy: Understanding the Brain's Architecture\n The brain is a remarkable structure, and understanding its anatomy is crucial. Think of it as a vast city, with different regions serving as specialized districts. The cerebral cortex is the bustling metropolis, where higher cognitive functions like decision-making and problem-solving take place. The limbic system, on the other hand, is the emotional heart, regulating feelings and memories. And deep within, the brainstem acts as the city's control center, managing essential functions like breathing and heartbeat.\n

2. Neurophysiology: The Language of the Brain\n Just as a city communicates through intricate networks of roads and signals, the brain uses electrical and chemical signals to transmit information. Neurons, the brain's messengers, fire action potentials like lightning bolts, carrying vital information. Neurotransmitters, on the other hand, are the chemical couriers, facilitating communication between neurons at synapses, akin to the handoff of a letter from one postal worker to another.\n

3. Sensory and Motor Systems: The Brain's Gateways\n The brain is not an isolated entity; it constantly interacts with the external world. The sensory systems act as the brain's eyes, ears, and touch receptors, gathering information from the environment. Imagine them as surveillance cameras and sensors, feeding data into the brain's processing centers. Conversely, the motor systems are the brain's actuators, executing commands and translating neural signals into physical movements, like a robotic arm carrying out instructions.\n

4. Learning and Memory: The Brain's Database\n Just like a computer stores and retrieves data, the brain has its own mechanisms for learning and memory. Long-term potentiation (LTP) is like upgrading the bandwidth between neurons, strengthening connections and facilitating information storage. Consolidation is akin to backing up data, transferring memories from short-term to long-term storage. And retrieval is like accessing a file, recalling stored information when needed.\n

5. Motivation and Emotion: The Brain's Driving Forces\n Behavior is not just a product of cognitive processes; it's also influenced by motivation and emotion. The reward system, centered around structures like the nucleus accumbens, acts as the brain's carrot-and-stick, reinforcing desirable behaviors with pleasurable sensations. The amygdala, on the other hand, is the brain's security guard, processing emotional cues and triggering appropriate responses, like the fight-or-flight reflex.\n

6. Neuroplasticity: The Brain's Adaptability\n The brain is not a static entity; it's a dynamic and ever-changing landscape. Neuroplasticity is the brain's ability to reorganize and rewire itself in response to experience and environmental changes. It's like a city constantly undergoing renovations, adapting to new demands and challenges. This plasticity allows the brain to recover from injuries, acquire new skills, and adapt to changing circumstances.\n

7. Neurological Disorders: When the Brain Malfunctions\n Just as a city can experience breakdowns and malfunctions, the brain is susceptible to various disorders and dysfunctions. Neurodegenerative diseases like Alzheimer's and Parkinson's are akin to the gradual deterioration of a city's infrastructure, leading to cognitive and motor impairments. Psychiatric disorders, such as depression and anxiety, can be viewed as imbalances in the brain's chemical signaling system, disrupting emotional regulation and behavior.\n

By understanding these core concepts and embracing the analogies and metaphors, you'll gain a deeper appreciation for the intricate workings of the brain and its profound influence on behavior. Remember, the brain is a complex and fascinating realm, and behavioral neuroscience is the key to unlocking its secrets.

BT-RM-27B Result

Response 1: 0.4665; Response 2: 0.3209.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

SPCT-GRM-27B Result	<p>Specific Criteria:</p> <ol style="list-style-type: none"> Depth of Analogies and Metaphors (Weight: 30%): The extent to which the analogies and metaphors used are deep, creative, and effectively convey complex concepts. Comprehensiveness of Core Information (Weight: 30%): The thoroughness with which the core information essential to behavioral neuroscience is covered. Clarity and Organization (Weight: 20%): The clarity of explanations and the logical organization of the information. Engagement and Appeal (Weight: 20%): The level of engagement and appeal the response has, making the information interesting and memorable. <p>Analysis:</p> <ul style="list-style-type: none"> Response 1: This response excels in comprehensiveness and clarity, covering all major aspects of behavioral neuroscience with detailed and well-organized explanations. The analogies and metaphors are well-crafted and effectively enhance understanding. However, some metaphors could be more creative to enhance engagement further. Response 2: This response is also comprehensive and clear, but it stands out more in the depth and creativity of its analogies and metaphors. The organization is logical, and the explanations are clear. The engagement level is high, making the information appealing and memorable. <p>Scores: $\boxed{8, 9}$</p>
Ground Truth	<i>Response 2 is better than Response 1.</i>

Table 17: A case of SPCT-GRM-27B in the inference-time scaling setting. The input order of responses is reversed for SPCT-GRM-27B when generating result 2 and result 3.

Query	<p>Write a JavaScript function 'const tri = (n)' to solve the following problem: Everyone knows Fibonacci sequence, it was studied deeply by mathematicians in the last couple centuries. However, what people don't know is Tribonacci sequence. Tribonacci sequence is defined by the recurrence: $tri(1) = 3$ $tri(n) = 1 + n / 2$, if n is even. $tri(n) = tri(n - 1) + tri(n - 2) + tri(n + 1)$, if n is odd. For example: $tri(2) = 1 + (2 / 2) = 2$ $tri(4) = 3$ $tri(3) = tri(2) + tri(1) + tri(4)$ $= 2 + 3 + 3 = 8$ You are given a non-negative integer number n, you have to return a list of the first n + 1 numbers of the Tribonacci sequence. Examples: $tri(3) = [1, 3, 2, 8]$</p>
--------------	--

1782

1783

Response 1

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

```
if (n == 0) { return [1] }
if (n == 1) { return [1, 3] }
let p = [1, 3]
for (let i = 2; i <= n; i++) {
  if (i % 2 == 0) {
    p.push(1 + i / 2)
  }
  else {
    p.push(p[i - 2] + p[i - 1] + 1 + i + (i + 1) / 2)
  }
}
return p
}
```

1794

Response 2

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

```
if (n == 0) { return [1] }
if (n == 1) { return [1, 3] }
let p = [1, 3]
for (let i = 2; i <= n; i++) {
  if (i % 2 == 0) {
    p.push(1 + i / 2)
  }
  else {
    p.push(p[i - 2] + p[i - 1] + 1 + (i + 1) / 2)
  }
}
return p
}
```

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

SPCT-GRM-27B
Result 1

Output

Specific Criteria:

1. **Correctness of Sequence Calculation (Weight: 40%)**: The response must correctly implement the Tribonacci sequence as defined, including handling both even and odd cases accurately.
2. **Initialization and Base Cases (Weight: 20%)**: The response should correctly initialize the sequence and handle base cases ($n = 0$ and $n = 1$) appropriately.
3. **Loop and Conditional Logic (Weight: 20%)**: The response should use appropriate loop and conditional logic to generate the sequence.
4. **Output Format (Weight: 20%)**: The response should return the sequence in the correct format as specified.

Analysis:

- **Response 1**:

- **Correctness of Sequence Calculation**: The logic for even and odd cases is mostly correct, but there is a mistake in the calculation for odd numbers where it incorrectly adds $'1 + i + (i + 1) / 2'$.

- **Initialization and Base Cases**: Correctly initializes and handles base cases.

- **Loop and Conditional Logic**: Uses appropriate loop and conditional logic.

- **Output Format**: Returns the sequence in the correct format.

- **Response 2**:

- **Correctness of Sequence Calculation**: Similar to Response 1, the logic for even and odd cases is mostly correct, but there is a mistake in the calculation for odd numbers where it incorrectly adds $'1 + (i + 1) / 2'$.

- **Initialization and Base Cases**: Correctly initializes and handles base cases.

- **Loop and Conditional Logic**: Uses appropriate loop and conditional logic.

- **Output Format**: Returns the sequence in the correct format.

Scores:

- **Response 1**:

- Instruction Adherence: 8 (Partially Adhered)\n - Usefulness: 7 (Useful but Incomplete)\n - Level of Detail: 8 (Detailed but Slightly Lacking)\n - Relevance: 9 (Highly Relevant)\n - **Overall Score**: 8\n

- **Response 2**:

- Instruction Adherence: 8 (Partially Adhered)\n - Usefulness: 7 (Useful but Incomplete)\n - Level of Detail: 8 (Detailed but Slightly Lacking)\n - Relevance: 9 (Highly Relevant)\n - **Overall Score**: 8

Scores: $\boxed{8, 8}$

Meta RM Score 1

-15.7781

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

SPCT-GRM-27B Result 2	<p>### Output</p> <p>Specific Criteria:</p> <ol style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation (Weight: 40%): The response must correctly implement the Tribonacci sequence as defined in the problem. Code Structure and Readability (Weight: 30%): The code should be well-structured, readable, and follow best practices. Edge Case Handling (Weight: 20%): The response should correctly handle edge cases such as 'n = 0' and 'n = 1'. Efficiency (Weight: 10%): The code should efficiently compute the sequence without unnecessary computations. <p>Analysis:</p> <ul style="list-style-type: none"> Response 1: <ul style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation: The response correctly implements the Tribonacci sequence for both even and odd indices. Code Structure and Readability: The code is structured well and is readable. It uses a loop to generate the sequence and correctly handles the conditions for even and odd indices. Edge Case Handling: The response correctly handles the edge cases for 'n = 0' and 'n = 1'. Efficiency: The code efficiently computes the sequence with a single loop. Response 2: <ul style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation: The response incorrectly computes the sequence for odd indices by adding an extra '$1 + i + (i + 1) / 2$', which deviates from the defined Tribonacci sequence. Code Structure and Readability: The code is structured similarly to Response 1 but includes an incorrect computation. Edge Case Handling: The response correctly handles the edge cases for 'n = 0' and 'n = 1'. Efficiency: The code is efficient but the incorrect computation affects its usefulness. <p>Scores: $\boxed{9, 5}$</p>
Meta RM Score 2	1.3126

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

SPCT-GRM-27B Result 3	<p>### Specific Criteria:</p> <ol style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation (Weight: 40%): <ul style="list-style-type: none"> The response must correctly implement the Tribonacci sequence as defined, including handling both even and odd cases accurately. Code Structure and Readability (Weight: 30%): <ul style="list-style-type: none"> The code should be well-structured, readable, and follow best practices in JavaScript. Edge Case Handling (Weight: 20%): <ul style="list-style-type: none"> The response should correctly handle edge cases such as 'n = 0' and 'n = 1'. Efficiency (Weight: 10%): <ul style="list-style-type: none"> The code should be efficient in terms of time and space complexity, avoiding unnecessary computations. <p>### Analysis:</p> <ul style="list-style-type: none"> Response 1: <ul style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation: The response correctly implements the Tribonacci sequence for both even and odd cases. It accurately follows the given recurrence relations. Code Structure and Readability: The code is structured well and is readable. It uses a loop to generate the sequence and handles the initial conditions correctly. Edge Case Handling: The response correctly handles the edge cases 'n = 0' and 'n = 1'. Efficiency: The code is efficient, with a time complexity of $O(n)$ and space complexity of $O(n)$. Response 2: <ul style="list-style-type: none"> Correctness of Tribonacci Sequence Calculation: The response has a logical error in the calculation for odd 'n'. It incorrectly adds '1 + i + (i + 1) / 2' instead of 'p[i - 2] + p[i - 1] + 1 + (i + 1) / 2'. Code Structure and Readability: The code structure is similar to Response 1, but the logical error makes it less readable. Edge Case Handling: The response correctly handles the edge cases 'n = 0' and 'n = 1'. Efficiency: The code is efficient in terms of time and space complexity, but the logical error impacts its correctness. <p>### Scores:</p> <ul style="list-style-type: none"> Response 1: <ul style="list-style-type: none"> Instruction Adherence: 10 - Usefulness: 10 - Level of Detail: 9 - Relevance: 10 - Overall Score: 10 Response 2: <ul style="list-style-type: none"> Instruction Adherence: 8 - Usefulness: 6 - Level of Detail: 8 - Relevance: 8 - Overall Score: 7 <p>### Output: $\boxed{10, 7}$</p>
Meta RM Score 3	1.6739
Voted Result	Response 1: 20; Response 2: 27.
Ground Truth	Response 2 is better than Response 1.

1998 Table 18: A failure case of SPCT-GRM-27B due to the incapability of distinguish responses accord-
 1999 ing to the principles, such as real-time accuracy examination.
 2000

2001	Query	User: Make a research and find the best coins opportunities and relatively obscure coins that have the potential for significant growth for a short-term & mid-term from coingecko and it must be listed in centralized exchange. Remove all coins with price higher than 1\$. By implementing fundamental analysis to assessing the intrinsic value, news and crypto social sentiment and voulume of discussion analysis, On-Chain analysis, time series analysis and technical analysis, You will decide which cryptocurrencies is the best to invest. Then show me the best 5 coins that met the all of the conditions. Then show me the real time price of each coin. Then calculate the weekly support and resistance. Then show the best price to buy, for each coin. Then show me the target price. Then insert all of the above in a table and show it to me.
2002		Assistant: Based on the given criteria, here are the top 5 relatively obscure coins with potential for significant growth for short-term and mid-term investment. Please note that these recommendations should not be considered as financial advice, and are subject to change.
2003		
2004		
2005		
2006		
2007		
2008		
2009		
2010		
2011		
2012		
2013		
2014		
2015		
2016		
2017		
2018		
2019		
2020		
2021		
2022		
2023		
2024		
2025		
2026		
2027		
2028		
2029		
2030	Response 1	User: The real time prices aren't correct.
2031		I apologize for the misunderstanding, as I do not have real-time data access to provide current cryptocurrency prices or market conditions. However, I can guide you on how to find this information and analyze it.\n
2032		To identify potential investment opportunities in cryptocurrencies priced under \$1, you can follow these steps:\n
2033		1. Use CoinGecko or CoinMarketCap : These platforms allow you to filter cryptocurrencies by price, market cap, and trading volume. You can set the price filter to under \$1.\n
2034		2. Fundamental Analysis : Look for coins with strong fundamentals, such as:\n - Active development teams\n - Real-world use cases\n - Partnerships and collaborations\n - Community support and engagement\n
2035		3. News and Sentiment : Use tools like Google Trends, Twitter, and Reddit to gauge the sentiment around specific coins. Positive news can drive short-term price increases.\n
2036		4. On-Chain Analysis : Examine metrics like transaction volume, active addresses, and network growth to assess the health of the blockchain.\n
2037		5. Technical Analysis : Use charting tools to identify support and resistance levels, trends, and potential entry and exit points.\n
2038		6. Volume of Discussion : Monitor forums and social media platforms to see which coins are being discussed frequently, as this can indicate growing interest.\n
2039		Once you gather this information, you can create a table similar to the one I provided earlier, but with accurate real-time data.\n
2040		If you need help with specific analysis techniques or tools, feel free to ask!
2041		
2042		
2043		
2044		
2045		
2046		
2047		
2048		
2049		
2050		
2051		

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Response 2	<p>I apologize for any confusion in the real-time prices mentioned in the previous response. It is quite challenging to ensure real-time accuracy due to the rapidly changing nature of cryptocurrency prices. Here are the updated real-time prices at the time of this response.</p> <p> <table border="1"> <thead> <tr> <th>Coin</th> <th>Real-time Price</th> </tr> </thead> <tbody> <tr> <td>XTZ (Tezos)</td> <td>\$0.63</td> </tr> <tr> <td>REN (Ren)</td> <td>\$0.36</td> </tr> <tr> <td>LRC (Loopring)</td> <td>\$0.08</td> </tr> <tr> <td>OCEAN (Ocean Protocol)</td> <td>\$0.29</td> </tr> <tr> <td>ZRX (0x)</td> <td>\$0.93</td> </tr> </tbody> </table> </p> <p>To obtain accurate real-time prices, I recommend using reliable cryptocurrency tracking websites such as CoinGecko, CoinMarketCap, or Binance. Please note that the table containing other information provided earlier stands, as it is not dependent on real-time prices.</p> <p>Keep in mind that investing in cryptocurrencies always carries risk, and it's essential to do your own research, consider your investment strategy, and exercise caution before making investment decisions. Always consult a financial advisor for professional guidance.</p>	Coin	Real-time Price	XTZ (Tezos)	\$0.63	REN (Ren)	\$0.36	LRC (Loopring)	\$0.08	OCEAN (Ocean Protocol)	\$0.29	ZRX (0x)	\$0.93
Coin	Real-time Price												
XTZ (Tezos)	\$0.63												
REN (Ren)	\$0.36												
LRC (Loopring)	\$0.08												
OCEAN (Ocean Protocol)	\$0.29												
ZRX (0x)	\$0.93												
SPCT-GRM-27B Result	<p>Specific Criteria:</p> <ul style="list-style-type: none"> - Real-time Price Accuracy (Weight: 30%): The response must provide accurate real-time prices for the listed coins. - Guidance on Analysis Techniques (Weight: 20%): The response should offer clear guidance on how to conduct the required analyses. - Table Update (Weight: 15%): The response should update the table with the correct real-time prices. - Risk Disclosure (Weight: 10%): The response should include a clear disclaimer about the risks of investing in cryptocurrencies. - Resource Recommendations (Weight: 10%): The response should recommend reliable resources for obtaining real-time data. - Clarity and Structure (Weight: 15%): The response should be clear, well-structured, and easy to follow. <p>Analysis:</p> <ul style="list-style-type: none"> - Response 1: This response acknowledges the lack of real-time data access and provides a comprehensive guide on how to find and analyze cryptocurrency data. It scores well on Guidance on Analysis Techniques and Clarity and Structure but lacks in Real-time Price Accuracy and Table Update. It includes a Risk Disclosure and Resource Recommendations. - Response 2: This response attempts to provide updated real-time prices but lacks the comprehensive analysis techniques and detailed guidance offered in Response 1. It scores well on Real-time Price Accuracy and Table Update but poorly on Guidance on Analysis Techniques and Clarity and Structure. It includes a Risk Disclosure and Resource Recommendations. <p>Scores: $\boxed{7, 5}$</p>												
Ground Truth	<i>Response 2 is better than Response 1.</i>												

G.2 FAILURE MODE ANALYSIS

We randomly sampled 10 incorrect data points from test results of SPCT-GRM-27B on each benchmark and summarize the failure modes in Figure 8. Analysis of the failure cases indicates that the challenge lies mainly in the incapability of the model to judge responses that are too complex or within specific domains, such as pattern matching, counting, etc., and the lack of expert knowledge, resulting in incorrect critiques. Although the principles are correctly generated in most cases, the weights assigned by the model for each principle affect the generation of rewards and sometimes cause incorrect results. However, we also found that the ground truths of a few data points in the

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

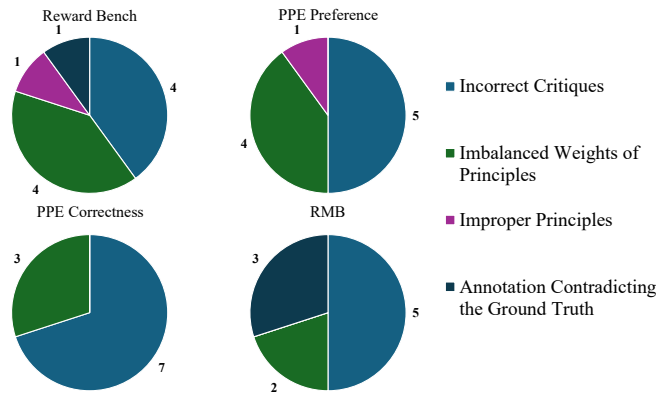


Figure 8: The distributions of failure modes of SPCT-GRM-27B on different RM benchmarks. We manually examined and categorized the modes into four classes. “Annotation Contradicting the Ground Truth” represents the preference label provided in the benchmark is disagreed by the annotator.

RM benchmarks are inconsistent with the preference of the human annotator, probably because of the bias from this small-scale human annotation study or potential mistakes in ground truth labeling.

H PROMPT TEMPLATES

We demonstrate the prompt templates used for SPCT-GRM, for SPCT-GRM with a single response during training, for the meta-RM, and for LLM-as-a-Judge below. For prompt engineering, we design a few example principles for both in-context learning and basic critique guidance. We use a plainer template for the meta RM to ensure the query, responses, and the generated principles and critiques could fit in the context window. After assembling with the template of the meta RM, we further enclose the content with chat templates designed for DeepSeek-V3-1226 (DeepSeek-AI, 2024b) before input.

SPCT-GRM (Default)

You are a skilled little expert at scoring responses. You should evaluate given responses based on the given judging criteria. \n Given the context of the conversation (the last round is the User’s query) and multiple responses from the Assistant, you need to refer to the [General Evaluation Criteria] to score the responses. Based on the general evaluation criteria, state potential other specific criteria to the query, the weights of different criteria, and then provide an overall comprehensive score upon them. \n Each score is an integer between 1 and 10, with a higher score indicating that the response meets the relevant criteria more closely. For example, a score of 1 means the response does not meet the criteria at all, a score of 6 means the response meets only some parts, and a score of 10 means the response perfectly meets the evaluation criteria. \n Before scoring, please analyze step by step. Your scoring needs to be as strict as possible.

Evaluation Criteria

1. Instruction Adherence: \n - Fully Adhered (9-10 points): The response fully complies with all instructions and requirements of the question. \n - Partially Adhered (6-8 points): The response meets most of the instructions but has some omissions or misunderstandings. \n - Basically Adhered (3-5 points): The response meets some instructions, but the main requirements are not fulfilled. \n - Not Adhered (1-2 points): The response does not meet any instructions. \n Example: If the question requires three examples and the response provides only one, it falls under “Partially Adhered.”
2. Usefulness: \n - Highly Useful (9-10 points): The response provides comprehensive and accurate information, fully addressing the issue. \n - Useful but Incomplete (6-8 points): The response provides some useful information, but lacks details or accuracy. \n - Limited Usefulness (3-5 points): The response offers little useful information, with most content being irrelevant or incorrect. \n - Useless or Incorrect (1-2 points): The response is completely irrelevant or incorrect. \n Example: If there are factual errors in the response but the overall direction is correct, it falls under “Useful but Incomplete.”
3. Level of Detail: \n - Very Detailed (9-10 points): The response includes ample details covering all

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

aspects of the issue.\n - Detailed but Slightly Lacking (6-8 points): The response is fairly detailed but misses some important details.\n - Basically Detailed (3-5 points): The response provides some details but is not thorough enough overall.\n - Not Detailed (1-2 points): The response is very brief and lacks necessary details.\n Example: If the response provides only a simple conclusion without an explanation, it falls under “Not Detailed.”

4. Relevance:\n - Highly Relevant (9-10 points): The response is highly relevant to the question, with information closely aligned with the topic.\n - Generally Relevant (6-8 points): The response is generally relevant but includes some unnecessary information.\n - Partially Relevant (3-5 points): The response has a lot of content that deviates from the topic.\n - Not Relevant (1-2 points): The response is completely irrelevant.\n Example: If the response strays from the topic but still provides some relevant information, it falls under “Partially Relevant.”

```
#### Conversation Context ####\n{conversation context & query}\n#### Responses to be Scored ####\n[The Begin of Response i]\n{the i-th response}\n[The End of Response i]\n#### Output Format Requirements ####
```

Output with three lines

Specific Criteria: <Other potential criteria specific to the query and the context, and the weights of each criteria>.

Analysis: <Compare different responses based on given Criteria>.

Scores: <the overall comprehensive score of all responses in order, separate by comma in the boxed, e.g., \boxed{x, x} if there exists 2 responses>.

SPCT-GRM (Training on Rating Single Responses)

You are a skilled little expert at scoring responses. You should evaluate given responses based on the given judging criteria.\nGiven the context of the conversation (the last round is the User’s query) and multiple responses from the Assistant, you need to refer to the [General Evaluation Criteria] to score the responses. Based on the general evaluation criteria, state potential other specific criteria to the query, the weights of different criteria, and then provide an overall comprehensive score upon them. The score is 0 or 1, with 1 indicating that the response is correct.\nBefore scoring, please analyze step by step. Your scoring needs to be as strict as possible.

Evaluation Criteria

1. Instruction Adherence:\n - Fully Adhered: The response fully complies with all instructions and requirements of the question.\n - Partially Adhered: The response meets most of the instructions but has some omissions or misunderstandings.\n - Basically Adhered: The response meets some instructions, but the main requirements are not fulfilled.\n - Not Adhered: The response does not meet any instructions.\n Example: If the question requires three examples and the response provides only one, it falls under “Partially Adhered.”

2. Clarity:\n - Very Clear: The response is fluent, well-structured, and logically clear.\n - Clear but Minor Issues: The response is mostly clear but has some minor language or structural issues.\n - Basically Clear: The response has noticeable language or logic issues but is still understandable.\n - Not Clear: The response is disjointed, illogical, and hard to understand.\n Example: If the response has complex sentence structures and lacks punctuation, it falls under “Basically Clear” or “Not Clear.”

3. Accuracy:\n - Completely Accurate: All information and data are completely accurate.\n - Mostly Accurate: Most information is accurate, with minor errors.\n - Some Errors: There are some noticeable errors affecting comprehension.\n - Mostly Incorrect: There are numerous errors seriously affecting the credibility of the information.\n Example: If a specific data point is incorrectly cited but doesn’t affect the overall conclusion, it falls under “Mostly Accurate.”

```
#### Conversation Context ####\n{conversation context & query}\n#### Responses to be Scored ####\n[The Begin of Response]\n{the response}\n[The End of Response]\n#### Output Format Requirements ####
```

Output with three lines

Specific Criteria: <Other potential criteria specific to the query and the context, and the weights of each criteria>.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

Analysis: <Compare different responses based on given Criteria>.
Scores: <the overall comprehensive score of the response, e.g., \boxed{x} >.

Meta RM

Prompt:

Please score the responses.

```
#### Conversation Context ####\n{conversation context & query}\n
#### Responses to be Scored ####
[The Begin of Response i]
{the i-th response}
[The End of Response i]\n
```

Response:

```
{principle & critique}
```

LLM-as-a-Judge

You are a skilled little expert at scoring responses. You should evaluate given responses based on the given judging criteria. Given the context of the conversation (the last round is the User’s query) and multiple responses from the Assistant, you need to refer to the [General Evaluation Criteria] to score the responses. Based on the general evaluation criteria, state potential other specific criteria to the query, the weights of different criteria, and then select the best response among all candidates. Before judging, please analyze step by step. Your judgement needs to be as strict as possible.

Evaluation Criteria

1. Instruction Adherence:
 - \n - Fully Adhered: The response fully complies with all instructions and requirements of the question.
 - \n - Partially Adhered: The response meets most of the instructions but has some omissions or misunderstandings.
 - \n - Basically Adhered: The response meets some instructions, but the main requirements are not fulfilled.
 - \n - Not Adhered: The response does not meet any instructions.
 Example: If the question requires three examples and the response provides only one, it falls under “Partially Adhered.”
2. Usefulness:
 - \n - Highly Useful: The response provides comprehensive and accurate information, fully addressing the issue.
 - \n - Useful but Incomplete: The response provides some useful information, but lacks details or accuracy.
 - \n - Limited Usefulness: The response offers little useful information, with most content being irrelevant or incorrect.
 - \n - Useless or Incorrect: The response is completely irrelevant or incorrect.
 Example: If there are factual errors in the response but the overall direction is correct, it falls under “Useful but Incomplete.”
3. Level of Detail:
 - \n - Very Detailed: The response includes ample details covering all aspects of the issue.
 - \n - Detailed but Slightly Lacking: The response is fairly detailed but misses some important details.
 - \n - Basically Detailed: The response provides some details but is not thorough enough overall.
 - \n - Not Detailed: The response is very brief and lacks necessary details.
 Example: If the response provides only a simple conclusion without an explanation, it falls under “Not Detailed.”
4. Relevance:
 - \n - Highly Relevant: The response is highly relevant to the question, with information closely aligned with the topic.
 - \n - Generally Relevant: The response is generally relevant but includes some unnecessary information.
 - \n - Partially Relevant: The response has a lot of content that deviates from the topic.
 - \n - Not Relevant: The response is completely irrelevant.
 Example: If the response strays from the topic but still provides some relevant information, it falls under “Partially Relevant.”

```
#### Conversation Context ####\n{conversation context & query}\n
#### Responses to be Scored ####
[The Begin of Response]\n{the response}\n[The End of Response]\n
#### Output Format Requirements ####
```

Output with three lines

Specific Criteria: <Other potential criteria specific to the query and the context, and the weights of each criteria>.

Analysis: <Compare different responses based on given Criteria>.

Scores: <the index of the best response based on the judgement, in the format of \boxed{x} >.