LEARNING ORDINAL PROBABILISTIC REWARD FROM PREFERENCES

Anonymous authorsPaper under double-blind review

ABSTRACT

Reward models are crucial for aligning large language models (LLMs) with human values and intentions. Existing approaches follow either Generative (GRMs) or Discriminative (DRMs) paradigms, yet both suffer from limitations: GRMs typically demand costly point-wise supervision, while DRMs produce uncalibrated relative scores that lack probabilistic interpretation. To address these challenges, we introduce a novel reward modeling paradigm: *Probabilistic Reward Model* (PRM). Instead of modeling reward as a deterministic scalar, our approach treats it as a random variable, learning a full probability distribution for the quality of each response. To make this paradigm practical, we present its closed-form, discrete realization: the *Ordinal Probabilistic Reward Model* (OPRM), which discretizes the quality score into a finite set of ordinal ratings. Building on OPRM, we propose a data-efficient training strategy called *Region Flooding Tuning* (RgFT). It enables rewards to better reflect absolute text quality by incorporating qualitylevel annotations, which guide the model to concentrate the probability mass within corresponding rating sub-regions. Experiments on various reward model benchmarks show that our method improves accuracy by $2.9\% \sim 7.4\%$ compared to prior reward models, demonstrating strong performance and data efficiency. Analysis of the score distribution provides evidence that our method captures not only relative rankings but also absolute quality. Our models, data, and code will be released and open-sourced.

1 Introduction

Reinforcement Learning from Human Feedback (RLHF) has emerged as a pivotal technique for aligning Large Language Models (LLMs) with human values and intentions (Achiam et al., 2023; Ouyang et al., 2022). As a critical component of the RLHF process (Bai et al., 2022), the reward model is trained to assign scores that quantify the degree of alignment between the model's outputs and human preferences. Recent advances (Guo et al., 2025a; Lightman et al., 2024) have shown that well-designed reward signals, whether applied during training or inference, can significantly enhance LLM performance across diverse domains (Shao et al., 2024; Huang et al., 2025; Jin et al., 2025). However, learning a reward model that can accurately capture human preference signals remains a significant challenge (Gao et al., 2023; Sun et al., 2025a; Zhong et al., 2025). Most recent efforts typically follow either the generative or discriminative paradigm, yet both approaches exhibit inherent limitations that hinder their effectiveness in practice.

Discriminative Reward Models (DRMs), which append an MLP-based value head to a base model, are commonly optimized with the Bradley-Terry objective to output a scalar reward (Liu et al., 2024b; Cai et al., 2024; Lou et al., 2024). A key limitation of this paradigm is that its reward scores reflect only relative preferences, not intrinsic quality. It indicates that one response is preferred but fails to explain why, making it difficult to establish a trusted acceptance threshold to discern high-quality responses from low-quality ones. In response, Generative Reward Models (GRMs) have emerged (Mahan et al., 2024; Zhang et al., 2024). These models leverage the native generative capabilities of LLMs to produce Chain-of-Thought critiques before rendering a preference judgment, conceptually aligning with the LLM-as-a-Judge paradigm (Zheng et al., 2023). While GRMs offer superior interpretability through their critique generation, they often rely on rigid pairwise input formats that limit flexibility in Best-of-N (BoN) scenarios. Moreover, achieving performance comparable to DRMs frequently requires costly pointwise supervision for calibration, substantially increasing the annotation burden.

Table 1: Comparison of OPRM with baseline reward models across multiple dimensions. Margin Sensitivity (whether distinguish samples with subtle preference differences), Require Training (whether requires training on preference data), Value Head Free (whether eliminates the need for additional value head), and Input Flexibility (whether supports rating single and multiple responses).

Baselines	Input Format	Output Format	Margin Sensitivity	Require Training	Value Head Free	Input Flexible
Bradley-Terry (Bradley & Terry, 1952)	Single Response	Continuous Score	/	/	X	✓
PairRM (Jiang et al., 2023)	Response Pairs	Continuous Score	✓	1	X	X
CLoud (Ankner et al., 2024)	Single Response	Critique + Continuous Score	✓	1	X	/
LLM-as-a-Judge (Zheng et al., 2023)	Response Pairs	Discrete Score	X	X	✓	X
Pointwise GRM (Liu et al., 2025)	Single Response	Critique + Discrete Score	X	1	✓	/
OPRM (Ours)	Single Response	Continuous Score	√	/	√	✓

Consequently, the field faces a critical trade-off: choosing between the efficiency of DRMs and the interpretability of GRMs, with neither approach offering a complete solution.

To transcend this trade-off, we introduce a novel reward modeling paradigm: *Probabilistic Reward Model* (PRM). Instead of approximating rewards with a deterministic scalar value like the Bradley-Terry model (Bradley & Terry, 1952), PRM reframes the task as learning a full probability distribution over the reward space. Since learning this continuous distribution is computationally intractable, we translate it into a discrete realization: *Ordinal Probabilistic Reward Model* (OPRM). Specifically, OPRM discretizes the reward space into a finite set of ordinal ratings, thereby replacing the intractable integration with a closed-form summation that makes our paradigm more practical. Thus, OPRM resolves the core trade-off. By providing a full reward distribution, it unlocks richer interpretability and uncertainty estimation than DRMs, while its flexible input format enables efficient scoring of single or multiple responses, making it better suited than GRMs for modern evaluations like Best-of-N (BoN). Table 1 summarizes the advantages of OPRM over existing reward modeling baselines.

Building upon the OPRM paradigm, we further propose *Region Flooding Tuning* (RgFT), a novel training strategy designed to calibrate the reward distribution to reflect absolute textual quality. The core principle of RgFT is to leverage quality-level annotations (i.e., **good**, **normal**, and **bad**) on preference data. Rather than optimizing over the full distribution of ordinal ratings, RgFT guides the model to concentrate probability mass within rating sub-regions corresponding to the quality-level labels. Critically, RgFT facilitates semi-supervised training by jointly leveraging a mixture of quality-labeled and preference-only data, obviating the need for costly large-scale annotation.

Empirically, OPRM achieves state-of-the-art performance on RewardBench (Lambert et al., 2024), PPE (Frick et al., 2024), RMB (Zhou et al., 2024), outperforming 70B, 340B, GPT-40, and Claude models by up to **2.9**%~**7.4**%. Moreover, when trained on domain-specific data featuring quality-level annotations, RgFT further boosts the performance of OPRM. Further analysis demonstrates the effectiveness of OPRM and RgFT across a variety of experimental settings.

Before delving into details, we summarize our contributions as follows:

- We propose a novel reward modeling paradigm, the Ordinal Probabilistic Reward Model. By learning a full probability distribution for a response's quality, OPRM mitigates the core trade-off between the efficiency of DRMs and the interpretability of GRMs.
- We design a data-efficient training strategy, Region Flooding Tuning, which grounds the reward distribution in an absolute quality scale by guiding the model to concentrate probability mass within correct rating sub-regions.
- We conduct extensive experiments on four benchmarks covering over ten tasks, demonstrating the
 effectiveness of OPRM in precise reward modeling across diverse scenarios. Additional studies
 confirm that RgFT significantly improves the accuracy, robustness, and interpretability of OPRM.

2 RELATED WORK

Discriminative Reward Model. Discriminative reward model typically consists of a base model and a MLP-based reward head (classifier) that outputs a scalar score for a given input. These models are commonly trained using the Bradley-Terry (BT) (Bradley & Terry, 1952) loss to maximize the reward margin between chosen and rejected responses. While the core BT loss remains a standard component, considerable research has focused on enhancing data quality and refining the modeling

framework (e.g. Skywork-reward (Liu et al., 2024b), InternLM2-reward (Cai et al., 2024), Helpsteer2-preference (Wang et al., 2024b), QRM (Dorka, 2024), URM (Lou et al., 2024), CLoud (Ankner et al., 2024), ArmoRM (Wang et al., 2024a), and PURM (Sun et al., 2025b)), further boosting DRM performance. Nonetheless, these methods are limited to learning a pairwise ranking, yielding scores that are unbounded and difficult to interpret. In contrast, our approach learns a probabilistic distribution over scores, which enables more reliable and calibrated outputs.

Generative Reward Model. Unlike discriminative models, generative reward models directly leverage LLM-generated outputs to evaluate preference data, which is conceptually aligned with the LLM-as-a-Judge paradigm (Zheng et al., 2023). These models output chain-of-thought (CoT) reasoning (critiques) before generating preference judgments or rewards (e.g., Critic-RM (Yu et al., 2024), PROMETHEUS (Kim et al., 2023), CLoud (Ankner et al., 2024), Generative Verifiers (Zhang et al., 2024), GenRM (Mahan et al., 2024), Synthetic Critique (Ye et al., 2024), and RISE (Yu et al., 2025)), enhancing the interpretability of the reward signals. The final reward can be derived either from training a separate reward head or from the LM head itself, which generates real-valued scores as tokens. Recent advances have employed reinforcement learning to construct reasoning-based reward models (DeepSeek-GRM (Liu et al., 2025), RM-R1 (Chen et al., 2025b), J1 (Whitehouse et al., 2025), RRM (Guo et al., 2025b), and JudgeLRM (Chen et al., 2025a)), demonstrating promising scalability in inference-time computation. However, these approaches face efficiency challenges and struggle to outperform DRMs under equivalent computational constraints. Conversely, our method maintains interpretability while matching the computational efficiency and performance of DRMs.

3 PRELIMINARIES

Preference data annotation. To annotate the preference data, the SFT model $\pi^{\rm SFT}$ is given prompts x to two distinct outputs $(y_1,y_2) \sim \pi^{\rm SFT}(y\mid x)$. These output pairs are then presented to human labelers, who express their preference for one output. This preference can be denoted as $y_{\rm c} \succ y_{\rm r} \mid x$, where $y_{\rm c}$ and $y_{\rm r}$ represent the chosen and rejected outputs, respectively, from the pair (y_1,y_2) .

Standard Bradley-Terry Reward Modeling. Following the Bradley-Terry model (Bradley & Terry, 1952), we model the probability of preferring response y_c over y_r based on their underlying scalar rewards, which are provided by a reward function $r_{\psi}(x,y)$. This preference distribution is formulated as follows:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \frac{\exp(r_{\psi}(x, y_{c}))}{\exp(r_{\psi}(x, y_{c})) + \exp(r_{\psi}(x, y_{r}))},$$

$$= \sigma(r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r})),$$
(1)

which σ is the logistic function. Treating the problem as a binary classification task yields the negative log-likelihood loss function:

$$\mathcal{L}(r_{\psi}) = -\mathbb{E}_{(x,y_{c},y_{r})\sim\mathcal{D}_{rm}}[\log P_{\psi}(y_{c} \succ y_{r} \mid x)],$$

$$= -\mathbb{E}_{(x,y_{c},y_{r})\sim\mathcal{D}_{rm}}[\log \sigma(r_{\psi}(x,y_{c}) - r_{\psi}(x,y_{r}))],$$
(2)

where dataset is composed of comparisons denoted as $\mathcal{D}_{\mathrm{rm}} = \{x^{(i)}, y_{\mathrm{c}}^{(i)}, y_{\mathrm{r}}^{(i)}\}_{i=1}^{N}$. In the realm of LLMs, the network $r_{\psi}(x,y)$ is often initialized using the SFT model $\pi^{\mathrm{SFT}}(y\mid x)$. It then incorporates an additional linear layer on the final transformer layer to generate a singular scalar prediction representing the reward value.

4 ORDINAL PROBABILISTIC REWARD MODEL

In this section, we introduce **Ordinal Probabilistic Reward Model**, a novel reward modeling paradigm that learns a probability distribution over response quality. We begin by outlining the continuous form of our reward modeling optimization paradigm, the **Probabilistic Reward Modeling** (§ 4.1). We then discretize this formulation into a tractable form, termed OPRM (§ 4.2) and conclude by presenting the complete training and inference pipeline for OPRM (§ 4.3).

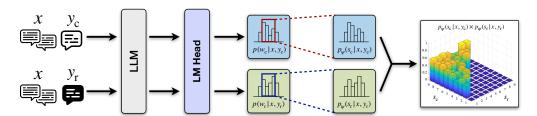


Figure 1: **The architectures of Ordinal Probabilistic Reward Model**. Given a problem and a pair of responses, designated as **chosen** and **rejected**, the OPRM utilizes its language model (LM) head to obtain the ordinal rating probabilities for each response. A joint probability matrix is then constructed by computing the Cartesian product of these two sets of probabilities for optimization.

4.1 PROBABILISTIC REWARD MODELING

Departing from the conventional Bradley-Terry reward model (Section 3), which estimates a single scalar value for each response, we propose a reward modeling objective derived from Random Utility Model theory (Manski, 1977; Cascetta, 2009). Our objective enables the model to learn a probability distribution over the quality of each response. Concretely, we model the quality score of a response y for a given input x as a continuous random variable S. This random variable is supported on a bounded interval $[a,b] \subset \mathbb{R}$. Our reward model, parameterized by ψ , learns the conditional probability density function (PDF) $p_{\psi}(s \mid x, y)$ of this variable, where s is a realization of S. This density must satisfy $\int_a^b p_{\psi}(s \mid x, y) \, ds = 1$.

Given a preference pair (y_c, y_r) with a chosen and a rejected response, we model their quality scores as two independent random variables, S_c and S_r . Their scores are drawn from the distributions defined by their respective conditional PDFs: $s_c \sim p_\psi(\cdot \mid x, y_c)$ and $s_r \sim p_\psi(\cdot \mid x, y_r)$. The probability of the preference $y_c \succ y_r$ is then modeled as the probability that the score of the chosen response exceeds that of the rejected one:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \mathbb{E}_{s_{c}, s_{r}} \left[\mathbb{1}(s_{c} > s_{r}) \right] \tag{3}$$

where $\mathbb{1}(\cdot)$ denotes the indicator function. Expanding the expectation in integral form over the bounded interval [a,b], we obtain:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \int_{a}^{b} \int_{a}^{b} \mathbb{1}(s_{c} > s_{r}) p_{\psi}(s_{c} \mid x, y_{c}) p_{\psi}(s_{r} \mid x, y_{r}) ds_{r} ds_{c}$$
(4)

This expression corresponds to computing the probability that a random score sampled from the chosen response exceeds a random score from the rejected response, integrating over their joint distribution constrained to the bounded interval. Since $\mathbb{1}(s_c > s_r) = 1$ only when $s_c > s_r$ and 0 otherwise, so it effectively truncates the integral domain to (a, s_c) for $p_{\psi}(s_r \mid x, y_r)$. Thus, we can equivalently restructure the double integral as follows:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \int_{a}^{b} p_{\psi}(s_{c} \mid x, y_{c}) \left(\int_{a}^{s_{c}} p_{\psi}(s_{r} \mid x, y_{r}) \, \mathrm{d}s_{r} \right) \, \mathrm{d}s_{c} \tag{5}$$

Finally, we can simply optimize the Eq. (5) by minimizing the negative log-likelihood loss. Notably, the Bradley-Terry model is a special case of the Probabilistic Reward Modeling framework, arising when the quality score distribution is constrained to a unimodal Gumbel distribution with fixed shape parameters (see Appendix B for a detailed proof). However, this objective lacks a closed-form analytical solution and requires estimation through Monte Carlo sampling. This computational challenge motivates our transition from the continuous formulation to a more tractable discrete one.

4.2 From Continuous to Discrete

To obtain a closed-form analytical solution, we adapt the continuous formulation in Eq. (5) by modeling the scores as discrete random variables over a finite set of ordinal ratings $\{a, a+1, \ldots, b\}$. This yields the following closed-form expression:

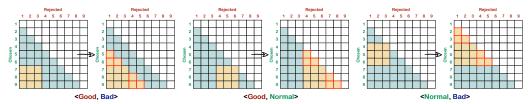


Figure 2: **Region Flooding Tuning**. To ensure the correctness of the reward modeling, region flooding is applied to each of the three partition combinations, resulting in a lower triangular form.

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \sum_{s_{c}=a}^{b} p_{\psi}(s_{c} \mid x, y_{c}) \left(\sum_{s_{r}=a}^{s_{c}-1} p_{\psi}(s_{r} \mid x, y_{r}) \right)$$
(6)

As observed in Eq. (2), the Bradley-Terry model maximizes the score gap between chosen (y_c) and rejected (y_r) responses, creating a steep reward landscape with pronounced gradients beneficial for RL optimization. Similarly, our optimization objective in Eq. (6) inherits and generalizes this desirable property. By operating over full reward distributions instead of single scalars, our objective naturally shifts probability mass upward for chosen responses and downward for rejected ones, thereby widening their separation. A detailed proof via gradient analysis can be found in Appendix C.

Ordinal Probabilistic Reward Modeling presents two key advantages: (1) *Quantifying Uncertainty*, the variance of the output distribution serves as a measure of model confidence—wide distributions for ambiguous comparisons indicate uncertainty, while sharp, peaked distributions reflect clear preferences, enhancing interpretability. Our method thus explicitly captures the inherent uncertainty in human preference judgments, a crucial aspect often overlooked by discriminative reward models. (2) *Handling Annotation Disagreement*, our method can represent multimodal score distributions (e.g., Mixture of Gaussians), enabling it to capture disagreements among annotators. By explicitly capturing conflicting signals within the score distribution, our model becomes robust to the performance degradation often caused by inconsistent preference data (Sun et al., 2025a). This contrasts sharply with traditional methods like the Bradley-Terry model, which are restricted to unimodal preferences.

4.3 PIPELINE

Our training pipeline, illustrated in Figure 1, begins by formatting the preference data pairs (x, y_c, y_r) (see Section 3) into a structured input using a prompt template. The details of this template and criteria are provided in Appendix F. As the next step in the pipeline, following the parameter-free technique from prior work (Cui et al., 2023), we compute the distribution over quality score $s \in \{a, a+1, \ldots, b\}$ (where $a, b \in \mathbb{Z}$) by directly repurposing the LM head's vocabulary probabilities, thus obviating the need for a separate prediction head and avoiding any new parameters. In our implementation, we set the quality score range from 1 to 9 (i.e., a=1,b=9) The score distribution is then formed by directly extracting the vocabulary probabilities of the corresponding numeric tokens (i.e., '1' to '9'). This approach allows the model to directly leverage its inherent ordinal knowledge of numbers.

In summary, both the chosen and rejected inputs are fed into the LLM backbone and its LM head, yielding the post-softmax vocabulary probability distributions $p(w_{\rm c}\mid x,y_{\rm c})$ and $p(w_{\rm r}\mid x,y_{\rm r})$ at the last token position. The probabilities of all numeric tokens are then normalized to form the distribution for our ordinal probabilistic reward modeling:

$$p_{\psi}(s_{c} = i \mid x, y_{c}) = \frac{p(w_{c} = \texttt{i} \mid x, y_{c})}{\sum_{j=1}^{9} p(w_{c} = \texttt{j} \mid x, y_{c})}, \ p_{\psi}(s_{r} = i \mid x, y_{r}) = \frac{p(w_{r} = \texttt{i} \mid x, y_{r})}{\sum_{j=1}^{9} p(w_{r} = \texttt{j} \mid x, y_{r})}$$
(7)

Finally, we substitute the obtained $p_{\psi}(s_{\rm c}=i\mid x,y_{\rm c})$ and $p_{\psi}(s_{\rm r}=i\mid x,y_{\rm r})$ into Eq. (6) and maximize $P_{\psi}(y_{\rm c}\succ y_{\rm r}\mid x)$ using the negative log-likelihood loss.

During the inference stage, we simply input a response y given prompt x to obtain a quality score distribution $p_{\psi}(s \mid x, y)$. We can derive a scalar reward score through either argmax or weighted averaging. For a discussion of other possible decoding strategies, see Appendix I.2. In our subsequent experiments, we adopt the straightforward weighted averaging approach to compute the reward score: $r_{\psi}(x,y) = \sum_{s=a}^{b} s \cdot p_{\psi}(s \mid x,y)$, avoiding the tie-prone argmax method.

Table 2: Overall results of different methods and models on four RM benchmarks. **bold numbers** indicate the best performance. <u>Underlined numbers</u> indicate the second best. The Overall* score is the average performance excluding Reward Bench due to its known data contamination issues.

Model	Reward Bench	PPE-P	PPE-C	RMB	Overall	Overall*	
Repo	Reported Results of Public Models						
Skywork-Reward-Gemma-2-27B	93.8	56.6	56.6	60.2	66.8	57.8	
DeepSeek-V2.5-0905	81.5	62.8	58.5	65.7	67.1	62.3	
Gemini-1.5-Pro	86.8	66.1	59.8	56.5	67.3	60.8	
ArmoRM-8B-v0.1	90.4	60.6	61.2	64.6	69.2	62.1	
InternLM2-20B-Reward	90.2	61.0	63.0	62.9	69.3	62.3	
LLaMA-3.1-70b-Instruct	84.1	65.3	59.2	68.9	69.4	64.5	
Claude-3.5-sonnet	84.2	65.3	58.8	70.6	69.7	63.2	
Nemotron-4-340B-Reward	92.0	59.3	60.8	69.9	70.5	63.3	
GPT-40	86.7	67.1	57.6	73.8	71.3	66.2	
Reproduced Results of Baseline Methods From DeepSeek							
LLM-as-a-Judge	83.4	64.2	58.8	64.8	67.8	62.6	
DeepSeek-BTRM-27B	81.7	68.3	66.7	57.9	68.6	64.3	
CLoud-Gemma-2-27B	82.0	67.1	62.4	63.4	68.7	64.3	
DeepSeek-PairRM-27B	87.1	65.8	64.8	58.2	69.0	62.9	
DeepSeek-GRM-27B-RFT	84.5	64.1	59.6	67.0	68.8	63.6	
DeepSeek-GRM-27B	86.9	64.7	59.8	69.0	69.9	64.5	
	Results of Our Me	thod					
OPRM-Qwen2.5-7B	87.8	61.1	61.3	71.5	70.4	64.6	
OPRM-Qwen2.5-14B	89.3	63.0	64.3	73.8	72.6	67.0	
OPRM-Qwen2.5-32B	91.3	63.9	66.1	75.6	74.2	68.5	
OPRM-Qwen2.5-72B	89.3	65.1	64.3	73.5	73.1	67.6	
Results of Our	Method (w/ Regio	n Floodin _ë	g Tuning)				
OPRM-RgFT-Qwen2.5-7B	86.2	62.3	62.4	70.1	70.3(\doldon10.1)	64.9(\(\daggered{10.3}\))	
OPRM-RgFT-Qwen2.5-14B	87.3	63.4	65.6	72.8	72.3(\\dot0.3)	67.3(\^0.3)	
OPRM-RgFT-Qwen2.5-32B	88.9	64.6	67.3	74.8	73.9(\\dot0.3)	68.9 (†0.4)	
OPRM-RgFT-Qwen2.5-72B	89.1	65.3	66.4	74.2	73.8(<u>†0.7)</u>	68.6(†1.0)	

5 REGION FLOODING TUNING

While OPRM effectively captures relative preferences, precisely aligning its scoring distribution with absolute quality judgments presents a further challenge. To address this, we introduce **Region Tuning** (RgT), a cost-effective method that enhances the model's fidelity to absolute quality scores using minimal annotations (§ 5.1). Subsequently, we refine RgT to preserve the desirable properties (as detailed in Appendix C), culminating in our final method: **Region Flooding Tuning** (RgFT) (§ 5.2).

5.1 REGION TUNING

Building upon the OPRM optimization objective from Eq. (6), which employs the finite set of ordinal ratings $S = \{1, 2, \dots, 9\}$ for all data, we introduce a more fine-grained partitioning based on the absolute quality of each response, a technique we term *Region Tuning* (RgT).

Specifically, we further partition the finite set into three quality levels, guiding the model to concentrate the probability mass within corresponding rating sub-region: $S_{bad} = \{1, 2, 3\}$, $S_{normal} = \{4, 5, 6\}$, and $S_{good} = \{7, 8, 9\}$. Consequently, for a single preference data point consisting of a chosen and a rejected response, there are six possible combinations of quality levels. These include pairs from different levels, as well as pairs where both responses fall into the same level, denoted as $< l_{chosen}, l_{rejected} > < cood, normal>, < cood, bad>, < normal, bad>, < cood, good, cood>, < normal, normal>, < bad, bad>.$

This partitioning allows us to redefine the preference probability by conditioning it on the quality levels of the chosen and rejected responses. Thus, the optimization objective is formulated as:

$$P_{\psi}(y_{c} \succ y_{r} \mid x, l_{\text{chosen}}, l_{\text{rejected}}) = \sum_{s_{c} \in S_{l_{\text{chosen}}}} p_{\psi}(s_{c} \mid x, y_{c}) \left(\sum_{s_{r} \in S_{l_{\text{rejected}}}} p_{\psi}(s_{r} \mid x, y_{r}) \mathbb{1}(s_{c} > s_{r}) \right)$$
(8)

Table 3: Detailed results of different methods on the PPE Correctness benchmark. **Bold numbers** indicate the best performance. Underlined numbers indicate the second best.

Model	MMLU-Pro	MATH	GPQA	MBPP-Plus	IFEval	PPE Correctness		
Results of Our Method								
OPRM-Qwen2.5-7B	65.2	70.1	56.3	59.0	56.1	61.3		
OPRM-Qwen2.5-14B	66.7	70.7	57.1	67.4	59.5	64.3		
OPRM-Qwen2.5-32B	71.2	73.2	57.9	66.2	62.2	66.1		
OPRM-Qwen2.5-72B	73.4	75.9	58.6	54.1	59.5	64.3		
	Results of Our Method (w/ Region Flooding Tuning)							
OPRM-RgFT-Qwen2.5-7B	64.8	71.6	55.9	63.0	56.8	62.4(†1.1)		
OPRM-RgFT-Qwen2.5-14B	69.5	74.0	57.3	67.0	60.0	65.6(†1.3)		
OPRM-RgFT-Qwen2.5-32B	73.3	76.8	58.5	67.2	60.6	67.3 (↑1.2)		
OPRM-RgFT-Qwen2.5-72B	72.8	77.1	59.0	62.0	61.2	<u>66.4</u> (<u>↑2.1</u>)		

5.2 From Region Tuning to Region Flooding Tuning

As shown in Figure 2, when $l_{\text{chosen}} \neq l_{\text{rejected}}$, Eq. (8) optimizes a square-shaped joint probability region, resulting in constant partial derivatives $\frac{\partial P}{\partial p_c(k)}$ and $\frac{\partial P}{\partial p_r(k)}$. In this case, the optimization objective no longer shifts the probability mass of the chosen response upwards and the rejected response downwards to increase their separation. This leads to the loss of a desirable property of OPRM, as mentioned in Section 4.2 (see Appendix C for a formal proof).

As shown in Figure 2, we propose region flooding to the optimized joint probability region, expanding it into a lower triangular shape to preserve the desired property. As its expansion process closely resembles breadth-first search algorithm, we term it *Region Flooding Tuning* (RgFT). RgFT provides three key advantages: (1) *Interpretability*, RgFT constrains the model to concentrate probability mass within the score regions correspond to pre-defined quality levels, enabling reward scores to more accurately reflect the absolute quality of responses. (2) *Semi-supervised Learning*, RgFT supports semi-supervised training by combining quality-labeled data with preference-only data. (3) *Customizability*, RgFT allows for the flexible tailoring of rating sub-regions to their corresponding quality levels, making the strategy adaptable to diverse application requirements (see Appendix I.1).

6 EXPERIMENTS

6.1 EXPERIMENTAL SETUP

In our experiments, we curate a dataset of 130k samples for reward model training, drawn primarily from publicly available open-source datasets: **Skywork Reward Preference 80K** (Liu et al., 2024b) and **UltraFeedback Binarized Preferences** (Cui et al., 2023). We employ the Qwen2.5-Instruction series of models (7B, 14B, 32B, and 72B) (Team, 2024) as the backbone for training the OPRM. We compare OPRM to different categories of baselines: **Discriminative RMs**, **Generative RMs** and **DeepSeek-RM**. Following prior work, we evaluate the performance of different methods on various RM benchmarks: **Reward Bench** (Lambert et al., 2024), **PPE** (Frick et al., 2024), **RMB** (Zhou et al., 2024). We use the standard pair accuracy and Best-of-N evaluation metrics for each benchmark. Detailed information on the data, baselines, benchmarks, and metrics is provided in Appendix D.

6.2 MAIN RESULTS

As shown in Table 2, we compares the overall results of OPRM with different baseline reward models on RM benchmarks. We present the performance of OPRM with the reported results of public models and the reproduced results of baseline methods from DeepSeek. We observe that OPRM outperforms the baseline methods in overall performance, and achieves competitive results against strong public RMs, such as Nemotron-4-340B-Reward and GPT-40. Notably, the 14B, 32b, and 72B models surpass all prior leading reward models, improving upon the previous best result by 1.3%, 2.9%, and 1.8%, respectively, despite being significantly smaller in scale. Moreover, the most significant performance enhancement is observed on the RMB and PPE-Correctness benchmarks, which utilize Best-of-N evaluation to better reflect practical effectiveness on downstream tasks. We attribute the 32B model's superior performance over the 72B model to the exceptional zero-shot capability of Qwen2.5-32B. This enables it to outperform larger models on the RM benchmark without any fine-tuning. The more detailed numbers on RewardBench, PPE Correctness, and RMB are in Table 5, Table 6, and Table 7 in Appendix G.

Table 4: Detailed results of Qwen2.5-7B with different methods on the Role Play benchmark. Bold numbers indicate the best performance. Underlined numbers indicate the second best.

Method	Pair-Accuracy	Best-of-N	Best-of-N-plus	Worst-of-N	Overall
Random Baseline	50.0	25.5	31.3	68.7	43.9
	Training	on Role Play	Data Only		
BT Model	70.4	48.6	51.3	83.6	63.5
BT Model - w/ Margin	71.0	49.3	52.3	84.2	64.2
OPRM (ours)	71.3	49.4	52.5	84.1	64.3
OPRM-RgFT (ours)	$72.1(\uparrow 0.8)$	50.7(↑ 1.3)	53.6(†1.1)	$85.1(\uparrow 1.0)$	$65.4(\uparrow 1.1)$
Tro	aining on Mixed Ro	ole Play and	General-Domain L	Data	
BT Model	73.8	51.2	54.3	86.0	66.3
BT Model - w/ Margin	75.3	53.4	55.7	87.2	67.9
OPRM (ours)	74.4	54.1	56.1	87.8	68.1
OPRM-RgFT (ours)	75.8 (<u>†1.4</u>)	55.8 (↑1.7)	59.3 (†3.2)	89.9 (†2.1)	70.2 (†2.1)

6.3 THE IMPACT OF REGION FLOODING TUNING

Building upon OPRM, we incorporate RgFT for further experimentation. We train the OPRM-RgFT series of models using the preference data, further enriched with three defined quality level annotations (good, normal, and bad). Consistent with the advantages of RgFT described in Section 5.

6.3.1 FEWER ANNOTATIONS, BETTER RESULTS

As presented in Table 2 and Table 3, our evaluation of OPRM-RgFT on four RM benchmarks reveals a notable performance divergence. On one hand, RgFT consistently improves performance across all model sizes on the PPE benchmarks. Notably, OPRM-RgFT-32B achieves SOTA accuracy of 67.3% on the PPE-Correctness benchmark, surpassing all prior leading reward models. On the other hand, its performance on other benchmarks is inconsistent. We hypothesize that this discrepancy stems from biases introduced by our annotation strategy for general data (see Appendix H.1). This process, involving coarse AI annotation with simple manual correction, is effective for verifiable tasks with explicit correctness labels like PPE-Correctness but likely introduces label noise for other tasks. Further supporting this claim, our subsequent experiments show that incorporating fine-grained manual annotations leads to consistent performance improvements.

6.3.2 TOWARDS HUMAN-ALIGNED SCORE DISTRIBUTIONS.

To evaluate the impact of Region Flooding Tuning on absolute quality assessment, we analyze the score distributions produced by our models. Specifically, we curated two distinct datasets for this analysis: an **Absolute-Good Set** with 100 high-quality prompt-response pairs and an **Absolute-Bad Set** with 100 poor-quality pairs. These pairs are manually selected by experts based on a multi-faceted evaluation across dimensions such as instruction following, factual accuracy, and helpfulness. We then score both datasets using three models: the baseline BTRM-32B, our base model OPRM-32B, and its RgFT-enhanced version. As illustrated in Figure 3, the base OPRM-32B already exhibits a basic capacity for absolute quality assessment: within its [1, 9] scoring range, it generally assigns scores above 5 to good responses and below 5 to bad ones. Crucially, OPRM-RgFT-32B significantly enhances this capabi

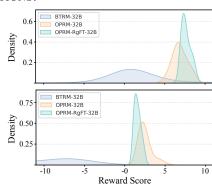


Figure 3: Comparison of score distributions for responses of high-quality (**Top**) and low-quality (**Bottom**).

cially, OPRM-RgFT-32B significantly enhances this capability. The RgFT-enhanced model polarizes the score distributions, pushing scores for high-quality responses into the [7,9] range while confining low-quality ones to [1,3]. This increased separation makes the score itself a more reliable and interpretable indicator of absolute quality. Case studies in Appendix J provide detailed scoring examples that further corroborate these findings and demonstrate the improved reliability of RgFT scores.

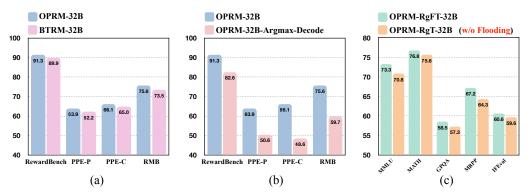


Figure 4: **Ablation Study**: (a) Assessing the superiority of OPRM over the BT Model. (b) Evaluating the efficacy of Weighted Average Decoding. (c) Validating the necessity of Region Flooding.

6.3.3 SEMI-SUPERVISED DOMAIN ADAPTATION

To simulate practical applications and reduce annotation costs, we investigate RgFT's effectiveness in a semi-supervised domain adaptation setting. Specifically, we curated a training set of 31K role-playing instances with quality-level labels (see Appendix H.2 for annotation details) and a mixed dataset by combining these with an equal volume of unlabeled general-domain preference data. For evaluation, we build a test set of 500 questions, each with 5-10 responses, and designed three new metrics: **Best-of-N** (top-scoring is **good**-level), **Worst-of-N** (bottom-scoring is **bad**-level), and **Best-of-N-plus**(top-scoring is not **bad**-level). As shown in Table 4, we benchmark our models against BT and BT-with-Margin baselines (see Appendix E for detailed formulas) under two settings: training on role-play data only, and on the mixed dataset. In both settings, OPRM surpasses the baselines, and OPRM-RgFT further improves upon OPRM. Crucially, incorporating unlabeled general-domain data significantly boosts the performance of OPRM and OPRM-RgFT from 64.3% to 68.1% and 65.4% to 70.2%, respectively. This demonstrates that RgFT can effectively leverage unlabeled preference data in a semi-supervised manner, offering a cost-effective path for domain adaptation.

6.4 ABLATION STUDY

We ablate key components of our OPRM and RgFT to validate their contributions. As shown in Figure 4, each component proves essential for optimal performance.

Effectiveness of OPRM Loss. We replace our OPRM loss with the standard Bradley-Terry (BT) loss, training an identical 32B model. Figure 4(a) shows this change caused a 1.1% to 2.1% performance drop across all benchmarks, validating the superiority of modeling reward as an ordinal variable.

Impact of Decoding Method. We compared our standard weighted averaging decoding with a simpler Argmax approach, which directly selects the token with the highest probability. As shown in Figure 4(b), Argmax decoding led to a substantial 8.7% to 17.5% performance drop. We attribute this to Argmax's inability to capture fine-grained quality differences, which resulted in excessive ties.

Necessity of the Flooding Mechanism. The flooding mechanism is designed to create desirable lower triangular score regions (see Appendix C). Removing it resulted in a 1.0% to 2.9% performance drop on the PPE Correctness benchmark (Figure 4(c)). The degradation are most pronounced when distinguishing between responses of similar quality, confirming the mechanism's critical role.

7 CONCLUSION

In this paper, we propose *Ordinal Probabilistic Reward Model*, a novel paradigm that learns a full probability distribution over an ordinal reward space. To better anchor these rewards to absolute quality, we further proposed *Region Flooding Tuning*, a training strategy that leverages quality-level annotations to calibrate the model's probability distribution. Extensive experiments on four diverse reward modeling benchmarks show that our approach consistently improves performance by 2.9% to 7.4%. Furthermore, detailed analysis reveals that OPRM is superior to the conventional Bradley-Terry model and that RgFT is crucial for discerning fine-grained quality differences. We believe OPRM with RgFT offer a powerful framework for developing more accurate and reliable reward models, a critical step towards building more capable and aligned large language models.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. URL https://arxiv.org/abs/2303.08774.
- Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D Chang, and Prithviraj Ammanabrolu. Critique-out-loud reward models. *arXiv preprint arXiv:2408.11791*, 2024. URL https://arxiv.org/abs/2408.11791.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. URL https://arxiv.org/abs/2204.05862.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024. URL https://arxiv.org/abs/2403.17297.
- Ennio Cascetta. Random utility theory. In *Transportation systems analysis: models and applications*, pp. 89–167. Springer, 2009.
- Nuo Chen, Zhiyuan Hu, Qingyun Zou, Jiaying Wu, Qian Wang, Bryan Hooi, and Bingsheng He. Judgelrm: Large reasoning models as a judge. *arXiv preprint arXiv:2504.00050*, 2025a. URL https://arxiv.org/abs/2504.00050.
- Xiusi Chen, Gaotang Li, Ziqi Wang, Bowen Jin, Cheng Qian, Yu Wang, Hongru Wang, Yu Zhang, Denghui Zhang, Tong Zhang, et al. Rm-r1: Reward modeling as reasoning. *arXiv* preprint *arXiv*:2505.02387, 2025b. URL https://arxiv.org/abs/2505.02387.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv* preprint arXiv:2310.01377, 2023. URL https://arxiv.org/abs/2310.01377.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023. URL https://arxiv.org/abs/2305.14233.
- Nicolai Dorka. Quantile regression for distributional reward models in rlhf. *arXiv preprint arXiv:2409.10164*, 2024. URL https://arxiv.org/abs/2409.10164.
- Evan Frick, Tianle Li, Connor Chen, Wei-Lin Chiang, Anastasios N Angelopoulos, Jiantao Jiao, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. How to evaluate reward models for rlhf. *arXiv* preprint arXiv:2410.14872, 2024. URL https://arxiv.org/abs/2410.14872.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023. URL https://proceedings.mlr.press/v202/gao23h.html.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. URL https://arxiv.org/abs/2407.21783.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025a. URL https://arxiv.org/abs/2501.12948.

```
Jiaxin Guo, Zewen Chi, Li Dong, Qingxiu Dong, Xun Wu, Shaohan Huang, and Furu Wei. Reward reasoning model. arXiv preprint arXiv:2505.14674, 2025b. URL https://arxiv.org/abs/2505.14674.
```

- Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. Won't get fooled again: Answering questions with false premises. *arXiv* preprint arXiv:2307.02394, 2023. URL https://arxiv.org/abs/2307.02394.
- Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models. arXiv preprint arXiv:2503.06749, 2025. URL https://arxiv.org/abs/2503.06749.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. URL https://arxiv.org/abs/2410.21276.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. LLM-blender: Ensembling large language models with pairwise ranking and generative fusion. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14165–14178, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.792. URL https://aclanthology.org/2023.acl-long.792/.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025. URL https://arxiv.org/abs/2503.09516.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, et al. Prometheus: Inducing fine-grained evaluation capability in language models. In *The Twelfth International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=8euJaTveKw.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*, 2024. URL https://arxiv.org/abs/2403.13787.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=v8L0pN6EOi.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. arXiv preprint arXiv:2109.07958, 2021. URL https://arxiv.org/abs/2109.07958.
- Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024a. URL https://arxiv.org/abs/2405.04434.
- Chris Yuhao Liu, Liang Zeng, Jiacai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. Skywork-reward: Bag of tricks for reward modeling in llms. *arXiv* preprint *arXiv*:2410.18451, 2024b. URL https://arxiv.org/abs/2410.18451.
- Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu. Inference-time scaling for generalist reward modeling. *arXiv preprint arXiv:2504.02495*, 2025. URL https://arxiv.org/abs/2504.02495.
- Xingzhou Lou, Dong Yan, Wei Shen, Yuzi Yan, Jian Xie, and Junge Zhang. Uncertainty-aware reward model: Teaching reward models to know what is unknown. *arXiv preprint arXiv:2410.00847*, 2024. URL https://arxiv.org/abs/2410.00847.

- Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. Generative reward models. *arXiv* preprint *arXiv*:2410.12832, 2024. URL https://arxiv.org/abs/2410.12832.
- Charles F Manski. The structure of random utility models. *Theory and decision*, 8(3):229, 1977.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. URL https://arxiv.org/abs/2402.03300.
- Hao Sun, Yunyi Shen, and Jean-Francois Ton. Rethinking reward modeling in preference-based large language model alignment. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL https://openreview.net/forum?id=rfdblE10qm.
- Wangtao Sun, Xiang Cheng, Xing Yu, Haotian Xu, Zhao Yang, Shizhu He, Jun Zhao, and Kang Liu. Probabilistic uncertain reward model. *arXiv preprint arXiv:2503.22480*, 2025b. URL https://arxiv.org/abs/2503.22480.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. URL https://arxiv.org/abs/2403.05530.
- Qwen Team. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024. URL https://metaso-static.oss-accelerate.aliyuncs.com/metaso/document/64a141da-f885-44ab-9883-94b03b737cdf.pdf.
- Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 10582–10592, Miami, Florida, USA, November 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.620. URL https://aclanthology.org/2024.findings-emnlp.620/.
- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. Helpsteer2-preference: Complementing ratings with preferences. *arXiv* preprint arXiv:2410.01257, 2024b. URL https://arxiv.org/abs/2410.01257.
- Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. Helpsteer2: Open-source dataset for training top-performing reward models. *arXiv preprint arXiv:2406.08673*, 2024c. URL https://arxiv.org/abs/2406.08673.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv* preprint *arXiv*:2109.01652, 2021. URL https://arxiv.org/abs/2109.01652.
- Chenxi Whitehouse, Tianlu Wang, Ping Yu, Xian Li, Jason Weston, Ilia Kulikov, and Swarnadeep Saha. J1: Incentivizing thinking in llm-as-a-judge via reinforcement learning. *arXiv* preprint arXiv:2505.10320, 2025. URL https://arxiv.org/abs/2505.10320.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=CfXh93NDgH.

- Minghao Yang, Chao Qu, and Xiaoyu Tan. Inf-orm-llama3.1-70b, 2024. URL [https://huggingface.co/infly/INF-ORM-Llama3.1-70B] (https://huggingface.co/infly/INF-ORM-Llama3.1-70B).
- Zihuiwen Ye, Fraser Greenlee-Scott, Max Bartolo, Phil Blunsom, Jon Ander Campos, and Matthias Gallé. Improving reward models with synthetic critiques. *arXiv preprint arXiv:2405.20850*, 2024. URL https://arxiv.org/abs/2405.20850.
- Jiachen Yu, Shaoning Sun, Xiaohui Hu, Jiaxu Yan, Kaidong Yu, and Xuelong Li. Improve Ilmas-a-judge ability as a general ability. *arXiv preprint arXiv:2502.11689*, 2025. URL https://arxiv.org/abs/2502.11689.
- Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan, Chenguang Zhu, Richard Yuanzhe Pang, Yundi Qian, Xuewei Wang, Suchin Gururangan, Chao Zhang, et al. Self-generated critiques boost reward modeling for language models. arXiv preprint arXiv:2411.16646, 2024. URL https://arxiv.org/abs/2411.16646.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. *arXiv* preprint arXiv:2408.15240, 2024. URL https://arxiv.org/abs/2408.15240.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_Benchmarks.html.
- Jialun Zhong, Wei Shen, Yanzeng Li, Songyang Gao, Hua Lu, Yicheng Chen, Yang Zhang, Wei Zhou, Jinjie Gu, and Lei Zou. A comprehensive survey of reward models: Taxonomy, applications, challenges, and future. *arXiv preprint arXiv:2504.12328*, 2025. URL https://arxiv.org/abs/2504.12328.
- Enyu Zhou, Guodong Zheng, Binghai Wang, Zhiheng Xi, Shihan Dou, Rong Bao, Wei Shen, Limao Xiong, Jessica Fan, Yurong Mou, et al. Rmb: Comprehensively benchmarking reward models in llm alignment. *arXiv preprint arXiv:2410.09893*, 2024. URL https://arxiv.org/abs/2410.09893.

APPENDIX

A THE USE OF LARGE LANGUAGE MODELS (LLMS)

We utilize Large Language Models (LLMs) to aid in the writing and polishing of this manuscript. Specifically, LLMs are employed to correct grammatical errors, improve sentence structure, and enhance the clarity and conciseness of the text. This process is primarily applied to the Introduction, Related Work, and Appendix sections. All scientific contributions, methodologies, and conclusions presented in this paper are the original work of the authors. The LLMs serve solely as a writing-enhancement tool.

B BRADLEY-TERRY AS A SPECIAL CASE OF PROBABILISTIC REWARD MODELING

In this section, we demonstrate that the Bradley-Terry model for pairwise preferences can be derived from a more general probabilistic reward modeling framework under a specific set of distributional assumptions.

Let $p_{\psi}(s \mid x, y)$ denote the probability density function of a score s assigned to a response y given a context x, where the scoring mechanism is parameterized by ψ . Consider two responses for the same context x: a chosen response y_c and a rejected response y_r . Let s_c and s_r be the random variables for their respective scores, with distributions $p_{\psi}(s_c \mid x, y_c)$ and $p_{\psi}(s_r \mid x, y_r)$. We assume s_c and s_r are conditionally independent given s_r , s_r , s_r , s_r .

The probability that y_c is preferred over y_r , denoted $P_{\psi}(y_c \succ y_r \mid x)$, is the probability that the score of the chosen response is greater than that of the rejected one, i.e., $P(s_c > s_r)$. This can be expressed generally as:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \int_{-\infty}^{\infty} p_{\psi}(s_{c} \mid x, y_{c}) \left(\int_{-\infty}^{s_{c}} p_{\psi}(s_{r} \mid x, y_{r}) \, \mathrm{d}s_{r} \right) \mathrm{d}s_{c}. \tag{9}$$

We show that this general formulation reduces to the Bradley-Terry model under specific assumptions. For clarity, we first establish the functional form of the Bradley-Terry model in terms of the sigmoid function.

Lemma B.1 (Bradley-Terry Model in Sigmoid Form). The Bradley-Terry (BT) model, which defines the preference probability based on underlying quality scores $r_{\psi}(x, y_{c})$ and $r_{\psi}(x, y_{r})$ as

$$P_{BT}(y_{c} \succ y_{r} \mid x) = \frac{\exp(r_{\psi}(x, y_{c}))}{\exp(r_{\psi}(x, y_{c})) + \exp(r_{\psi}(x, y_{r}))},$$
(10)

is equivalent to the sigmoid function of the difference in scores:

$$P_{BT}(y_{c} \succ y_{r} \mid x) = \sigma(r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r})),$$
 (11)

where $\sigma(z) = 1/(1 + e^{-z})$ is the standard logistic sigmoid function.

Proof. We start from the standard definition of the BT model and manipulate it algebraically. By dividing the numerator and the denominator of Eq. (equation 10) by $\exp(r_{\psi}(x, y_{r}))$, we obtain:

$$\begin{split} P_{\text{BT}}(y_{\text{c}} \succ y_{\text{r}} \mid x) &= \frac{\exp(r_{\psi}(x, y_{\text{c}})) \cdot \exp(-r_{\psi}(x, y_{\text{r}}))}{(\exp(r_{\psi}(x, y_{\text{c}})) + \exp(r_{\psi}(x, y_{\text{r}}))) \cdot \exp(-r_{\psi}(x, y_{\text{r}}))} \\ &= \frac{\exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))}{\exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}})) + 1} \\ &= \frac{\exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))}{1 + \exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))}. \end{split}$$

To bring this into the form of the sigmoid function $\sigma(z)$, we can divide the numerator and denominator by $\exp(r_{\psi}(x, y_{\rm c}) - r_{\psi}(x, y_{\rm r}))$:

$$\begin{split} P_{\text{BT}}(y_{\text{c}} \succ y_{\text{r}} \mid x) &= \frac{1}{\frac{1 + \exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))}{\exp(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))}} \\ &= \frac{1}{\exp(-(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}}))) + 1} \\ &= \sigma(r_{\psi}(x, y_{\text{c}}) - r_{\psi}(x, y_{\text{r}})). \end{split}$$

This completes the proof of the lemma.

With this lemma, we can prove the main proposition.

Proposition B.2. The general preference probability $P_{\psi}(y_c \succ y_r \mid x)$ defined in Eq. (9) is equivalent to the Bradley-Terry model if the following assumptions hold:

- 1. The score difference $\Delta s_{\psi} \triangleq s_{\rm c} s_{\rm r}$ follows a logistic distribution.
- 2. The mean of this logistic distribution is the difference of deterministic underlying quality scores: $\mu = r_{\psi}(x, y_c) r_{\psi}(x, y_r)$.
- 3. The scale parameter of the logistic distribution is unity (s = 1).

Proof. The preference probability is the probability that the score of the chosen response exceeds that of the rejected one. This can be expressed in terms of the score difference random variable $\Delta s_{\psi} = s_{\rm c} - s_{\rm r}$:

$$P_{\psi}(y_c \succ y_r \mid x) = P(s_c > s_r) = P(\Delta s_{\psi} > 0).$$
 (12)

Assumption 1 states that Δs_{ψ} follows a logistic distribution. The cumulative distribution function (CDF) of a logistic random variable Z with mean μ and scale s is given by $F_Z(z) = (1 + e^{-(z-\mu)/s})^{-1}$. Therefore, we can compute the preference probability as:

$$P(\Delta s_{\psi} > 0) = 1 - P(\Delta s_{\psi} \le 0)$$

$$= 1 - F_{\Delta s_{\psi}}(0)$$

$$= 1 - \frac{1}{1 + e^{-(0 - \mu)/s}}$$

$$= 1 - \frac{1}{1 + e^{\mu/s}}$$

$$= \frac{(1 + e^{\mu/s}) - 1}{1 + e^{\mu/s}} = \frac{e^{\mu/s}}{1 + e^{\mu/s}}$$

$$= \frac{1}{1 + e^{-\mu/s}}.$$
(13)

This final expression is precisely the sigmoid function, $\sigma(\mu/s)$.

Then, we apply the remaining assumptions. **Assumption 2** posits that the mean of the distribution is $\mu = r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r})$. **Assumption 3** sets the scale parameter to unity, s = 1. Substituting these into our result from Eq. (13) yields:

$$P_{\psi}(y_{c} \succ y_{r} \mid x) = \frac{1}{1 + e^{-(r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r}))}} = \sigma(r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r})). \tag{14}$$

From Lemma B.1, we know that the Bradley-Terry model also simplifies to $\sigma(r_{\psi}(x,y_{c})-r_{\psi}(x,y_{r}))$. Since the probabilistic reward model under the specified assumptions and the Bradley-Terry model both yield the identical functional form, we have shown that the latter is a special case of the former.

C GRADIENT ANALYSIS OF THE PREFERENCE PROBABILITY

In this section, we conduct a formal gradient-based analysis to demonstrate that maximizing the preference probability, $P_{\psi}(y_c \succ y_r \mid x)$, incentivizes the underlying probabilistic model to maximally separate the score distributions of the chosen and rejected responses.

Proposition C.1 (Optimization Incentive of Preference Maximization). Let the scores for responses be drawn from a discrete set $\{a, a+1, \ldots, b\}$. Let $p_{\rm c}(k) \triangleq p_{\psi}(s_{\rm c}=k \mid x, y_{\rm c})$ and $p_{\rm r}(k) \triangleq p_{\psi}(s_{\rm r}=k \mid x, y_{\rm r})$ be the respective probability mass functions (PMFs). Maximizing the preference probability $P(s_{\rm c}>s_{\rm r})$ with respect to the variables $\{p_{\rm c}(k)\}$ and $\{p_{\rm r}(k)\}$ under the constraints $\sum_k p_{\rm c}(k)=1$ and $\sum_k p_{\rm r}(k)=1$ creates the following incentives:

- 1. For the chosen response y_c , shifting probability mass from any score k to a higher score k+1 will increase or maintain the objective value.
- 2. For the rejected response y_r , shifting probability mass from any score k+1 to a lower score k will increase or maintain the objective value.

This implies that the optimization process drives the PMF of y_c towards the maximum score b and the PMF of y_r towards the minimum score a.

Proof. The preference probability $P \triangleq P_{\psi}(y_c \succ y_r \mid x)$ for discrete scores is given by:

$$P = \sum_{i=a}^{b} p_{c}(i) P(s_{r} < i) = \sum_{i=a}^{b} p_{c}(i) \left(\sum_{j=a}^{i-1} p_{r}(j) \right).$$
 (15)

We analyze the gradient of P with respect to the probability mass at each score for y_c and y_r separately.

Part 1: Incentive for the Chosen Response Score (y_c) . We first compute the partial derivative of P with respect to $p_c(k)$ for some score $k \in \{a, \ldots, b\}$. From Eq. (15), only the term where i = k depends on $p_c(k)$, so:

$$\frac{\partial P}{\partial p_{c}(k)} = \frac{\partial}{\partial p_{c}(k)} \left[p_{c}(k) \sum_{j=a}^{k-1} p_{r}(j) \right] = \sum_{j=a}^{k-1} p_{r}(j) = P(s_{r} < k). \tag{16}$$

This derivative represents the sensitivity of the objective to an increase in probability mass at score k. To understand the incentive for shifting mass, consider moving an infinitesimal probability mass $\epsilon > 0$ from a score k to a higher score k+1. This corresponds to a change in the PMF: $p_c(k) \to p_c(k) - \epsilon$ and $p_c(k+1) \to p_c(k+1) + \epsilon$. The resulting change in P, denoted ΔP , can be approximated by the first-order Taylor expansion (which is exact since P is linear in p_c):

$$\Delta P \approx \epsilon \frac{\partial P}{\partial p_{c}(k+1)} - \epsilon \frac{\partial P}{\partial p_{c}(k)}$$

$$= \epsilon \left(P(s_{r} < k+1) - P(s_{r} < k) \right) \quad \text{(using Eq. equation 16)}$$

$$= \epsilon \cdot P(s_{r} = k)$$

$$= \epsilon \cdot p_{r}(k). \tag{17}$$

Since probabilities are non-negative, $p_{\rm r}(k) \geq 0$, and we defined $\epsilon > 0$, it follows that $\Delta P \geq 0$. This demonstrates that any shift of probability mass to a higher score for $y_{\rm c}$ is guaranteed to be a non-decreasing change in the objective function. This creates a persistent optimization pressure to move the entire distribution $p_{\rm c}$ towards the maximum score b.

Part 2: Incentive for the Rejected Response Score (y_r) . To analyze the effect of $p_r(k)$, it is more convenient to rewrite Eq. (15) by swapping the order of summation:

$$P = \sum_{j=a}^{b-1} p_{\rm r}(j) \left(\sum_{i=j+1}^{b} p_{\rm c}(i) \right).$$
 (18)

The partial derivative of P with respect to $p_r(k)$ for $k \in \{a, ..., b-1\}$ is:

$$\frac{\partial P}{\partial p_{\rm r}(k)} = \frac{\partial}{\partial p_{\rm r}(k)} \left[p_{\rm r}(k) \sum_{i=k+1}^b p_{\rm c}(i) \right] = \sum_{i=k+1}^b p_{\rm c}(i) = P(s_{\rm c} > k). \tag{19}$$

Now, consider shifting an infinitesimal probability mass $\epsilon>0$ from a score k+1 to a *lower* score k. This corresponds to the change: $p_{\rm r}(k)\to p_{\rm r}(k)+\epsilon$ and $p_{\rm r}(k+1)\to p_{\rm r}(k+1)-\epsilon$. The resulting change in P is:

$$\begin{split} \Delta P &\approx \epsilon \frac{\partial P}{\partial p_{\rm r}(k)} - \epsilon \frac{\partial P}{\partial p_{\rm r}(k+1)} \\ &= \epsilon \left(P(s_{\rm c} > k) - P(s_{\rm c} > k+1) \right) \quad \text{(using Eq. equation 19)} \\ &= \epsilon \cdot P(s_{\rm c} = k+1) \\ &= \epsilon \cdot p_{\rm c}(k+1). \end{split} \tag{20}$$

Since $p_{\rm c}(k+1) \geq 0$ and $\epsilon > 0$, we have $\Delta P \geq 0$. This shows that shifting probability mass to a lower score for $y_{\rm r}$ is always a non-decreasing change. This creates a consistent optimization pressure to move the distribution $p_{\rm r}$ towards the minimum score a.

Combining both parts, we have formally shown that maximizing the preference probability $P(s_c > s_r)$ drives the model to separate the score distributions by pushing the mass of p_c towards the highest possible score and the mass of p_r towards the lowest possible score.

D DETAILED EXPERIMENTAL SETUP

Training Preference Data. We curate a dataset of 130k samples for reward model training, drawn primarily from publicly available open-source datasets: Skywork Reward Preference 80K (Liu et al., 2024b) is a high-quality, pairwise preference dataset that spans multiple domains, including chat, safety, mathematics, and code. It employs advanced data filtering techniques to ensure the reliability of preferences across different tasks. UltraFeedback Binarized Preferences (Cui et al., 2023) is a large-scale, fine-grained, and diverse preference dataset designed for training powerful reward and critic models. It comprises approximately 64k prompts from various sources, including UltraChat (Ding et al., 2023), ShareGPT, Evol-Instruct (Xu et al., 2024), TruthfulQA (Lin et al., 2021), FalseQA (Hu et al., 2023), and FLAN (Wei et al., 2021). Each prompt is used to query multiple LLMs to generate four distinct responses, resulting in a total of 256k samples.

Baselines. In our main experiments, we employ the Qwen2.5-Instruction series of models (7B, 14B, 32B, and 72B) (Team, 2024) as the backbone for training the OPRM. We compare OPRM to different categories of baselines: (1) Discriminative RMs, including Skywork-Reward (Liu et al., 2024b), ArmoRM (Wang et al., 2024a), InternLM-20B-Reward (Cai et al., 2024), and Nemotron-4-340B-Reward (Wang et al., 2024c). (2) Generative RMs, including DeepSeek-V2.5 (Liu et al., 2024a), Gemini-1.5-Pro (Team et al., 2024), LLaMA-3.1-70B (Grattafiori et al., 2024), Claude-3.5-sonnet, and GPT-4o (Hurst et al., 2024). (3) DeepSeek-RM, a collection of baselines re-implemented by DeepSeek, including LLM-as-A-Judge (Zheng et al., 2023), DeepSeek-BTRM (Bradley & Terry, 1952), DeepSeek-PairRM (Jiang et al., 2023), CLoud-Gemma-2 (Ankner et al., 2024) and DeepSeek-GRM (Liu et al., 2025).

Benchmarks and Evaluation Metrics. Following prior work, we evaluate the performance of different methods on various RM benchmarks of different domains: Reward Bench (Lambert et al., 2024), PPE-Preference, PPE-Correctness (Frick et al., 2024), RMB (Zhou et al., 2024). We use the standard Best-of-N evaluation metrics for each benchmark: the accuracy of picking the best response from a set of responses. Specifically, Reward Bench and PPE Preference involve pairwise comparisons, with each prompt featuring two candidate responses. In contrast, PPE Correctness is designed for a large-scale Best-of-N evaluation, presenting 32 responses for each prompt. RMB is a hybrid, incorporating both pairwise comparison tasks and a Best-of-5 selection format.

E BRADLEY-TERRY LOSS WITH MARGIN

Inspired by INF-ORM (Yang et al., 2024), which employs GPT-40 to evaluate the preference margin between chosen and rejected responses, we annotate each pair in our dataset with a margin label. The original evaluation in INF-ORM follows these rules: (1) If the chosen answer is much better than rejected answer, set margin to 10; (2) If the chosen answer is better than the rejected answer, set margin to 3; (3) If the chosen answer is slightly better than rejected answer, set margin to 1.

Analogously, we define margins based on the combination of quality-level annotations. Specifically, pairs with the same quality level, such as **<good**, **good>**, **<normal**, **normal>**, and **<bad>**, are assigned a margin of 1. Pairs with adjacent quality levels, namely **<good**, **normal>** and **<normal**, **bad>**, are assigned a margin of 3. Finally, a margin of 10 is assigned to pairs with distant quality levels, like **<good**, **bad>**.

After that, the Bradley-Terry Loss with Margin is defined as:

$$\mathcal{L}(r_{\psi}) = -\mathbb{E}_{(x, y_{c}, y_{r}) \sim \mathcal{D}_{rm}}[m(x, y_{c}, y_{r}) \cdot \log \sigma(r_{\psi}(x, y_{c}) - r_{\psi}(x, y_{r}))], \tag{21}$$

Here, $m(x, y_c, y_r)$ stands for the margin value between chosen and rejected responses. This formula helps the model to better understand which responses are preferred over others, based on the scores we gave them.

F OPRM'S PROMPT TEMPLATES

As shown in Prompt F.1, we present the prompt templates used during the OPRM training process. Furthermore, we designed distinct scoring criteria for general-domain and domain-specific data. Specifically, the general-domain preference data is evaluated using the OpenAI 3H Criteria, as detailed in Prompt F.2. In contrast, the role-play data used in our experiments is assessed against our custom Role Play Criteria, presented in Prompt F.2.

F.1 PROMPT

```
Ordinal Distributional Reward Modeling Prompt
## Evaluation Guidelines
You are an expert text quality evaluator.
Your task is to carefully assess and score the AI-generated candidate response based on the
specific criteria provided below.
## Scoring Criteria
{criteria}
## Context of Interaction
The following represents a historical conversation between a user and an AI assistant:
{prompt}
## AI Response to Evaluate
Below is the candidate response generated by the AI that needs to be evaluated:
{response}
## Evaluation Instructions
Please score the response according to the criteria on a scale from 1 to 9, where:
- 1-3: Bad performance
- 4-6: Normal performance
- 7-9: Good performance
Use the following format for your final assessment: "Score: <total points>".
Score:
```

F.2 CRITERIA

Role Play Criteria

Relevant: Responses should not be irrelevant. Responses should be coherent and reasonable within the given context, without containing logical flaws.

Credible: Responses should be credible, meet the prompt requirements and provide accurate information. In role-playing or storytelling, responses should be consistent with the given character and background.

Interesting: Responses should be engaging and able to capture the user's attention. Responses that are self-aware, thought-provoking, or captivating will be more interesting.

Openai 3H Criteria

Helpfulness: The content generated should aim to be of assistance to the user. It should provide accurate information, answer questions to the best of the model's knowledge, and offer solutions that are practical and relevant to the user's queries. The model should not generate content that is intentionally misleading or unhelpful.

Truthfulness: The model should strive to provide information that is truthful and based on the most accurate data available up to its knowledge cutoff. It should not fabricate information or present opinions as facts. When uncertain, the model should acknowledge the limitations of its knowledge.

Harmlessness: The content should not cause harm to individuals or groups. It should avoid generating content that is offensive, discriminatory, or promotes harmful behaviors. The model should be programmed to refuse to generate content that could be used for malicious purposes, including but not limited to generating false information, engaging in deception, or promoting illegal activities.

G DETAILED EXPERIMENT RESULTS

We report the detailed per-subset experiment results on RewardBench 5, PPE Correctness 6, and RMB 7. The results for the baseline methods are sourced from their original papers.

H QUALITY-LEVEL ANNOTATION

H.1 GENERAL-DOMAIN DATA

For the acquisition of quality-level annotations, we follow the methodology of UltraFeedback (Cui et al., 2023). The process involves two main steps. First, we employ the gpt-40 model (Hurst et al., 2024) to annotate each prompt-response pair with fine-grained scores across multiple dimensions, such as instruction-following, truthfulness, and helpfulness. Second, these scores are averaged, and the resulting value is mapped to one of our three predefined quality levels—good, normal, or bad—based on specific score intervals.

Following this automatic annotation, we perform a manual verification step. For verifiable tasks, such as mathematics and coding, we check the responses against the ground truth. If a response is found to be incorrect, its quality level is manually downgraded to **bad**. We acknowledge that for more subjective tasks, this per-instance verification is not always feasible, which may introduce some annotation noise.

Finally, to ensure logical consistency, we filter out all pairs where the chosen response is not strictly better than the rejected one. This includes invalid combinations of $\langle l_{\text{chosen}}, l_{\text{rejected}} \rangle$ such as $\langle normal, good \rangle$, $\langle bad, normal \rangle$, and $\langle bad, good \rangle$. The remaining data constitutes our final training set.

H.2 ROLE-PLAY DATA

For the domain-specific data, we employe a team of six human experts to perform accurate quality-level annotation. The data consists of Role-Play Dialogues, for which the experts assessed each prompt-response pair against four core dimensions: **Core Role-Playing Consistency, Interactivity**

Table 5: Detailed results of different methods on the Reward Bench benchmark.

Method	Chat	Chat Hard	Safety	Reasoning	Reward Bench			
Reported Results of Public Models								
Skywork-Reward-Gemma-2-27B	95.8	91.4	91.9	96.1	93.8			
DeepSeek-V2.5-0905	-	-	-	-	81.5			
Gemini-1.5-Pro	94.1	77.0	85.8	90.2	86.8			
ArmoRM-8B-v0.1	96.9	76.8	90.5	97.3	90.4			
InternLM2-20B-Reward	98.9	76.5	89.5	95.8	90.2			
LLaMA-3.1-70b-Instruct	97.2	70.2	82.8	86.0	84.1			
Claude-3.5-sonnet	96.4	74.0	81.6	84.7	84.2			
Nemotron-4-340B-Reward	95.8	87.1	91.5	93.6	92.0			
GPT-40	96.1	76.1	88.1	86.6	86.7			
Reproduced Results of Baseline Methods From DeepSeek								
LLM-as-a-Judge	96.7	69.3	83.5	84.3	83.4			
DeepSeek-BTRM-27B	96.7	86.2	75.7	89.8	81.7			
CLoud-Gemma-2-27B	96.7	69.3	83.5	84.3	82.0			
DeepSeek-PairRM-27B	95.5	86.8	52.3	92.0	87.1			
DeepSeek-GRM-27B-RFT	94.7	77.2	87.0	79.2	84.5			
DeepSeek-GRM-27B	94.1	78.3	88.0	83.8	86.0			
	Resu	lts of Our Meth	hod					
OPRM-Qwen2.5-7B	96.4	76.3	86.2	92.2	87.8			
OPRM-Qwen2.5-14B	96.6	78.1	86.1	96.2	89.3			
OPRM-Qwen2.5-32B	96.9	81.8	89.6	96.7	91.3			
OPRM-Qwen2.5-72B	96.4	79.6	88.1	93.0	89.3			
Results of C	Results of Our Method (w/ Region Flooding Tuning)							
OPRM-RgFT-Qwen2.5-7B	95.5	76.5	86.4	86.5	86.2			
OPRM-RgFT-Qwen2.5-14B	96.9	79.4	88.1	84.6	87.3			
OPRM-RgFT-Qwen2.5-32B	95.3	82.7	89.2	88.4	88.9			
OPRM-RgFT-Qwen2.5-72B	96.9	82.7	89.7	87.1	89.1			

& Narrative Progression, Fundamental Linguistic Quality, and Immersion. Based on this multidimensional evaluation, they directly assigned a quality level of **good**, **normal**, or **bad** to each prompt-response pair. To ensure high annotation quality and consistency, a label was only accepted if at least two experts reached a consensus. We consider this high-fidelity annotation process to be crucial for fully leveraging the capabilities of our Region Flooding Tuning method.

I FUTURE WORK

In this section, we outline several promising directions for future research that build upon the OPRM framework, such as customized Region Flooding Tuning method and customized decoding method.

I.1 CUSTOMIZED REGION FLOODING TUNING

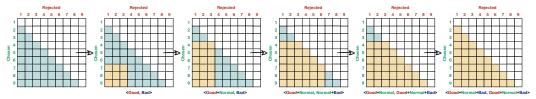


Figure 5: **Annotation Region Flooding Tuning**. As annotation ambiguity increases, the target optimization region "floods" to encompass a wider set of plausible outcomes. A more uncertain annotation results in a larger target region than a more certain one.

As discussed in Section 5.2, a core advantage of Region Flooding Tuning is its customizability. In this section, we demonstrate a novel application of this feature. Annotator inconsistency is a well-known

Table 6: Detailed results of different methods on the PPE Correctness benchmark.

Method	MMLU-Pro	MATH	GPQA	MBPP-Plus	IFEval	PPE Correctness			
Reported Results of Public Models									
Skywork-Reward-Gemma-2-27B	54.0	63.0	53.0	59.0	54.0	56.6			
DeepSeek-V2.5-0905	-	-	-	-	-	58.5			
Gemini-1.5-Pro	-	-	-	-	-	59.8			
ArmoRM-8B-v0.1	66.0	71.0	57.0	54.0	58.0	61.2			
InternLM2-20B-Reward	68.0	70.0	57.0	58.0	62.0	63.0			
LLaMA-3.1-70b-Instruct	-	-	-	-	-	59.2			
Claude-3.5-sonnet	66.0	63.0	56.0	52.0	57.0	58.8			
Nemotron-4-340B-Reward	70.0	65.0	57.0	49.0	63.0	60.8			
GPT-40	-	-	-	-	-	57.6			
Reproduced Results of Baseline Methods From DeepSeek									
LLM-as-a-Judge	66.0	68.0	52.8	50.2	56.8	58.8			
DeepSeek-BTRM-27B	68.8	73.2	56.8	68.8	66.0	66.7			
CLoud-Gemma-2-27B	68.7	68.8	53.5	59.0	62.0	62.4			
DeepSeek-PairRM-27B	68.3	74.7	55.0	63.1	62.9	64.8			
DeepSeek-GRM-27B-RFT	64.8	68.7	55.5	49.0	60.2	59.6			
DeepSeek-GRM-27B	64.8	68.8	55.6	50.1	59.8	59.8			
	Res	sults of Our	· Method						
OPRM-Qwen2.5-7B	65.2	70.1	56.3	59.0	56.1	61.3			
OPRM-Qwen2.5-14B	66.7	70.7	57.1	67.4	59.5	64.3			
OPRM-Qwen2.5-32B	71.2	73.2	57.9	66.2	62.2	66.1			
OPRM-Qwen2.5-72B	73.4	75.9	58.6	54.1	59.5	64.3			
Re	Results of Our Method (w/ Region Flooding Tuning)								
OPRM-RgFT-Qwen2.5-7B	64.8	71.6	55.9	63.0	56.8	62.4			
OPRM-RgFT-Qwen2.5-14B	69.5	74.0	57.3	67.0	60.0	65.6			
OPRM-RgFT-Qwen2.5-32B	73.3	76.8	58.5	67.2	60.6	67.3			
OPRM-RgFT-Qwen2.5-72B	72.8	77.1	59.0	62.0	61.2	66.4			

Table 7: Detailed results of different methods on the RMB benchmark.

Method	Helpfulness BoN	Helpfulness Pair	Harmlessness BoN	Harmlessness Pair	RMB				
Reported Results of Public Models									
Skywork-Reward-Gemma-2-27B	47.2	65.3	56.1	72.1	60.2				
DeepSeek-V2.5-0905	-	-	-	-	65.7				
Gemini-1.5-Pro	53.6	76.3	29.9	66.1	56.5				
ArmoRM-8B-v0.1	63.6	78.7	49.7	66.3	64.6				
InternLM2-20B-Reward	58.5	76.3	49.9	67.0	62.9				
LLaMA-3.1-70b-Instruct	64.8	81.1	55.8	73.9	68.9				
Claude-3.5-sonnet	70.5	83.8	51.8	76.4	70.6				
Nemotron-4-340B-Reward	-	-	-	-	69.9				
GPT-40	63.9	81.5	68.2	81.4	73.8				
	Reproduced Results	s of Baseline Method.	s From DeepSeek						
LLM-as-a-Judge	55.8	78.5	50.8	73.9	64.8				
DeepSeek-BTRM-27B	64.0	83.0	33.6	51.0	57.9				
CLoud-Gemma-2-27B	64.7	81.1	41.7	66.1	63.4				
DeepSeek-PairRM-27B	59.9	83.3	34.1	55.5	58.2				
DeepSeek-GRM-27B-RFT	58.4	79.3	54.2	76.0	67.0				
DeepSeek-GRM-27B	62.3	80.5	57.0	76.1	69.0				
	Re	sults of Our Method							
OPRM-Qwen2.5-7B	63.1	78.4	65.7	78.8	71.5				
OPRM-Qwen2.5-14B	65.8	80.7	68.2	80.5	73.8				
OPRM-Qwen2.5-32B	69.2	82.1	68.9	82.0	75.6				
OPRM-Qwen2.5-72B	68.7	82.4	64.2	78.5	73.5				
	Results of Our Method (w/ Region Flooding Tuning)								
OPRM-RgFT-Qwen2.5-7B	63.4	79.0	62.4	75.6	70.1				
OPRM-RgFT-Qwen2.5-14B	66.3	81.3	65.3	78.2	72.8				
OPRM-RgFT-Qwen2.5-32B	67.6	81.4	69.1	81.2	74.8				
OPRM-RgFT-Qwen2.5-72B	67.9	82.3	67.1	79.5	74.2				

challenge in preference data collection, an issue that is exacerbated at finer annotation granularities and leads to increased label ambiguity. To address this, we propose a method that explicitly models this ambiguity by permitting a single prompt-response pair to be associated with multiple potential quality levels. For instance, a response might be considered both **good** and **normal**, or even all three levels in cases of extreme uncertainty. As illustrated in Figure 5, our approach handles this by

optimizing over an expanded set of joint probabilities, corresponding to all plausible quality-level assignments for a given pair.

I.2 Customized Decoding Method

The full probability distribution $p_{\psi}(s \mid x, y)$ produced by OPRM allows us to go beyond a simple expected score. To fully leverage this rich distributional information, we introduce **Uncertainty-aware Decoding**. This method adjusts the expected score by penalizing predictive uncertainty, thereby favoring responses that are predicted to be high-quality with high confidence. The final reward score $r_{\psi}(x,y)$ is calculated as:

$$r_{\psi}(x,y) = \underbrace{\sum_{s=a}^{b} s \cdot p_{\psi}(s \mid x, y)}_{\text{Expected Score}} - \underbrace{\lambda \cdot u(x, y)}_{\text{Uncertainty Term}}$$
(22)

where the first term is the standard expected score. The second term, u(x,y), is an uncertainty measure of the distribution, such as its **Shannon entropy** or **variance**. The hyperparameter $\lambda \geq 0$ controls the strength of the uncertainty penalty.

J CASE STUDY

We present several case studies in Table 8, Table 9, and Table 10, along with the scoring results from BTRM, OPRM, and OPRM-RgFT for comparative analysis. A key finding is that by being grounded in absolute quality assessment, OPRM and OPRM-RgFT exhibit a reduced susceptibility to length bias, a known issue in purely relative models like BTRM. Moreover, for tasks with verifiable answers, our models consistently prioritize factual correctness, whereas BTRM can be misled by other stylistic attributes.

1188 1189 1190 1191 1192 1193 1194 Table 8: Case 1 1195 1196 **Example** 1197 Question: You are a highly skilled mathematician tasked with solving complex math problems. Your goal is to 1198 provide clear, step-by-step solutions that can be easily parsed and evaluated. 1199 Here is the math problem you need to solve: cyroblem> The length of a rectangle is twice its width. Given the length of the diagonal is $5\sqrt{5}$, find the area of the rectangle. </problem> 1201 Box your final answer using LaTeX, for example: $x = \sqrt{[Your final numerical or algebraic answer]}$ \$. 1202 Now, please solve the given math problem and provide your solution in the specified format. 1203 **Chosen Response**: Here's my solution to the problem: Let the width of the rectangle be w. Since the length is twice the width, the length of the rectangle is 2w. 1205 We are given that the length of the diagonal is $5\sqrt{5}$. Using the Pythagorean theorem, we can write: $(5\sqrt{5})^2 =$ 1206 $w^2 + (2w)^2$ 1207 Simplifying the equation, we get: $25 \cdot 5 = w^2 + 4w^2$ 1208 Combine like terms: $25 \cdot 5 = 5w^2$ 1209 Divide both sides by 5: $25 = w^2$ 1210 Take the square root of both sides: $w = \sqrt{25} = 5$ 1211 Now that we have the width, we can find the length: $2w = 2 \cdot 5 = 10$ The area of the rectangle is the product of the width and length: $A = w \cdot 2w = 5 \cdot 10 = |50|$ 1212 Therefore, the area of the rectangle is 50. 1213 1214 Rejected Response: To solve this problem, we can start by using the Pythagorean theorem to relate the length and width of the rectangle to the length of the diagonal: $d^2 = l^2 + w^2$ 1215 where d is the length of the diagonal, l is the length of the rectangle, and w is the width of the rectangle. Since 1216 the length of the rectangle is twice its width, we can write: l=2w1217 Substituting this into the Pythagorean theorem, we get: $d^2 = (2w)^2 + w^2$ 1218 Simplifying this expression, we get: $d^2 = 5w^2$ 1219 Since we are given that the length of the diagonal is $5\sqrt{5}$, we can set up the equation: $d = 5\sqrt{5}$ Substituting this into the equation above, we get: $5\sqrt{5}^2 = 5w^2$ Simplifying this expression, we get: $25 \cdot 5 = 5w^2$ 1221 Simplifying further, we get: $125 = w^2$ 1222 Taking the square root of both sides, we get: $w = \sqrt{125} = 5\sqrt{5}$ 1223 Now that we have found the width, we can find the length by substituting this value back into the equation: 1224 $l = 2w = 2(5\sqrt{5}) = 10\sqrt{5}$ 1225 The area of the rectangle is given by: $A = lw = (10\sqrt{5})(5\sqrt{5}) = 50\sqrt{5}^2 = 50 \cdot 5 = 250$ 1226 Final Answer: The final answer is 250. 1227 BTRM-Qwen2.5-32B 1228 Chosen Score: 0.03 1229 Rejected Score: -2.70 1230 OPRM-Qwen2.5-32B (Ours) 1231 Chosen Score: 7.04 1232 Rejected Score: 3.63 1233 OPRM-RgFT-Qwen2.5-32B (Ours) 1234 Chosen Score: 8.06 Rejected Score: 3.13 1236 1237

1242 1243 1245 1246 1247 Table 9: Case 2 1248 1249 Example 1250 **Question**: is $(x^y)^z = (x^z)^y$? 1251 **Chosen Response:** Yes, the equation $(x^y)^z = (x^z)^y$ is indeed true, and it is a direct consequence of the 1252 properties of exponents in mathematics. 1253 To see why this is true, let's consider the equation on both sides: 1254 Left side: $(x^y)^z$ 1255 Right side: $(x^z)^y$ First, we can simplify the left side by applying the power of a power rule, which states that $(a^m)^n = a^{(mn)}$: 1256 Left side: $(x^y)^z = x^{(yz)}$ 1257 Now, let's simplify the right side using the same power of a power rule: 1258 Right side: $(x^z)^y = x^{(zy)}$ 1259 Now we can see that both sides of the equation are equal: 1260 $x^{(yz)} = x^{(zy)}$ 1261 Since the bases (x) are the same and the exponents (yz) and (zy) are also the same, we can conclude that the 1262 equation $(x^y)^z = (x^z)^y$ is true. This property holds for any real numbers x, y, and z, as long as $x \neq 0$ and y1263 and z are non-zero. 1264 Rejected Response: A classic question about exponentiation! 1265 Unfortunately, the answer is no, $(x^y)^z$ is not necessarily equal to $(x^z)^y$ in general. 1266 To see why, let's break it down: $(x^y)^z$ means "raise x to the power of y, and then raise the result to the power of z". This can be written as: 1267 $(x^y)^z = x^{(y*z)}$ 1268 On the other hand, $(x^z)^y$ means "raise x to the power of z, and then raise the result to the power of y". This can 1269 be written as: 1270 $(x^z)^y = x^{(z*y)}$ 1271 As you can see, the two expressions are not equal in general, because y * z is not necessarily equal to z * y. For 1272 example, if x = 2, y = 3, and z = 4, then: $(2^3)^4 = 2^{(3*4)} = 2^{12} = 4096$ 1273 while $(2^4)^3 = 2^{(4*3)} = 2^{12} = 4096$ 1275 But in this case, it just so happens that y * z = z * y, so the two expressions are equal. However, this is not 1276 always the case. 1277 For example, if x = 2, y = 2, and z = 3, then: $(2^2)^3 = 2^{(2*3)} = 2^6 = 64$ 1278 while 1279 $(2^3)^2 = 2^{(3*2)} = 2^6 = 64$ 1280 Again, it just so happens that the two expressions are equal in this case, but this is not a general rule. 1281 So, to summarize: $(x^y)^z$ is not necessarily equal to $(x^z)^y$, unless y and z commute, i.e., y*z=z*y. 1282 BTRM-Qwen2.5-32B 1283 Chosen Score: -0.96 1284 Rejected Score: -0.39 1285 OPRM-Qwen2.5-32B (Ours) 1286 Chosen Score: 7.85 1287 Rejected Score: 3.84 1288 OPRM-RgFT-Qwen2.5-32B (Ours) 1289 Chosen Score: 8.03 1290 Rejected Score: 2.05 1291

1344

1297 1298 1299 1300 1301 1302 1303 1304 1305 1306 Table 10: Case 3 1307 1308 **Example** 1309 **Question**: Anita ____ some application letters to the company but they never call her. 1310 A.has send 1311 B.have sent 1312 C.have send 1313 D.have sending 1314 E.has sent 1315 Chosen Response: The correct answer is E. "has sent". 1316 Explanation: 1317 In this sentence, we need to use the present perfect tense, which is formed by "has/have + past participle". The past participle of the verb "send" is "sent". 1318 A. "has send" is incorrect because "send" is the base form of the verb, not the past participle. 1319 B. "have sent" is incorrect because the subject "Anita" is singular, so we should use "has" instead of "have". 1320 C. "have send" is incorrect for two reasons: "have" is used with a plural subject, and "send" is not the past 1321 participle. 1322 D. "have sending" is incorrect because "sending" is the present participle, not the past participle. E. "has sent" is the correct form, using "has" for the singular subject "Anita" and the past participle "sent". 1323 Therefore, the correct sentence should be: "Anita has sent some application letters to the company, but they 1324 never call her." 1325 **Rejected Response**: The correct option is E. "has sent." 1326 The sentence should read: "Anita has sent some application letters to the company but they never call her." 1327 The verb "has sent" is used because the subject of the sentence, "Anita," is singular and the action of sending the 1328 letters happened in the past. The present perfect tense is used to indicate that the action was completed in the 1329 past, but its effects are still being felt in the present. Therefore, "has sent" is the appropriate form of the verb to 1330 use in this sentence. 1331 BTRM-Owen2.5-32B 1332 Chosen Score: 0.24 Rejected Score: -1.12 1333 1334 OPRM-Qwen2.5-32B (Ours) Chosen Score: 8.08 1335 Rejected Score: 6.76 1336 OPRM-RgFT-Qwen2.5-32B (Ours) 1337 Chosen Score: 8.81 1338 Rejected Score: 8.14 1339 1340 1341 1343