

---

# Pre-Trained Policy Discriminators are General Reward Models

---

Shihan Dou<sup>1,2\*</sup>, Shichun Liu<sup>1,2\*†</sup>, Yuming Yang<sup>1,2\*</sup>, Yicheng Zou<sup>1\*†</sup>, Yunhua Zhou<sup>1</sup>, Shuhao Xing<sup>1</sup>, Chenhao Huang<sup>2</sup>, Qiming Ge<sup>1</sup>, Demin Song<sup>1</sup>, Haijun Lv<sup>1</sup>, Songyang Gao<sup>1</sup>, Chengqi Lv<sup>1</sup>, Enyu Zhou<sup>2</sup>, Honglin Guo<sup>2</sup>, Zhiheng Xi<sup>2</sup>, Wenwei Zhang<sup>1</sup>, Qipeng Guo<sup>1</sup>, Qi Zhang<sup>2</sup>, Xipeng Qiu<sup>2</sup>, Xuanjing Huang<sup>2</sup>, Tao Gui<sup>2†</sup>, Kai Chen<sup>1†</sup>

<sup>1</sup>**Shanghai AI Laboratory**, <sup>2</sup>**Fudan University**  
 {zouyicheng, chenkai}@pjlab.org.cn, tgui@fudan.edu.cn

Q <https://github.com/InternLM/POLAR>

## Abstract

We offer a novel perspective on reward modeling by formulating it as a policy discriminator, which quantifies the difference between two policies to generate a reward signal, guiding the training policy towards a target policy with desired behaviors. Based on this conceptual insight, we propose a scalable pre-training method named **POL**icy DiscriminAtive LeaRning (**POLAR**), which trains a reward model (RM) to discern identical policies and discriminate different ones. Unlike traditional reward modeling methods relying on absolute preferences, POLAR captures the relative difference between one policy and an arbitrary target policy, which is a scalable, high-level optimization objective suitable for modeling generic ranking relationships. Leveraging the POLAR pre-training paradigm, we present a series of RMs with parameter scales from 1.8B to 7B. Empirical results show that POLAR substantially outperforms traditional non-pre-trained methods, significantly enhancing RM performance. For instance, POLAR-7B could improve preference accuracy from 54.8% to 81.0% on STEM tasks and from 57.9% to 85.5% on creative writing tasks compared to SOTA baselines. POLAR also shows robust generalization capabilities in RLHF using Reinforcement Fine-tuning (RFT), providing reliable reward signals and markedly enhancing policy performance—improving LLaMa3.1-8B from an average of 47.36% to 56.33% and Qwen2.5-32B from 64.49% to 70.47% on 20 benchmarks. Moreover, scaling experiments reveal a clear power-law relationship between computation and performance, supported by linear correlation coefficients approaching 0.99. The impressive performance, strong generalization, and scaling properties suggest that POLAR is a promising direction for developing general and strong reward models.

## 1 Introduction

Reinforcement learning (RL) plays a crucial role in the post-training of large language models (LLMs) [126; 78; 4]. Its success hinges on the reward model’s (RM) ability to provide precise and stable feedback to the policy model [107; 28]. Although recent approaches successfully leverage labeled preference pairs to train RMs for alignment with human preferences, these methods often face challenges in terms of scalability and generalization [116; 64; 72; 49; 105]. The former is limited by the difficulty of acquiring large volumes of high-quality labeled pairs [23; 21], while the latter stems from the fact that this subjective approach to modeling human preferences makes RMs vulnerable to reward hacking [15; 12; 120]. On the other hand, several works, such as DeepSeek’s R1 [34], utilize

---

\*Equal contributions. †Corresponding authors. ‡Work done during an internship at Shanghai AI Laboratory.

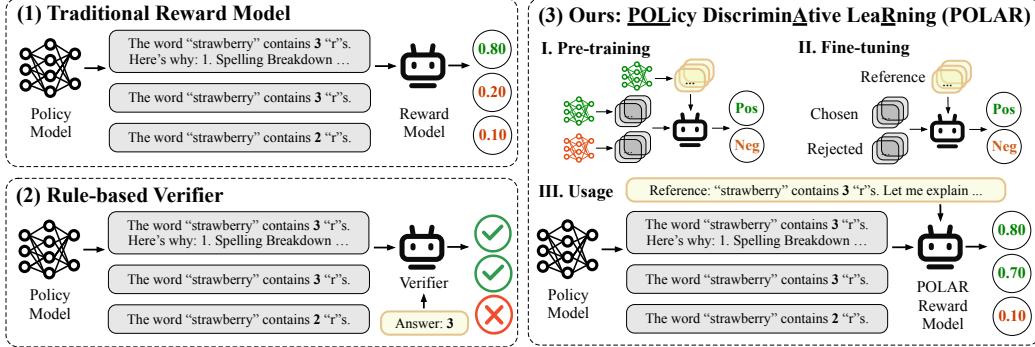


Figure 1: Comparison of three reward modeling methods: (1) traditional methods incorporate absolute preferences into RMs, which directly assess the quality of trajectories; (2) rule-based verifier validate the candidate trajectory through the gold answer and predefined rules; (3) POLAR pre-trains an RM to recognize identical policies and discriminate different ones, enabling it to measure the difference in trajectories between a training policy and a target policy with desired behaviors.

rules to verify the correctness of model outputs and provide accurate reward signals for RL. However, these rule-based verifiers can only be applied in scenarios where model outputs can be automatically verified by pre-defined rules, such as in reasoning and coding tasks [30; 110; 63; 81]. In open-ended domains like writing and translation, rules are usually complicated and difficult to design in advance, making the rule-based verifiers hard to extend to general tasks.

Before delving into reward modeling, it is instructive to revisit the widespread success of LLMs. By adopting a unified Next Token Prediction (NTP) optimization target [85; 86; 9], LLMs effectively harmonize diverse NLP tasks under a common objective, addressing the challenge of cross-task generalization. This inspires us to reconsider the training paradigm of RMs. Traditional RMs heavily rely on absolute, manually-defined criteria to generate a preference score. Analogous to how LLMs unified NLP tasks, we should identify a fundamental, criterion-agnostic objective for RM pre-training.

Instead of traditional absolute preference modeling, we propose redefining a reward model as a “policy discriminator”. Specifically, by quantifying the difference between candidate policies and a given target policy, we establish a criterion-agnostic objective, which naturally assigns higher scores to policies that are more “similar” to the desired target policy. This reward signal could guide the training policy toward desired behaviors during RL. Furthermore, since target policies can be arbitrarily chosen, this objective eliminates reliance on manually defined preferences and is applicable to any scenario, thus offering a scalable and fundamental pre-training paradigm for RMs. We refer to this training objective as Policy Discriminative Learning (POLAR), as illustrated in Figure 1.

Starting from a diverse and extensive collection of policy models, we construct a large-scale synthetic corpus by sampling trajectories from these policies. We then formulate the pre-training task as a contrastive learning objective utilizing Bradley-Terry (BT) loss [8], which encourages the RM to recognize trajectories derived from identical policies while distinguishing those originating from different ones. Consequently, the pre-trained RM learns to assign higher rewards to trajectories exhibiting greater consistency with the target policy and generalizes this discrimination capability to unseen policies. After pre-training, analogous to supervised fine-tuning (SFT) in LLMs [78; 85] that enables rapid adaptation to specific tasks and instructions, we also introduce an SFT procedure for POLAR RMs tailored to align with human-defined criteria. This fine-tuning process requires only a small set of reference trajectories generated from a given target policy, accompanied by candidate trajectories annotated with ranking labels reflecting their difference relative to the target. The reference trajectories can be directly annotated by humans—or alternatively, they could be derived from high-performing LLMs. Such a flexible and lightweight fine-tuning approach allows the RM to rapidly adapt to new domains or criteria.

Leveraging the POLAR pre-training method, we present a series of reward models with parameter scales ranging from 1.8 to 7 billion. Our empirical results demonstrate that POLAR effectively discriminates among diverse policies, and the fine-tuned POLAR RMs could significantly outperform traditional preference modeling methods. Specifically, in preference evaluation tasks, POLAR-7B surpasses the SOTA 72B-parameter WorldPM [108], achieving an average improvement of 5.8%

points despite being approximately  $10\times$  smaller. Additionally, when applied within RLHF using Reinforcement Fine-tuning (RFT) [68; 77], POLAR RMs deliver more accurate reward signals and exhibit superior generalization across various downstream tasks, substantially enhancing the performance of popular policy models such as Qwen2.5 [117], LLaMa3.1 [26], and InternLM3 [11]. Moreover, POLAR exhibits scaling laws similar to those observed in LLMs, highlighting its significant potential for developing increasingly powerful reward models. In summary, our main contributions are as follows:

1. We propose POLAR, a novel criterion-agnostic pre-training paradigm for reward modeling based on a scalable training objective—policy discrimination.
2. Scaling experiments reveal promising scaling laws, highlighting the significant potential of POLAR for enhancing the upper bound of reward modeling and developing stronger and more generalizable RMs.
3. We developed the POLAR series of reward models. They substantially outperform traditional RMs in empirical evaluations, achieving higher preference accuracy and better generalization than considerably larger RMs. This advancement expands the potential and applicability of RL algorithms, such as RFT, paving the way for more innovative and diverse applications.

## 2 Related Work

**Reward Modeling in Reinforcement Learning** RL has emerged as a pivotal technique in the post-training phase of LLMs [52; 50; 106; 4; 82; 34; 78; 126; 121]. All RL approaches crucially depend on precise reward signals provided by the reward function. Reward functions broadly fall into model-based and rule-based categories, with our work aligning closely with the model-based reward function paradigm [28; 106; 51]. Existing model-based methods typically train RMs using labeled pairs to approximate a preference distribution [4; 89; 108]. However, the scarcity of extensive labeled preference pairs poses a significant challenge [23; 21; 5; 43]. Additionally, such methods often exhibit limited generalization, struggling to predict out-of-distribution preferences robustly, thus weakening their effectiveness in RL training [15; 12; 120]. Another line of research investigates generative reward models, leveraging LLMs themselves as verifiers [125; 42; 70; 123; 66; 14; 35]. Yet, these approaches inherently reflect the biases and preferences of the models used, frequently resulting in unreliable or biased reward signals [33]. Alternatively, some studies utilize rule-based verifiers to provide feedback [112; 63; 34; 102], but this approach struggles to address general tasks that are challenging to verify through predefined rules.

**Pre-training and Scaling Laws** Our work closely relates to unsupervised learning, a critical component of pre-training focused on deriving essential and generalizable knowledge from extensive unlabeled datasets [80; 9; 86; 36; 79; 20; 88; 85]. Unsupervised pre-training is foundational to modern advances in language modeling, powering the latest sophisticated AI systems [84; 24; 111; 54; 109; 67]. Concurrently, scaling laws describe the power-law relationship among computation and performance in neural language models, offering crucial insights into unsupervised pre-training [48; 10; 44; 39; 92; 114; 7]. Extensive research indicates that scaling up data size and model size consistently enhances LLM performance [26; 19; 83; 74; 3; 37; 85; 48; 124]. Recent studies further demonstrate the capability of scaling laws to predict the performance of larger models based on smaller ones, enabling efficient resource allocation [92; 44; 39]. The success of scaling laws in language modeling has significantly influenced theoretical understandings and practical advancements, inspiring similar explorations in vision [87; 94; 98; 62], multimodal learning [62; 1; 22; 29; 104], and reinforcement learning [40; 73]. Gao et al. [28] specifically expanded scaling laws into the over-optimization of preference-based reward models, further broadening their applicability in RMs.

## 3 Method

In this section, we first revisit the optimization objectives of RL in LLMs and show that the RM essentially functions as a policy discriminator. Building on this insight, we introduce the Policy Discriminative Learning (POLAR) as illustrated in Figure 2. During the first stage, the RM acquires the ability to discriminate among various policies and quantify their differences. Subsequently, in the supervised fine-tuning stage, we fine-tune the RM on trajectory ranking data derived from human judgments, explicitly capturing preferences regarding policy behaviors.

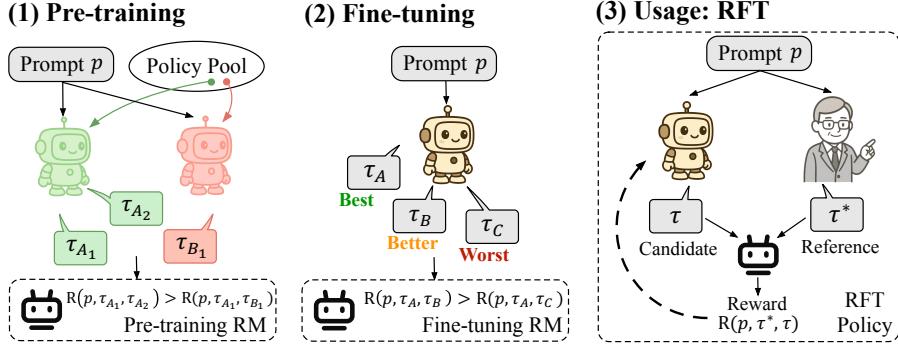


Figure 2: Overview of Policy Discriminative Learning (POLAR). **Stage 1:** In pre-training, the RM learns criterion-agnostic policy differences by assigning higher rewards to trajectory pairs from consistent policies. **Stage 2:** During fine-tuning, human annotators rank trajectories from the same policy, implicitly defining human criteria, to align RM evaluations with human standards. **Usage:** In Reinforcement Fine-Tuning (RFT), the fine-tuned RM provides reward signals comparing candidate trajectories with human-preferred references, guiding policy training toward desired behaviors.

### 3.1 Preliminary: Policy Optimization in Reinforcement Learning

Consider a training policy model  $\pi_\phi$  parameterized by  $\phi$  and a prompt distribution  $\mathcal{D}_x$ . Let  $\tau$  denote trajectories generated by this policy given prompts sampled from  $\mathcal{D}_x$ . The resulting policy distribution  $\mathcal{D}_{\pi_\phi}$  can thus be empirically approximated by sampling prompt-trajectory pairs  $(x, \tau)$ . Let  $r_\theta : (x, \tau) \rightarrow \mathbb{R}$  denote the reward model parameterized by  $\theta$ , which assesses policy performance by assigning scores to prompt-trajectory pairs  $(x, \tau)$ . Given a regularization coefficient  $\beta$  that controls the strength of the KL-divergence penalty relative to the initial policy  $\pi_{\text{init}}$ , the RL optimization objective can be formulated as follows:

$$\mathcal{O}_{\text{RL}}(\pi) = \mathbb{E}_{(x, \tau) \sim \mathcal{D}_\pi} [r_\theta(x, \tau) - \beta D_{\text{KL}}(\pi(\tau|x) || \pi_{\text{init}}(\tau|x))]. \quad (1)$$

The optimal policy  $\pi^*$  then admits a closed-form solution:

$$\pi^*(\tau|x) = \frac{1}{Z(x)} \pi_{\text{init}}(\tau|x) \exp\left(\frac{r_\theta(x, \tau)}{\beta}\right), \quad (2)$$

where  $Z(\cdot)$  is the partition function. This formulation provides a key insight: the reward model implicitly encodes a continuous operator that maps an initial policy to its optimal form through RL optimization, captured by the reward scores. Consequently, when a policy  $\pi_\phi$  is trained against this reward signal under KL constraints, it effectively learns to approximate this implicit mapping within the defined reward distribution.

### 3.2 Unsupervised Reward Pre-training via Distributional Alignment

While the Bradley-Terry framework [8] facilitates reward modeling based on explicit pairwise comparisons (e.g., harmlessness in safety alignment), it inherently presupposes an absolute human criterion. However, this assumption becomes problematic when considering broader, more general classes of criteria. Specifically, given an arbitrary criterion  $p$  inducing a partial ordering over responses, i.e.,  $\{r(a_1) \geq r(a_2) \geq \dots \geq r(a_n)\}$  for responses  $a_i$ , there always exists a complementary criterion  $\neg p$  that reverses this ordering. Consequently, exhaustively enumerating all possible criteria is not only computationally infeasible but also theoretically ill-posed.

To overcome this challenge, we propose an unsupervised reward pre-training paradigm that provides a criterion-agnostic initialization for the reward model. Specifically, let  $\pi^*$  and  $\pi_{\text{init}}$  represent two policy distributions, where  $\pi^*$  denotes the optimal policy for a given downstream task and  $\pi_{\text{init}}$  denotes the initial policy. Under the KL-constrained RL objective  $\mathcal{O}_{\text{RL}}$ , we observe that the reward function can be uniquely defined (up to an additive constant) by the density ratio between  $\pi^*$  and  $\pi_{\text{init}}$ , resulting in the following relationship:

$$r_\theta(x, \tau) \triangleq \beta \log \frac{\pi^*(\tau|x)}{\pi_{\text{init}}(\tau|x)} + \beta \log Z(x). \quad (3)$$

Since the partition function is a constant independent of trajectories, we focus only on the former item, whose expected reward spans the entire trajectory space  $\tau$ :

$$\mathbb{E}_{\tau \sim \pi^*(\cdot|x)} [r_\theta(x, \tau)] \sim \underbrace{\beta \mathbb{E}_{\pi^*} \left[ \log \frac{\pi^*}{\pi_{\text{init}}} \right]}_{D_{\text{KL}}(\pi^* \parallel \pi_{\text{init}})}. \quad (4)$$

This term converges to zero when the sampling distribution precisely matches the optimal distribution. Motivated by this observation, we replace direct reward regression with a distributional distance minimization approach. Given prior knowledge of desired behaviors  $\tau^*$  (e.g., human references), we train  $r_\theta$  to measure the divergence between  $\pi^*$  (approximated through  $\tau^*$ ) and  $\pi_{\text{init}}$  (approximated through policy rollouts).

Under this perspective, the role of the RM fundamentally shifts: instead of merely assessing the performance of individual trajectories, it now serves as a measure of policy differentiation, quantifying differences between the training policy and the desired target policy. In practice, given two different trajectories generated from the same prompt, the pre-trained RM estimates the divergence between their underlying sampling policies, which can be formalized as follows:

$$D(\pi_\phi, \pi^*) = \mathbb{E}_{(x, \tau) \sim \mathcal{D}_{\pi_\phi}, (\tau^*) \sim \mathcal{D}_{\pi^*}} [r_\theta(\tau, \tau^*|x)]. \quad (5)$$

This perspective naturally suggests viewing the reward model as a **policy discriminator**: it learns to distinguish between the training and target policies, quantifying their degree of difference. A smaller difference yields a larger assigned reward, thus incentivizing the training policy to progressively align more closely with the desired target policy through RL optimization. Here, we employ Bradley-Terry (BT) loss [8; 97] to pre-train the reward model:

$$\mathcal{L}_{\text{pre-train}}(\theta) = -\mathbb{E}_{(p, \tau_{A_1}, \tau_{A_2}, \tau_{B_1}) \sim \mathcal{D}_{\text{pre-train}}} [\log \sigma(r_\theta(p, \tau_{A_1}, \tau_{A_2}) - r_\theta(p, \tau_{A_1}, \tau_{B_1}))]. \quad (6)$$

Here  $p$  denotes a prompt sampled from the pre-training dataset  $\mathcal{D}_{\text{pre-train}}$ .  $\tau_{A_1}$  and  $\tau_{A_2}$  are generated by the same policy, while  $\tau_{B_1}$  originates from a different policy. These policies are randomly selected from a diverse policy pool comprising LLMs varying in architectures and parameters. The pre-training objective encourages the RM to capture nuanced distinctions between policies, assigning higher rewards to trajectory pairs drawn from more closely aligned or “similar” policies.

### 3.3 Supervised Fine-tuning with Human Criteria

After unsupervised pre-training, the reward model acquires a criterion-agnostic capability to discriminate between different policies. For practical usage, it is essential to align this discriminative ability with human-defined standards and preferences. To this end, we introduce a supervised fine-tuning stage, which is designed to efficiently adapt the pre-trained RM to human judgment criteria.

Ideally, this fine-tuning stage would utilize reference trajectories generated by a desired target policy, accompanied by candidate trajectories annotated with ranking labels indicating their relative differences. In practice, however, we adopt a simplified yet effective approach: given a prompt  $p$  from downstream tasks, we generate three distinct trajectories from the same policy. Human annotators then rank these trajectories from best to worst, denoted as  $(\tau_A \succ \tau_B \succ \tau_C)$ , where  $\tau_A$  is the most preferred,  $\tau_B$  is the second-best, and  $\tau_C$  is the least preferred. Although these trajectories originate from a single policy, the introduction of human evaluations implicitly imposes different human criteria, effectively treating these trajectories as if they were drawn from different underlying policies that reflect varying degrees of alignment with human standards. During fine-tuning, we employ the following supervised ranking loss, consistent with the Bradley-Terry framework:

$$\mathcal{L}_{\text{fine-tune}}(\theta) = -\mathbb{E}_{(p, \tau_A, \tau_B, \tau_C) \sim \mathcal{D}_{\text{fine-tune}}} [\log \sigma(r_\theta(p, \tau_A, \tau_B) - r_\theta(p, \tau_A, \tau_C))], \quad (7)$$

where  $\tau_A$  denotes the trajectory ranked highest by humans and  $\tau_C$  is ranked lowest. By optimizing this loss, the reward model quickly adapts to human preferences, effectively measuring the policy differences implied by human judgments, thereby closely aligning its policy discriminative capabilities with human evaluation criteria. This SFT stage effectively bridges the gap between the criterion-agnostic discrimination learned during pre-training and human-aligned evaluations, yielding a reward model robustly generalized to human judgment scenarios with minimal annotation overhead.

## 4 Experiments

We train reward models with parameter sizes of 1.8B and 7B using POLAR, denoted as POLAR-1.8B and POLAR-7B, respectively, and benchmark their performance against SOTA baseline RMs. In Section 4.1, we specify the model and training details. Section 4.2 presents the evaluation results on preference modeling benchmarks. In Section 4.3, we demonstrate the efficacy of POLAR within the RLHF framework. To further highlight POLAR’s potential and scalability, we scale up the model sizes and analyze the corresponding scaling laws in Section 4.4. Lastly, Section 4.5 presents comprehensive ablation studies to validate the contributions of different training stages in POLAR.

### 4.1 Model and Training Details

**Model Architecture** The architecture of POLAR RMs are based on an autoregressive Transformer [85; 86; 9], similar to models in the GPT series [85; 86], augmented with a linear prediction head—a common design adopted in traditional preference-based reward modeling methods [61; 11]. Specifically, the decoder configuration used in POLAR RMs aligns with that of the InternLM-2.5 series [11]. For traditional RMs, the prompt and trajectory are concatenated to form the model input, from which the model directly outputs a reward value. In contrast, as outlined in our training objective, POLAR takes a prompt along with two trajectories—a reference and a candidate—as input. We utilize special tokens to combine these three elements into a single input sequence:

prompt + reference <split\_token> prompt + candidate <reward\_token>

The linear head then processes the hidden state corresponding to the <reward\_token> token from the model’s final layer to produce the reward value.

**Training Details** During the pre-training stage, we follow the standard process widely adopted in LLM pre-training [11; 60; 117] and conduct extensive experiments on hyperparameters, deriving clear scaling laws for hyperparameter configurations. Appendix A.1.2 provides an in-depth analysis of how hyperparameters are chosen for POLAR. Training details for the pre-training and supervised fine-tuning stages, including data composition, are provided in Appendices A.1 and A.2, respectively.

### 4.2 Performance in Preference Evaluation

We first evaluate POLAR on the human preference prediction task, which measures the ability of RMs to accurately identify responses preferred by humans.

**Evaluation Setup** We compare POLAR with 5 SOTA RM baselines: InternLM2-Reward-7B [11], InternLM2-Reward-20B [11], Skywork-Reward-8B [61], Skywork-Reward-27B [61], and WorldPM-72B-UltraFeedback [108] (details in Appendix B.1). Our primary evaluation uses the RMB benchmark [127], containing 3,162 questions, each with multiple trajectories ranked by preference scores. The top-ranked trajectories are treated as references, representing samples drawn from a target policy. The task is to identify whether RMs correctly prefer the second-ranked trajectory over the third-ranked one. Additionally, we create another evaluation set from real user queries collected through online platforms and manually annotate trajectory rankings (see Appendix B.2). We carefully remove overlaps with the training data to maintain independence. Existing RM baselines typically assess trajectories without considering references. To ensure fairness, we evaluate baselines using two methods: (a) standard scoring without references, and (b) including references explicitly in the prompt (Figure 10). We report the best results for each baseline across these two settings.

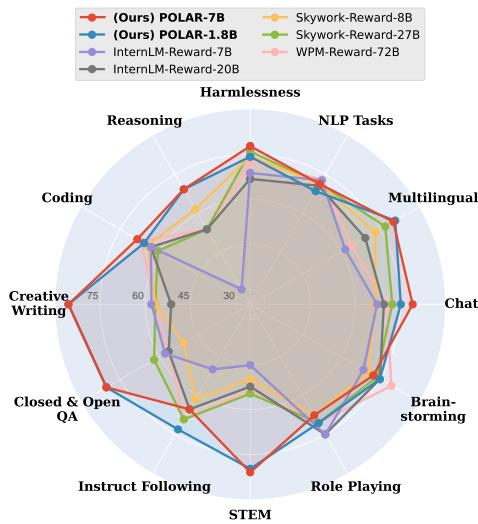


Figure 3: Comparison of POLAR and baselines on human preference prediction.

Table 1: Reward model performance comparison in RLHF training. **Baseline** denotes the initial policy model without RLHF. We compare POLAR against SOTA reward models across 20 benchmarks within the RLHF framework. Complete results are detailed in Tables 7, 8, 9, and 10.

Policy Model	Reward Model	General Task	Instruct Following	Coding	General Reasoning	Math	Knowledge	Average
InternLM3-8B-Instruct	Baseline	24.07	62.65	74.40	64.37	83.11	60.94	56.49
	InternLM2-Reward-7B	28.02	64.45	78.63	64.84	79.96	60.43	57.82
	Skywork-Reward-8B	29.21	63.75	74.66	64.82	83.36	59.95	57.92
	InternLM2-Reward-20B	28.76	66.75	74.16	64.97	82.20	60.65	58.09
	Skywork-Reward-27B	30.20	64.95	74.35	65.18	83.23	59.91	58.34
	WorldPM-72B-UltraFeedback	34.89	67.90	77.13	65.56	84.29	61.08	60.49
	POLAR-1.8B (Ours)	<b>37.50</b>	72.70	78.24	66.79	84.33	64.40	62.60
	POLAR-7B (Ours)	37.35	<b>73.25</b>	<b>79.63</b>	<b>67.89</b>	<b>85.18</b>	<b>64.46</b>	<b>63.18</b>
Llama-3.1-8B-Instruct	Baseline	15.59	63.35	70.69	52.95	67.60	49.39	47.36
	InternLM2-Reward-7B	25.37	60.80	59.24	54.15	65.21	46.35	48.06
	Skywork-Reward-8B	24.80	61.80	67.53	53.54	66.23	49.36	49.22
	InternLM2-Reward-20B	26.52	62.85	58.57	52.41	64.45	45.09	47.70
	Skywork-Reward-27B	24.57	61.70	66.34	54.58	66.25	49.97	49.44
	WorldPM-72B-UltraFeedback	21.36	63.85	70.86	54.74	69.56	49.70	49.64
	POLAR-1.8B (Ours)	27.96	65.20	71.35	57.52	71.11	51.30	52.71
	POLAR-7B (Ours)	<b>37.02</b>	<b>69.30</b>	<b>72.14</b>	<b>59.85</b>	<b>72.20</b>	<b>51.69</b>	<b>56.33</b>
Qwen2.5-7B-Instruct	Baseline	26.52	66.05	79.24	53.83	83.47	61.98	54.95
	InternLM2-Reward-7B	31.99	64.05	72.80	56.48	80.35	55.24	54.95
	Skywork-Reward-8B	32.44	68.00	76.71	58.09	83.13	58.12	57.04
	InternLM2-Reward-20B	33.05	68.40	74.06	55.41	82.62	58.36	56.15
	Skywork-Reward-27B	34.28	69.45	78.21	57.46	83.58	59.43	57.84
	WorldPM-72B-UltraFeedback	35.72	70.55	77.48	59.48	83.35	59.48	58.83
	POLAR-1.8B (Ours)	35.76	<b>71.35</b>	80.40	<b>62.52</b>	84.19	60.39	60.35
	POLAR-7B (Ours)	<b>37.70</b>	71.15	<b>81.15</b>	61.30	<b>84.70</b>	<b>62.57</b>	<b>60.90</b>
Qwen2.5-32B-Instruct	Baseline	31.07	75.50	86.56	69.74	89.35	71.07	64.49
	InternLM2-Reward-7B	36.10	75.70	83.13	69.72	87.20	64.99	64.29
	Skywork-Reward-8B	36.44	79.60	84.42	71.07	89.29	68.77	66.08
	InternLM2-Reward-20B	37.98	74.45	85.76	69.32	89.35	66.68	65.25
	Skywork-Reward-27B	38.43	80.10	83.95	71.93	86.78	69.15	66.64
	WorldPM-72B-UltraFeedback	40.59	78.65	86.79	70.38	89.90	69.05	67.15
	POLAR-1.8B (Ours)	40.24	80.25	87.47	72.23	90.03	<b>73.67</b>	68.55
	POLAR-7B (Ours)	<b>45.98</b>	<b>80.50</b>	<b>88.92</b>	<b>73.17</b>	<b>90.39</b>	73.59	<b>70.47</b>

**Results Comparison** To thoroughly illustrate RM performance, we follow the evaluation approach from prior studies [51; 127], categorizing the evaluation sets by task type, as shown in Figure 3. POLAR exhibits outstanding generalization, consistently outperforming baseline RMs across most tasks. Notably, on the STEM task, POLAR-1.8B and POLAR-7B surpass the best baseline by over 24.9 and 26.2 percentage points, respectively. POLAR also accurately identifies subtle distinctions in trajectories for challenging tasks like Reasoning and general tasks such as Chat and Creative Writing, accurately predicting human preferences. The superior performance of POLAR can be attributed to its unique training paradigm. Unlike traditional RMs that rely on absolute preferences and thus perform well primarily within distribution, POLAR RMs leverage large-scale criterion-agnostic pre-training to learn nuanced differences between policies, resulting in robust out-of-distribution generalization. Notably, POLAR-1.8B achieves competitive results with just 1.8B parameters, which is only 1/15th of Skywork-Reward-27B and 1/40th of WorldPM-72B-UltraFeedback. These results highlight POLAR’s efficiency and scalability.

We observed that POLAR-1.8B and POLAR-7B exhibit similar performance in preference evaluations. However, in downstream RL experiments (see Section 4.3), POLAR-7B demonstrates a notable advantage over POLAR-1.8B. This discrepancy between preference evaluations and actual RL tasks highlights an important limitation in traditional evaluation methodologies. Specifically, standard preference evaluation datasets may inadequately reflect the full spectrum of capabilities and nuanced distinctions reward models possess [127; 106; 90], underscoring the need for more comprehensive and representative evaluation frameworks in future studies.

#### 4.3 Performance in RLHF Training

**Setup and Implementations** We select four open-source LLMs as policies: InternLM3-8B-Instruct [11], LLaMa-3.1-8B-Instruct [26], Qwen2.5-7B-Instruct [117], and Qwen2.5-32B-Instruct [117]. Policy optimization employs the Proximal Policy Optimization (PPO) algorithm [93]. Unlike traditional baseline RMs, which directly evaluate policy trajectories without references, POLAR requires a reference trajectory. During RL training, we assess trajectories by measuring their consistency with

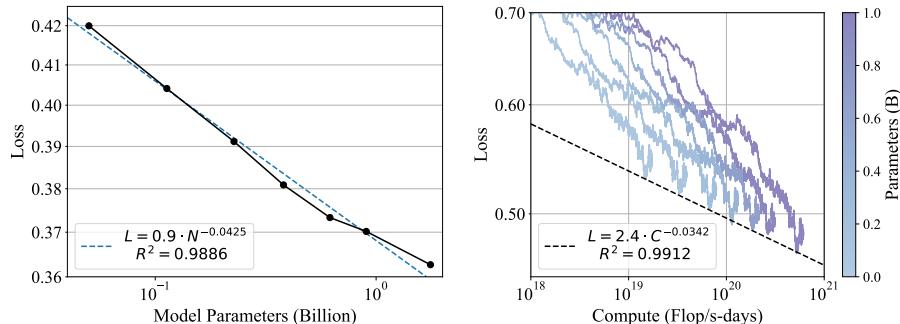


Figure 4: Scaling laws in POLAR. Validation loss vs. (left) model parameters  $N$  and (right) optimal training compute  $C$ . Dashed lines show the power-law fit, with  $R^2 = 0.9886$  (left) and  $R^2 = 0.9912$  (right). Results show a predictable decrease in validation loss as model size or compute increases.

the provided reference, adopting the Reinforcement Fine-Tuning (RFT). Comprehensive training details are available in Appendix B.3. Following previous studies [55; 11; 118], we evaluate policy performance on 20 popular benchmarks (Table 6) using OpenCompass [17], employing the optimal settings recommended by OpenCompass for fair comparison.

**Results Comparison** Tables 1, 7, 8, 9, and 10 summarize the results from RLHF experiments using POLAR and baseline RMs. POLAR RMs consistently outperform traditional non-pre-trained RMs. For instance, the Llama-3.1 fine-tuned using POLAR-7B achieves an average improvement of 9.0% over the initial policy model and 6.7% over the policy model optimized by WorldPM-72B-UltraFeedback across all benchmarks. These results align with our findings from preference evaluations. The enhanced generalization and reduced bias in POLAR’s reward signals primarily stem from its novel pre-training paradigm, which allows RMs to learn subtle distinctions between policies from extensive pre-training data rather than relying solely on labeled preference pairs. Moreover, the incorporation of reference trajectories further clarifies optimization objectives, making policy training directions more explicit and stable, thus reducing deviations and improving robustness across diverse and potentially out-of-distribution tasks. Notably, although POLAR-1.8B and POLAR-7B exhibit similar performance in preference evaluations, POLAR-7B demonstrates a significant advantage in downstream RL applications. The substantial performance gains from 1.8B to 7B parameters further illustrate the significant scaling effects achievable with the POLAR paradigm.

#### 4.4 Scaling Laws in POLAR

Previous studies [26; 85] have demonstrated that scaling up neural language models consistently reduces the validation loss  $L$  in a predictable manner. This reduction follows a power-law relationship relative to model parameters ( $N$ ), training tokens ( $D$ ), or computational resources ( $C$ ) [48; 10]:  $L = \beta \cdot X^\alpha$ , where  $X$  represents  $N$ ,  $D$ , or  $C$ . In this formulation,  $\alpha$  denotes the scaling exponent, indicating the rate at which validation loss decreases as scale  $X$  increases, while  $\beta$  is a normalization constant that determines the baseline level of the loss curve. Scaling laws thus facilitate the prediction of performance for larger models based on small-scale experiments, serving as an invaluable guide for efficient resource allocation and streamlined model development [92; 44; 39].

**Setup** We explore whether POLAR also exhibits scaling laws, which would underscore its scalability and potential. Specifically, we analyze how validation loss varies with respect to model parameters ( $N$ ) and optimal training compute ( $C$ ), provided that sufficient data is available. To this end, we train five RMs of different sizes, ranging from 50M to 1B parameters, using up to 54B training tokens.

**Empirical Results** We first investigate the scaling behavior of POLAR as the number of model parameters  $N$  increases. As shown in the left side of Figure 4, we observe a clear power-law relationship between validation loss and model size. The fitted scaling law is given by:

$$L = 0.9 \cdot N^{-0.0425}. \quad (8)$$

The high  $R^2$  value (i.e., 0.9886) of the fit indicates an excellent match between the empirical data and the power-law trend. These findings highlight the scalability of POLAR, showing that expanding the model size consistently leads to enhanced performance.

Table 2: Ablation study in RLHF training. **Baseline** denotes the initial policy model without RLHF, i.e., InternLM3-8B-Instruct [11]. **w/o PT** denotes the RM fine-tuned solely on human criteria, without any pre-training phase. **w/o PT & Ref** represents the RM trained via the traditional non-pre-trained method without reference trajectories. More detailed results are demonstrated in Table 12.

Reward Model	General Task	Instruct Following	Coding	General Reasoning	Math	Knowledge	Average
Baseline	24.07	62.65	74.40	64.37	83.11	60.94	56.49
1.8B	<b>POLAR</b>	37.50	72.70	78.24	66.79	84.33	64.40
	w/o PT	30.37	68.05	76.96	65.23	83.30	63.04
	w/o PT & Ref	27.67	63.65	75.49	64.89	82.50	60.36
7B	<b>POLAR</b>	37.35	73.25	79.63	67.89	85.18	64.46
	w/o PT	31.14	68.25	78.99	67.53	83.80	64.06
	w/o PT & Ref	31.07	68.30	76.74	66.16	83.05	62.06

We also examine the relationship between validation loss and optimal training compute  $C$ . For each model, we track the validation loss as a function of  $C$  throughout training. The right panel of Figure 4 shows that validation loss follows a power-law scaling trend with respect to compute:

$$L = 2.4 \cdot C^{-0.0342}. \quad (9)$$

The high  $R^2 = 0.9912$  of the fit indicates that POLAR reliably benefits from increased compute, further validating the applicability of scaling laws to our approach. This suggests that allocating more computational resources consistently yields better RM performance. Overall, these findings demonstrate that POLAR exhibits clear scaling laws similar to those observed in LLMs [40; 38]. Such predictable improvements underscore the strong potential of POLAR as a foundation for building more general reward models.

#### 4.5 Ablation Study

**Impact of POLAR Pre-training** We first investigate whether the strong performance of POLAR primarily stems from its pre-training stage. To do this, we train reward models under identical experimental settings but only utilizing the fine-tuning stage without the POLAR pre-training. Tables 11 and 2 report the performance results on preference evaluation and RLHF training, respectively. Results show these fine-tuning-only RMs achieve competitive results in preference evaluation but demonstrate a substantial performance decline in RLHF training. Consistent with prior studies [127; 106; 90], strong performance in preference evaluation does not guarantee effective reward signals for RLHF training. These observations underscore the critical role of RLHF effectiveness as an indicator of reward model quality, highlighting that POLAR’s pre-training stage significantly enhances both RM performance and generalization capability.

We also trained a traditional non-pre-trained reward model without reference trajectories using the same human criteria data. The experimental results show that, across nearly all tasks, the policy model fine-tuned with the RM **w/o PT** consistently outperforms the policy model fine-tuned with the RM **w/o PT & Ref**. This indicates that, even in the absence of pre-training, reference trajectories provide crucial guidance, simplifying the reward model’s evaluation task and thereby yielding clearer and more stable directions for policy training.

**RFT vs. SFT** In RLHF training, POLAR RMs utilize provided references to generate reward signals for policy training, a process referred to as RFT. To assess whether performance improvements stem specifically from POLAR or merely from the reference data, we directly fine-tune policies using a straightforward SFT approach on the same reference data, as shown in Table 3. Results reveal a significant performance decline for policies trained with SFT compared to those trained with POLAR RMs using RFT. This clearly demonstrates that RL exploits training data more robustly, and POLAR RMs serve effectively as graders to provide more accurate and robust supervision signals during the training process, reliably evaluating and guiding policy optimization.

## 5 Limitations and Future Work

**Pros and Cons of Reference Trajectories** On the one hand, the use of reference trajectories significantly enhances the accuracy and reliability of reward models. Ideally, incorporating multiple

Table 3: Ablation study on RFT vs. SFT. **RFT** denotes the reinforcement fine-tuning using POLAR. **SFT** denotes the straightforward supervised fine-tuning. Two training processes employ the same prompt-reference data. More detailed results are shown in Table 13.

Policy Model	Method	General Task	Instruct Following	Coding	General Reasoning	Math	Knowledge	Average
InternLM3-8B-Instruct	Baseline	24.07	62.65	74.40	64.37	83.11	60.94	56.49
	SFT	25.27	67.55	75.21	62.26	80.38	60.89	56.44
	RFT <sub>POLAR-7B</sub>	37.35	73.25	79.63	67.89	85.18	64.46	<b>63.18</b>
Qwen2.5-7B-Instruct	Baseline	26.52	66.05	79.24	53.83	83.47	61.98	54.95
	SFT	26.36	66.86	72.94	57.76	80.12	58.36	54.66
	RFT <sub>POLAR-7B</sub>	37.70	71.15	81.15	61.30	84.70	62.57	<b>60.90</b>
Qwen2.5-32B-Instruct	Baseline	31.07	75.45	86.56	69.74	89.35	71.07	64.49
	SFT	39.39	78.39	83.18	67.86	87.93	71.09	65.82
	RFT <sub>POLAR-7B</sub>	45.98	80.50	88.92	73.17	90.39	73.59	<b>70.47</b>

references per prompt could greatly reduce reward variance, particularly in open-ended scenarios. However, due to practical cost constraints, we currently provide only one single reference trajectory for each prompt. In future studies, we intend to explore the application of POLAR within contexts featuring multiple reference trajectories to further improve performance on open-ended tasks. On the other hand, dependence on reference trajectories increases the overall cost associated with reinforcement learning. Although references serve as essential guides, generating and annotating these trajectories incurs additional resource expenditure. Consequently, reducing reliance on references without sacrificing performance constitutes a critical direction for future research. Investigating methods to effectively leverage trajectories generated from alternative prompts, rather than creating separate trajectories for each individual prompt, could significantly decrease annotation costs while potentially enhancing the flexibility and broader applicability of the reward model.

**Integrating POLAR with Test-Time Scaling** While POLAR has demonstrated effectiveness in scaling RMs for policy discrimination, it primarily focuses on pre-training paradigms. Recent advancements in test-time scaling techniques have showcased the potential of dynamically refining reasoning processes during inference via mechanisms like extended computation and self-reflection, such as OpenAI o-series reasoning models and Deepseek R1 [34; 45; 71; 46]. Such methods enable models to adaptively allocate computational resources, significantly improving decision-making quality in challenging scenarios. In future work, we aim to explore the integration of POLAR’s pre-training strategies with test-time scaling techniques, investigating how such a combination can synergistically enhance RM’s performance and generalization in complex tasks.

**Exploring Scaling Potential of POLAR** Given the observed scaling law behavior, we anticipate that the current POLAR series has substantial room for further performance improvements. In future research, we plan to leverage greater computational resources to train larger-scale POLAR RMs. The data-preparing strategy employed by POLAR can be effectively scaled up to extensive pre-training datasets; however, this process inherently demands substantial policy sampling. Compared to traditional LLM data preparation, generating sufficient training data for POLAR is likely to incur considerably higher computational costs. By scaling up model size and computational resources, we aim to thoroughly investigate the limits of POLAR and release stronger, open-source models, thus facilitating continued advancements within the research community.

## 6 Conclusion

We propose a novel perspective on RM by reformulating it as a policy discriminator and introduce a scalable approach named Policy Discriminative Learning (POLAR). Leveraging POLAR, the RM attains robust capabilities in distinguishing among diverse policies and accurately quantifying their differences. These measured differences serve as effective reward signals to guide subsequent policy optimization. Extensive evaluations show that, despite having only small parameters, POLAR RMs significantly surpass larger SOTA baseline RMs, consistently achieving higher preference accuracy across multiple tasks. Additionally, POLAR shows substantial effectiveness in RLHF, providing reliable and generalizable reward signals even in diverse and potentially out-of-domain scenarios. Furthermore, empirical analysis confirms that POLAR exhibits clear power-law scaling behaviors, underscoring its considerable potential for future enhancements.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35: 23716–23736, 2022.
- [2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021. URL <https://arxiv.org/abs/2108.07732>.
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [4] 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.
- [5] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI: Harmlessness from AI feedback, 2022.
- [6] Xiao Bi, Deli Chen, Guanting Chen, Shanhua Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- [7] Gérard Biau. Analysis of a random forests model. *The Journal of Machine Learning Research*, 13:1063–1095, 2012.
- [8] 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.
- [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.
- [10] Ethan Caballero, Kshitij Gupta, Irina Rish, and David Krueger. Broken neural scaling laws. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=sckjveqlCZ>.
- [11] Zheng Cai, Maosong Cao, Haojong 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.
- [12] Lichang Chen, Chen Zhu, Davit Soselia, Juhai Chen, Tianyi Zhou, Tom Goldstein, Heng Huang, Mohammad Shoeybi, and Bryan Catanzaro. Odin: Disentangled reward mitigates hacking in rlhf. *arXiv preprint arXiv:2402.07319*, 2024.
- [13] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex

- Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [14] 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*, 2025.
  - [15] JayMari Chua, Yun Li, Shiyi Yang, Chen Wang, and Lina Yao. Ai safety in generative ai large language models: A survey. *arXiv preprint arXiv:2407.18369*, 2024.
  - [16] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
  - [17] OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models, 2023.
  - [18] 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, 2023.
  - [19] DeepSeek-AI. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024. URL <https://github.com/deepseek-ai/DeepSeek-LLM>.
  - [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
  - [21] Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Luu Anh Tuan, and Shafiq Joty. Data augmentation using llms: Data perspectives, learning paradigms and challenges. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 1679–1705, 2024.
  - [22] Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. *arXiv preprint arXiv:2309.11499*, 2023.
  - [23] Shihan Dou, Yan Liu, Enyu Zhou, Tianlong Li, Haoxiang Jia, Limao Xiong, Xin Zhao, Junjie Ye, Rui Zheng, Tao Gui, et al. Metarm: Shifted distributions alignment via meta-learning. *arXiv preprint arXiv:2405.00438*, 2024.
  - [24] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*, 2022.
  - [25] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246. URL <https://aclanthology.org/N19-1246/>.
  - [26] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- [27] Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El-Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned, 2022. URL <https://arxiv.org/abs/2209.07858>.
- [28] Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.
- [29] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.
- [30] Panagiotis Giadikiaroglou, Maria Lymperaiou, Giorgos Filandrianos, and Giorgos Stamou. Puzzle solving using reasoning of large language models: A survey. *arXiv preprint arXiv:2402.11291*, 2024.
- [31] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024.
- [32] 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.
- [33] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.
- [34] 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 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [35] Jiaxin Guo, Zewen Chi, Li Dong, Qingxiu Dong, Xun Wu, Shaohan Huang, and Furu Wei. Reward reasoning model. *arXiv preprint arXiv:2505.14674*, 2025.
- [36] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- [37] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Hee-woo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- [38] Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. Scaling laws and interpretability of learning from repeated data. *arXiv preprint arXiv:2205.10487*, 2022.
- [39] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostafa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.
- [40] Jacob Hilton, Jie Tang, and John Schulman. Scaling laws for single-agent reinforcement learning. *arXiv preprint arXiv:2301.13442*, 2023.

- [41] Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*, 2024.
- [42] Hui Huang, Yingqi Qu, Jing Liu, Muyun Yang, and Tiejun Zhao. An empirical study of llm-as-a-judge for llm evaluation: Fine-tuned judge models are task-specific classifiers. *arXiv e-prints*, pages arXiv–2403, 2024.
- [43] Saffron Huang, Divya Siddarth, Liane Lovitt, Thomas I Liao, Esin Durmus, Alex Tamkin, and Deep Ganguli. Collective constitutional ai: Aligning a language model with public input. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 1395–1417, 2024.
- [44] Berivan Isik, Natalia Ponomareva, Hussein Hazimeh, Dimitris Paparas, Sergei Vassilvitskii, and Sanmi Koyejo. Scaling laws for downstream task performance of large language models. *arXiv preprint arXiv:2402.04177*, 2024.
- [45] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [46] Qile Jiang, Zhiwei Gao, and George Em Karniadakis. Deepseek vs. chatgpt vs. claude: A comparative study for scientific computing and scientific machine learning tasks. *Theoretical and Applied Mechanics Letters*, 15(3):100583, 2025.
- [47] Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. FollowBench: A multi-level fine-grained constraints following benchmark for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikanth, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4667–4688, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.257>.
- [48] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [49] Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip HS Torr, Fahad Shahbaz Khan, and Salman Khan. Llm post-training: A deep dive into reasoning large language models. *arXiv preprint arXiv:2502.21321*, 2025.
- [50] Sandipan Kundu, Yuntao Bai, Saurav Kadavath, Amanda Askell, Andrew Callahan, Anna Chen, Anna Goldie, Avital Balwit, Azalia Mirhoseini, Brayden McLean, et al. Specific versus general principles for constitutional ai. *arXiv preprint arXiv:2310.13798*, 2023.
- [51] 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.
- [52] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- [53] Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 11260–11285, 2024.
- [54] Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Pre-trained language models for text generation: A survey. *ACM Computing Surveys*, 56(9):1–39, 2024.

- [55] Shiyao Li, Xuefei Ning, Luning Wang, Tengxuan Liu, Xiangsheng Shi, Shengen Yan, Guohao Dai, Huazhong Yang, and Yu Wang. Evaluating quantized large language models. In *Forty-first International Conference on Machine Learning*, 2024.
- [56] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024.
- [57] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval), 5 2023.
- [58] 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=v8L0pN6E0i>.
- [59] Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. Wildbench: Benchmarking llms with challenging tasks from real users in the wild, 2024. URL <https://arxiv.org/abs/2406.04770>.
- [60] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [61] 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*, 2024.
- [62] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- [63] Jiate Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. Rltf: Reinforcement learning from unit test feedback. *arXiv preprint arXiv:2307.04349*, 2023.
- [64] Tianqi Liu, Wei Xiong, Jie Ren, Lichang Chen, Junru Wu, Rishabh Joshi, Yang Gao, Jiaming Shen, Zhen Qin, Tianhe Yu, et al. Rrm: Robust reward model training mitigates reward hacking. *arXiv preprint arXiv:2409.13156*, 2024.
- [65] Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Andrew Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, Xiaohan Zhang, Lichao Sun, Xiaotao Gu, Hongning Wang, Jing Zhang, Minlie Huang, Yuxiao Dong, and Jie Tang. Alignbench: Benchmarking chinese alignment of large language models. In *ACL (1)*, pages 11621–11640, 2024. URL <https://aclanthology.org/2024.acl-long.624>.
- [66] 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.
- [67] Xudong Luo, Zhiqi Deng, Binxia Yang, and Michael Y Luo. Pre-trained language models in medicine: A survey. *Artificial Intelligence in Medicine*, page 102904, 2024.
- [68] Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Reasoning with reinforced fine-tuning. *arXiv preprint arXiv:2401.08967*, 3, 2024.
- [69] Kaijing Ma, Xeron Du, Yunran Wang, Haoran Zhang, ZhoufutuWen, Xingwei Qu, Jian Yang, Jiaheng Liu, minghao liu, Xiang Yue, Wenhao Huang, and Ge Zhang. KOR-bench: Benchmarking language models on knowledge-orthogonal reasoning tasks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=SVRRQ8goQo>.

- [70] 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.
- [71] Sarah Mercer, Samuel Spillard, and Daniel P Martin. Brief analysis of deepseek r1 and it’s implications for generative ai. *arXiv preprint arXiv:2502.02523*, 2025.
- [72] Yuchun Miao, Sen Zhang, Liang Ding, Rong Bao, Lefei Zhang, and Dacheng Tao. Inform: Mitigating reward hacking in rlhf via information-theoretic reward modeling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [73] Oren Neumann and Claudius Gros. Scaling laws for a multi-agent reinforcement learning model. *arXiv preprint arXiv:2210.00849*, 2022.
- [74] OpenAI. Gpt-4 technical report, 2023.
- [75] OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- [76] OpenAI. Openai o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.
- [77] OpenAI. Reinforcement fine-tuning, 2024. URL <https://platform.openai.com/docs/guides/rft-use-cases>.
- [78] 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. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744, 2022.
- [79] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [80] Silviu Pitis, Ziang Xiao, Nicolas Le Roux, and Alessandro Sordoni. Improving context-aware preference modeling for language models. *arXiv preprint arXiv:2407.14916*, 2024.
- [81] Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. Reasoning with large language models, a survey. *arXiv preprint arXiv:2407.11511*, 2024.
- [82] Moschoula Pternea, Prerna Singh, Abir Chakraborty, Yagna Oruganti, Mirco Milletari, Sayli Bapat, and Kebei Jiang. The rl/llm taxonomy tree: Reviewing synergies between reinforcement learning and large language models. *Journal of Artificial Intelligence Research*, 80:1525–1573, 2024.
- [83] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. Is chatgpt a general-purpose natural language processing task solver? In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [84] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained models for natural language processing: A survey. *Science China technological sciences*, 63(10):1872–1897, 2020.
- [85] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. *arXiv preprint arXiv:1611.02683*, 2018.
- [86] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [87] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [88] Prajit Ramachandran, Peter J Liu, and Quoc V Le. Unsupervised pretraining for sequence to sequence learning. *arXiv preprint arXiv:1611.02683*, 2016.

- [89] Alexandre Rame, Nino Vieillard, Leonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier Bachem, and Johan Ferret. WARM: On the benefits of weight averaged reward models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=s7RDnNUJy6>.
- [90] Noam Razin, Zixuan Wang, Hubert Strauss, Stanley Wei, Jason D Lee, and Sanjeev Arora. What makes a reward model a good teacher? an optimization perspective. *arXiv preprint arXiv:2503.15477*, 2025.
- [91] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- [92] Yangjun Ruan, Chris J Maddison, and Tatsunori B Hashimoto. Observational scaling laws and the predictability of langauge model performance. *Advances in Neural Information Processing Systems*, 37:15841–15892, 2025.
- [93] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [94] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. How much can clip benefit vision-and-language tasks? *arXiv preprint arXiv:2107.06383*, 2021.
- [95] Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Wei-Yin Ko, Madeline Smith, Antoine Bosselut, Alice Oh, Andre F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker. Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation, 2024. URL <https://arxiv.org/abs/2412.03304>.
- [96] Zayne Rea Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. MuSR: Testing the limits of chain-of-thought with multistep soft reasoning. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=jenyYQzue1>.
- [97] Hao Sun, Yunyi Shen, and Jean-Francois Ton. Rethinking bradley-terry models in preference-based reward modeling: Foundations, theory, and alternatives. *arXiv preprint arXiv:2411.04991*, 2024.
- [98] Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Alpha-clip: A clip model focusing on wherever you want. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13019–13029, 2024.
- [99] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In Anna Rogers, Jordan Boyd-Graber, and Naoki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.824. URL <https://aclanthology.org/2023.findings-acl.824/>.
- [100] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [101] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. <https://github.com/InternLM/InternLM-techreport>, 2023.

- [102] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [103] Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- [104] Changyao Tian, Xizhou Zhu, Yuwen Xiong, Weiyun Wang, Zhe Chen, Wenhui Wang, Yuntao Chen, Lewei Lu, Tong Lu, Jie Zhou, et al. Mm-interleaved: Interleaved image-text generative modeling via multi-modal feature synchronizer. *arXiv preprint arXiv:2401.10208*, 2024.
- [105] Guiyao Tie, Zeli Zhao, Dingjie Song, Fuyang Wei, Rong Zhou, Yurou Dai, Wen Yin, Zhejian Yang, Jiangyue Yan, Yao Su, et al. A survey on post-training of large language models. *arXiv preprint arXiv:2503.06072*, 2025.
- [106] Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, et al. Secrets of rlhf in large language models part ii: Reward modeling. *arXiv preprint arXiv:2401.06080*, 2024.
- [107] Binghai Wang, Rui Zheng, Lu Chen, Zhiheng Xi, Wei Shen, Yuhao Zhou, Dong Yan, Tao Gui, Qi Zhang, and Xuan-Jing Huang. Reward modeling requires automatic adjustment based on data quality. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4041–4064, 2024.
- [108] Binghai Wang, Runji Lin, Keming Lu, Le Yu, Zhenru Zhang, Fei Huang, Chujie Zheng, Kai Dang, Yang Fan, Xingzhang Ren, et al. Worldpm: Scaling human preference modeling. *arXiv preprint arXiv:2505.10527*, 2025.
- [109] Haifeng Wang, Jiwei Li, Hua Wu, Eduard Hovy, and Yu Sun. Pre-trained language models and their applications. *Engineering*, 25:51–65, 2023.
- [110] Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024. URL <https://arxiv.org/abs/2312.08935>.
- [111] Xiao Wang, Guangyao Chen, Guangwu Qian, Pengcheng Gao, Xiao-Yong Wei, Yaowei Wang, Yonghong Tian, and Wen Gao. Large-scale multi-modal pre-trained models: A comprehensive survey. *Machine Intelligence Research*, 20(4):447–482, 2023.
- [112] Xin Wang, Yasheng Wang, Yao Wan, Fei Mi, Yitong Li, Pingyi Zhou, Jin Liu, Hao Wu, Xin Jiang, and Qun Liu. Compilable neural code generation with compiler feedback. *arXiv preprint arXiv:2203.05132*, 2022.
- [113] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. MMLU-pro: A more robust and challenging multi-task language understanding benchmark. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=y10DM6R2r3>.
- [114] Larry Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.
- [115] Congying Xia, Chen Xing, Jiangshu Du, Xinyi Yang, Yihao Feng, Ran Xu, Wenpeng Yin, and Caiming Xiong. FOFO: A benchmark to evaluate LLMs’ format-following capability. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 680–699, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.40. URL <https://aclanthology.org/2024.acl-long.40/>.
- [116] Adam X Yang, Maxime Robeyns, Thomas Coste, Jun Wang, Haitham Bou Ammar, and Laurence Aitchison. Bayesian reward models for llm alignment. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*, 2024.

- [117] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [118] Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, Jiawei Hong, Kuikun Liu, Ziyi Wang, et al. Internlm-math: Open math large language models toward verifiable reasoning. *arXiv preprint arXiv:2402.06332*, 2024.
- [119] Alex Young, Bei Chen, Chao Li, Chengan Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.
- [120] Yinlong Yuan, Zhu Liang Yu, Zhenghui Gu, Xiaoyan Deng, and Yuanqing Li. A novel multi-step reinforcement learning method for solving reward hacking. *Applied Intelligence*, 49:2874–2888, 2019.
- [121] Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, Xiangpeng Wei, et al. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025.
- [122] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, 2019.
- [123] 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.
- [124] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
- [125] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=uccHPGDlao>.
- [126] Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Yuhao Zhou, Limao Xiong, et al. Delve into ppo: Implementation matters for stable rlhf. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023.
- [127] Enyu Zhou, Guodong Zheng, Binghai Wang, Zhiheng Xi, Shihan Dou, Rong Bao, Wei Shen, Limao Xiong, Jessica Fan, Yurong Mou, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. RMB: Comprehensively benchmarking reward models in LLM alignment. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=kmgr1G9TR0>.

## Appendix

### A Training Details

#### A.1 Pre-training

##### A.1.1 Pre-training Data

We construct a synthetic pre-training dataset for POLAR by generating prompt-trajectory pairs using a diverse collection of LLMs. Specifically, we first sample prefixes from the pre-training corpus typically used for LLM training. For each piece of text from the corpus, we randomly select an initial segment as the prefix, with lengths varying randomly between 1 and 1024 tokens. Subsequently, based on each prefix, various policies autoregressively generate trajectories constrained to a maximum length of 4096 tokens. For each prefix, we randomly select two different policies from a policy pool, which contains LLMs with varying architectures and parameter scales. The first policy generates two trajectories: one serves as the reference, and the other serves as the positive trajectory. Conversely, the second policy generates a single trajectory from the same prefix, acting as a negative sample. Additionally, to enrich dataset diversity, we include a small proportion of prompts that represent human instructions, which are derived from instruction-tuning datasets. Instruction-tuned LLMs are then used to generate responses as trajectories for these instruction-based prompts.

Table 4 summarizes the composition of trajectories used for RM pre-training, totaling 3.60T tokens. Table 5 enumerates the 53 open-source pre-trained (base) LLMs and 53 open-source instruction-tuned (chat) LLMs in our policy pool for generating these trajectories. To further enhance policy diversity and ensure broader distributional coverage, we additionally incorporate 78 intermediate training checkpoints from a single pre-trained LLM—InternLM3-8B-base.

All policies span the period from December 2023 to the present and were specifically chosen based on variations in pre-training data versions and the diversity exhibited in their generated responses.

During data construction, we encountered several issues that affected data quality. Firstly, certain policies tended to fall into repetitive loops, continuously generating identical or very similar content without meaningful progression. Instead of completely filtering out these outputs, we truncated repeated sections while preserving their loop-prone characteristics. Secondly, some generated trajectories were excessively long and lacked a clearly defined endpoint. To resolve this, we imposed a maximum output length of 4096 tokens and truncated any incomplete endings, ensuring each generated sequence was self-contained. Lastly, we observed that a low sampling temperature resulted in insufficient contrast for effective contrastive learning, while a high sampling temperature introduced biases, adversely impacting policy characterization. Based on these insights, we set the sampling temperature to 1.0, top-p to 0.9, and top-k to 50 in our trajectory sampling process.

##### A.1.2 Hyperparameter and Implementations

Our primary experiments are conducted using RMs containing 1.8 billion parameters and 7 billion parameters, denoted as POLAR-1.8B and POLAR-7B, respectively. The model architecture is detailed in Section 4.1, and we adopt the XTuner<sup>2</sup> framework for pre-training and fine-tuning. Rather than training from scratch, we initialize the RM from a pre-trained InternLM2.5-series model and perform one additional epoch of POLAR pre-training.

To identify optimal hyperparameters for pre-training POLAR RMs, we carried out scaling experiments designed to establish data-driven scaling laws. These scaling laws relate the optimal hyperparameters to the model size ( $N$ ), base model pre-training data size ( $D_p$ ), and reward model pre-training data size ( $D_{rm}$ ). The results of these scaling experiments are illustrated in Figures 5 and 6.

In the first stage, we select base models pre-trained with various combinations of  $N$  and  $D_p$  and train them on reward model datasets of varying sizes ( $D_{rm}$ ). We fix the batch size at 256 and

<sup>2</sup><https://github.com/InternLM/xtuner>

Table 5: List of Base and Chat LLMs used for generating trajectories for gathering pre-training data. We mainly utilized the Llama [32] series, Qwen [3; 103; 117] series, Yi [119] series, ChatGLM [31] series, DeepSeek[6] series, Gemma[100] series, and InternLM[101; 11] series models.

Base LLMs			
Llama-3.2-1B	Llama-3.2-3B	Meta-Llama-3-70B	Meta-Llama-3-8B
Meta-Llama-3.1-70B	Meta-Llama-3.1-8B	Qwen-14B	Qwen-1_8B
Qwen-72B	Qwen-7B	Qwen1.5-0.5B	Qwen1.5-1.8B
Qwen1.5-14B	Qwen1.5-32B	Qwen1.5-4B	Qwen1.5-72B
Qwen1.5-7B	Qwen2-0.5B	Qwen2-1.5B	Qwen2-72B
Qwen2-7B	Qwen2.5-0.5B	Qwen2.5-1.5B	Qwen2.5-14B
Qwen2.5-32B	Qwen2.5-3B	Qwen2.5-72B	Qwen2.5-7B
Yi-1.5-34B	Yi-1.5-6B	Yi-1.5-9B	Yi-34B
Yi-34B-200K	Yi-6B	Yi-6B-200K	Yi-9B
Yi-9B-200K	chatglm3-6b-base	deepseek-llm-67b-base	deepseek-llm-7b-base
gemma-2b	gemma-7b	glm-4-9b	internlm-20b
internlm-7b	internlm2-1_8b	internlm2-20b	internlm2-7b
internlm2-base-20b	internlm2-base-7b	internlm2_5-1_8b	internlm2_5-20b
internlm2_5-7b			
Chat LLMs			
Llama-3.2-1B-Instruct	Llama-3.2-3B-Instruct	Meta-Llama-3-70B-Instruct	Meta-Llama-3.1-70B-Instruct
Meta-Llama-3.1-8B-Instruct	Qwen-14B-Chat	Qwen-1_8B-Chat	Qwen-72B-Chat
Qwen-7B-Chat	Qwen1.5-0.5B-Chat	Qwen1.5-1.8B-Chat	Qwen1.5-14B-Chat
Qwen1.5-32B-Chat	Qwen1.5-4B-Chat	Qwen1.5-72B-Chat	Qwen1.5-7B-Chat
Qwen2-0.5B-Instruct	Qwen2-1.5B-Instruct	Qwen2-72B-Instruct	Qwen2-7B-Instruct
Qwen2.5-0.5B-Instruct	Qwen2.5-1.5B-Instruct	Qwen2.5-14B-Instruct	Qwen2.5-32B-Instruct
Qwen2.5-3B-Instruct	Qwen2.5-72B-Instruct	Qwen2.5-7B-Instruct	Yi-1.5-34B-Chat
Yi-1.5-34B-Chat-16K	Yi-1.5-6B-Chat	Yi-1.5-9B-Chat	Yi-34B-Chat
Yi-6B-Chat	chatglm3-6b	chatglm3-6b-32k	deepseek-llm-67b-chat
deepseek-llm-7b-chat	gemma-2b-it	gemma-7b-it	glm-4-9b-chat
glm-4-9b-chat-1m	internlm-chat-20b	internlm-chat-7b	internlm2-chat-1_8b
internlm2-chat-1_8b-sft	internlm2-chat-20b	internlm2-chat-20b-sft	internlm2-chat-7b
internlm2-chat-7b-sft	internlm2_5-1_8b-chat	internlm2_5-20b-chat	internlm2_5-7b-chat
internlm2_5-7b-chat-1m			

explore learning rates (LR) ranging from 5e-6 to 1e-4 for models with parameters between 100M and 924M. For each configuration, we calculate the validation loss and fit a quartic polynomial to the loss – log(LR) relationship. We identify the optimal learning rate corresponding to the minimum loss for each combination of  $N$ ,  $D_p$ , and  $D_{rm}$ , resulting in the empirical formula:

$$LR = 0.0002306 \cdot N^{0.01125} \cdot D_p^{-0.66587} \cdot D_{rm}^{0.33916}$$

where  $N$ ,  $D_p$ , and  $D_{rm}$  are all expressed in millions. In the second stage, we search for the optimal batch size within the range of 128 to 1024. Learning rates are computed using the above formula and proportionally scaled based on batch size adjustments relative to the baseline of 256. Employing a similar method, we fit a quadratic polynomial to the loss – log(batch size) curve and obtain:

$$\text{Batch Size} = 31.9032 \cdot N^{0.06944} \cdot D_{rm}^{0.52997}$$

Finally, for POLAR-1.8B, by substituting the values  $N = 1.8B$ ,  $D_p = 2.5T$ , and  $D_{rm} = 0.94T$ , we determine the final learning rate to be 1.4e-5 and the batch size as 1940. The pre-training process is conducted on 320 NVIDIA H800 GPUs for a total duration of 57 hours.

For pre-training POLAR-7B, we set  $N = 7B$ ,  $D_p = 4.0T$ , and  $D_{rm} = 3.6T$ . Then we get the learning rate 1.67e-5 and the batch size 4343. The pre-training process is conducted on 912 NVIDIA H800 GPUs for a total duration of 175 hours.

## A.2 Supervised Fine-tuning

### A.2.1 Supervised Fine-tuning Data

We construct a dataset comprising 150K manually labeled examples to fine-tune the pre-trained reward model. Each example includes a single prompt associated with three candidate outputs. The prompts are primarily sourced from widely used open-source preference-pair datasets such

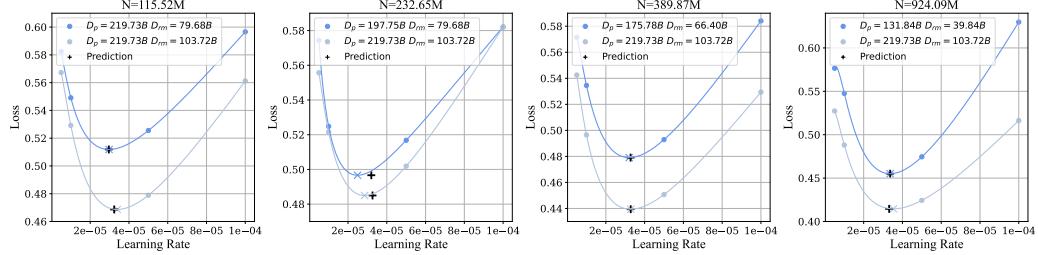


Figure 5: Scaling law of the learning rate with respect to model size and data scale in pre-training.

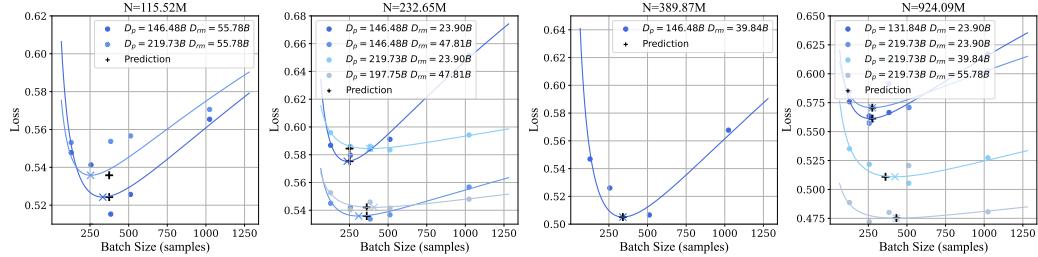


Figure 6: Scaling law of the batch size with respect to model size and data scale in pre-training.

as UltraFeedback [18] and HH-RLHF [4; 27], with a small subset derived from real user queries submitted to online chat platforms.

For prompts obtained from open-source datasets that originally contain two candidate outputs, we generate a third candidate using a state-of-the-art LLM randomly selected from GPT-4o [75], OpenAI o1 [76], Deepseek-R1 [34], and Deepseek-V3 [60], to mitigate potential distribution biases. Human annotators then reorder these three outputs based on preference and adherence to instructions. For prompts sourced from actual user queries, we generate three candidate outputs utilizing LLMs selected either from the models listed in Table 5 or from the aforementioned state-of-the-art models. Human annotators subsequently evaluate and annotate relative preferences among these outputs. Consequently, each prompt is associated with three outputs ranked from best to worst. The top two outputs constitute a positive pair, while the second and third-ranked outputs form a negative pair.

All annotations were performed by company staff possessing relevant professional expertise and compensated at standard salary rates. To safeguard user privacy, we filtered all prompts in the training set to remove personally identifiable information.

### A.2.2 Hyperparameter and Implementations

For supervised fine-tuning of POLAR models, we set the learning rate to 1e-5 for the 1.8B model and 2e-5 for the 7B model, use a batch size of 320, and train for one epoch. We experiment with multiple hyperparameter configurations and select the optimal one that avoids overoptimization. Each round of supervised fine-tuning runs on 16 NVIDIA H800 GPUs for approximately 0.5 hours.

For supervised fine-tuning of models in the ablation study, we adopt the same data and settings as POLAR, with two exceptions: The fine-tuning-only reward model (w/o PT) directly uses the InternLM2.5-series as the backbone, instead of POLAR’s pre-trained models. The traditional reward model (w/o PT & Ref) is also fine-tuned on InternLM2.5-series and trained on the same preference pairs as POLAR, but without incorporating reference trajectories, following the mainstream reward model training paradigm.

## B Details of Evaluation Setup

### B.1 Baseline Reward Models

**InternLM2-7B-Reward** and **InternLM2-20B-Reward** [11] are reward models built on InternLM2-Chat-7B-SFT and InternLM2-Chat-20B-SFT, respectively. They are trained on over 2.4 million preference samples, comprising both human-annotated and AI-generated data.

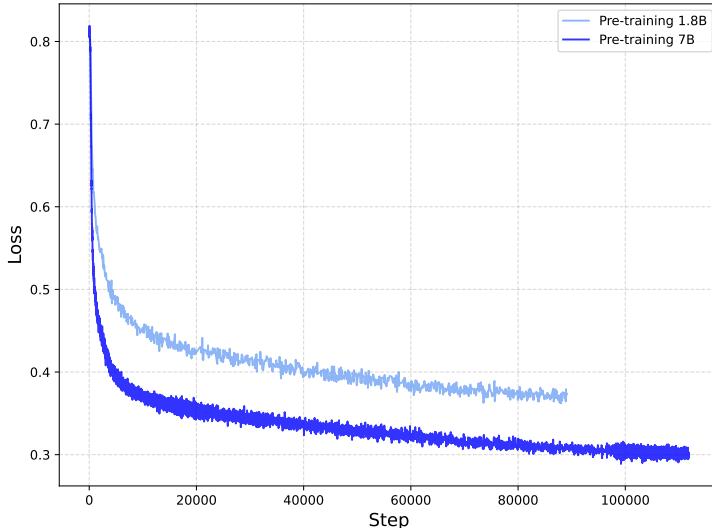


Figure 7: Pre-training loss curves for POLAR-1.8B and POLAR-7B.

**Skywork-Reward-Gemma-2-27B** and **Skywork-Reward-Llama-3.1-8B** [61] are reward models based on the gemma-2-27b-it and Meta-Llama-3.1-8B-Instruct architectures, respectively. Both models are trained using the Skywork Reward Data Collection, which consists of 80K high-quality preference pairs curated from publicly available sources. Notably, Skywork-Reward-Gemma-2-27B achieves state-of-the-art performance in several prior studies.

**WorldPM-72B-UltraFeedback** [108] is a SOTA reward model recently released by the Qwen Group. It is initially trained on 15M preference pairs collected from public forums spanning diverse user communities, capturing a unified representation of human preferences. The model is then fine-tuned on 100K preference pairs from UltraFeedback [18], a fine-grained and diverse preference dataset.

## B.2 Details of Preference Evaluation

For the preference evaluation set constructed from real user queries submitted to online platforms (as mentioned in Section 4.2), we generate three candidate trajectories per query using either the LLMs listed in Table 5 or the state-of-the-art models mentioned earlier. Human annotators are then instructed to rank these candidate outputs based on quality. To ensure a fair evaluation, we carefully exclude any examples overlapping with the training dataset. Following the same evaluation protocol as in the RMB set, we assess the reward model’s capability to correctly identify the better of the two remaining trajectories, given the highest-ranked output as a reference.

## B.3 Details of RLHF Training and Evaluation

**Training Data** We construct a dataset containing 1.25 million prompts paired with reference trajectories to support RLHF training for policy models. The prompts are primarily drawn from widely adopted open-source instruction datasets such as UltraFeedback [18] and HH-RLHF [4; 27], supplemented by a smaller subset collected from real user queries on online chat platforms. To promote better generalization, we ensure no overlap exists between prompts used for RLHF and those utilized during SFT. For each prompt, a reference trajectory is generated by randomly selecting a state-of-the-art LLM from GPT-4o [75], OpenAI o1 [76], Deepseek-R1 [34], and Deepseek-V3 [60]. Notably, our RLHF dataset is constructed entirely without the involvement of human annotators.

**Hyperparameter and Implementations** For RLHF experiments, we train policy models using the PPO algorithm implemented in OpenRLHF [41], guided by both our proposed reward model and several baseline reward models. To ensure robustness, experiments are conducted across multiple random seeds, and we report the results averaged over these runs. To maintain fairness in evaluation, each baseline reward model is tested under two scoring conditions (similar to preference evaluation)

when assigning rewards during PPO training: (1) standard scoring without reference trajectories, and (2) scoring with reference trajectories explicitly included in the prompt (Figure 10). For each baseline, we report the best performance observed across these two settings. Notably, we consistently find that standard scoring without references produces superior results, likely due to its closer alignment with the original training objectives of the baseline reward models.

During PPO training, we set the actor learning rate to  $1e-6$ , the critic learning rate to  $1e-5$ , the training batch size to 1024, the rollout batch size to 1024, and the number of epochs to 1. For all other hyperparameters, we generally follow OpenRLHF’s recommended configurations. PPO experiments for all policy models, except Qwen2.5-32B-Instruct, are conducted using 32 NVIDIA H800 GPUs per run, each taking approximately 48 hours. For Qwen2.5-32B-Instruct, we utilize 64 NVIDIA H800 GPUs per run, with each run lasting roughly 72 hours.

## C Additional Results

### C.1 Pre-Training Loss Curve

Figure 7 presents the relationship between training steps and loss for POLAR-1.8B and POLAR-7B, demonstrating a steady decline in loss and indicating stable convergence during pre-training.

### C.2 Case Study

To showcase the effectiveness of our method, Figure 8 and Figure 9 illustrate two representative examples of how POLAR-7B assigns rewards to trajectories. Detailed explanations are provided in the respective comments and captions.

Reward Scoring Example 1 for POLAR-7B

**User Prompt:**  
How many 'r's are there in the word 'strawberry'?

**Reference Response:**  
There are 3 'r's in the word 'strawberry'.

**Candidate Response 1:**  
There are 3 'r's in the word 'strawberry'.  
# ✓ Same as the reference response.  
# Reward: 0.690

**Candidate Response 2:**  
There are 2 'r's in the word 'strawberry'.  
# Wrong answer.  
# Reward: -9.523

**Candidate Response 3:**  
Let's count the 'r's in 'strawberry': 's', 't', 'r', 'a', 'w', 'b', 'e', 'r', 'r', 'y'. There are 3 'r's, so the answer is 3.  
# ✓ Correct answer with correct thoughts.  
# Reward: -6.476

**Candidate Response 4:**  
Let's count the 'r's in 'strawberry': 's', 't', 'r', 'a', 'w', 'b', 'e', 'r', 'r', 'y'. There are 2 'r's, so the answer is 3.  
# Correct answer with wrong thoughts.  
# Reward: -9.289

**Candidate Response 5:**  
Let's count the 'r's in 'strawberry': 's', 't', 'r', 'a', 'w', 'b', 'e', 'r', 'r', 'y'. There are 3 'r's, so the answer is 2.

# Wrong answer with correct thoughts.

# Reward: -9.562

**Candidate Response 6:**

Let's count the 'r's in 'strawberry': 's', 't', 'r', 'a', 'w', 'b', 'e', 'r', 'r', 'y'. There are 2 'r's, so the answer is 2.

# Wrong answer with wrong thoughts.

# Reward: -10.937

**Candidate Response 7:**

The word 'strawberry' contains 3 'r's.

# ✓ Correct answer.

# Reward: -1.341

**Candidate Response 8:**

The word 'strawberry' contains 2 'r's.

# Wrong answer.

# Reward: -9.929

Figure 8: An example of reward model scoring applied to the popular “strawberry” test. Despite variations in response phrasing, our reward model accurately identifies the correctness of the final answer and distinguishes the validity of the reasoning process.

Reward Scoring Example 2 for POLAR-7B

**User Prompt:**

Summarize the first book of Frank Herbert’s Dune in one witty short sentence.

**Reference Response:**

Royal teen discovers that life’s a beach—minus the ocean, plus spice, giant sandworms and deadly politics.

**Candidate Response 1:**

Royal teen discovers that life’s a beach—minus the ocean, plus spice, giant sandworms and deadly politics.

# ✓ Same as the reference response.

# Reward: 0.467

**Candidate Response 2:**

Royal teen discovers that life’s a beach—minus the ocean, plus magic, dark wizards and deadly politics.

# Closely resembles the reference response but includes factual errors.

# Reward: -7.711

**Candidate Response 3:**

Young noble’s move to desert planet turns into galactic Game of Thrones with fewer dragons, more worms.

# ✓ A distinct yet concise and witty summary that draws analogies from other dramas—markedly different from the reference response.

# Reward: -6.918

**Candidate Response 4:**

A noble family’s fall sparks a young heir’s rise as a leader on a harsh desert planet governed by prophecy and survival.

# A concise summary, but lacking wit—fails to meet the requirement.

# Reward: -8.648

**Candidate Response 5:**

Paul Atreides loses his father, gains prophetic powers, learns to ride a sandworm, leads a holy war, and discovers that being the chosen one comes with a lot of blood, sand, and questionable decisions.

# A witty summary, but overly long—fails to meet the requirement.

# Reward: -8.429

**Candidate Response 6:**

Boy gets planet, becomes god, loses soul — family drama ensues across galaxies.

# A concise and witty summary that draws from multiple Dune books rather than just the first—fails to follow the instruction.

# Reward: -10.359

Figure 9: An example of reward model scoring applied to a summarization task. For open-ended questions, our model accurately evaluates the alignment of the core content of responses with the prompt and reference answer, considering multiple evaluation criteria rather than merely relying on semantic similarity to the reference response.

### C.3 Additional Results of POLAR Evaluation

Figure 10 illustrates the prompt format used by baseline reward models to score trajectories when a reference trajectory is provided, as mentioned in Section 4.2.

Table 6 shows the details of our employed benchmarks for policy evaluation. Table 7, 8, 9, and 10 show the whole results of four policies in RLHF. Table 11, 12 and 13 present additional ablation results mentioned in Section 4.5.

Prompt including reference for reward model baselines

**System:**

You are a helpful AI assistant, designed to provide useful answers to customers. The user will provide you with a question and a reference answer. Please refer to this desired output and provide your own unique response.

**User:**

You will be provided with the following information:

[The Start of conversation history] {dialog history} [The end of conversation History]

[The start of question] {question} [The end of question]

[The start of desired output] {reference} [The end of desired output]

Please think carefully about this question and the desired output, and then provide your response to the question.

**Assistant:**

{response}

Figure 10: The prompt format used by baseline reward models to score trajectories with a provided reference, as mentioned in Section 4.2. Specifically, we construct a dialogue scenario in which the assistant generates a response to the prompt with access to the reference, and the baseline reward models then assign scores to the response. To ensure a fair comparison with our POLAR RMs, we report the best performance of each baseline across both the reference-included and traditional (no-reference) settings.

Table 6: Benchmarks we used in policy evaluation.

Benchmark	Category	Description
<b>AlignBench</b> [65]	General Task	A comprehensive Chinese benchmark for LLM alignment evaluation.
<b>AlpacaEval</b> [57]	General Task	An automated benchmark for evaluating instruction-following ability.
<b>Arena-Hard</b> [56]	General Task	A dynamic benchmark on challenging, real-world tasks.
<b>WildBench</b> [59]	General Task	An automated benchmark on real-world user tasks.
<b>MT-Bench</b> [125]	General Task	A benchmark on multi-turn conversational tasks.
<b>FollowBench</b> [47]	Instruction Following	A benchmark for evaluating LLMs' ability to follow instructions with multi-level, fine-grained constraints.
<b>FoFo</b> [115]	Instruction Following	A benchmark for evaluating LLMs' ability to follow complex, domain-specific formats and formatting rules.
<b>HumanEval</b> [13]	Coding	A benchmark for code generation tasks measuring functional correctness for synthesizing programs from docstrings.
<b>MBPP</b> [2]	Coding	A code generation benchmark including 1,000 entry-level Python programming problems.
<b>DROP</b> [25]	General Reasoning	A reading comprehension benchmark requiring discrete reasoning over paragraphs.
<b>BBH</b> [99]	General Reasoning	A diverse benchmark comprising 23 challenging BIG-Bench tasks.
<b>GPQA</b> [91]	General Reasoning	A graduate-level Google-Proof Q&A benchmark covering the fields of biology, physics, and chemistry.
<b>HellaSwag</b> [122]	General Reasoning	A challenging dataset for evaluating commonsense NLI.
<b>MuSR</b> [96]	General Reasoning	A dataset on multistep soft reasoning tasks specified in a natural language narrative.
<b>KOR-Bench</b> [69]	General Reasoning	A challenging benchmark comprising five tasks that introduce new rules independent of prior knowledge.
<b>GSM8K</b> [16]	Math	A dataset of 8.5K high quality linguistically diverse grade school math word problems .
<b>MATH-500</b> [58]	Math	A challenging dataset of 500 high school math competition problems.
<b>CMMLU</b> [53]	Knowledge	A comprehensive Chinese benchmark for assessing knowledge and reasoning abilities in Chinese contexts.
<b>MMLU-Pro</b> [113]	Knowledge	An enhanced version benchmark containing over 12,000 academic exam and textbook questions across 14 domains.
<b>MMMLU-Lite</b> [95]	Knowledge	A massive multilingual multitask language understanding dataset spanning 15 languages.

Table 7: Reward model comparison in RLHF training. The policy model is initialized from InternLM3-8B-Instruct [11] (i.e., baseline). Sky8B and Sky27B denote Skywork-Reward-8B and Skywork-Reward-8B [61], respectively. Int7B and Int20B denote InternLM2-Reward-7B and InternLM2-Reward-20B [11], respectively. WPM72B denotes WorldPM-72B-UltraFeedback [108].

Type	Dataset	Policy model: InternLM3-8B-Instruct							
		Baseline	Int7B	Sky8B	Int20B	Sky27B	WPM72B	POLAR-1.8B	POLAR-7B
General Task	AlignBench	6.13	6.24	6.09	6.37	6.27	6.43	6.54	6.53
	AlpacaEval	51.80	57.02	61.37	58.39	63.48	70.93	68.57	68.70
	Arena-Hard	40.79	45.96	45.01	47.33	45.56	53.69	63.47	63.34
	WildBench	13.67	22.64	25.30	23.42	27.40	35.00	40.34	39.75
	MT-Bench	7.97	8.24	8.27	8.29	8.27	8.40	8.57	8.43
Instruct Following	FollowBench	82.20	83.20	83.60	85.30	84.80	86.20	91.10	90.40
	FoFo	43.10	45.70	43.90	48.20	45.10	49.60	54.30	56.10
Coding	HumanEval	81.10	82.93	78.88	78.66	78.66	81.10	82.93	84.15
	MBPP	67.70	74.32	70.43	69.65	70.04	73.15	73.54	75.10
General Reasoning	DROP	82.82	82.92	82.70	83.13	82.88	83.34	86.64	85.63
	BBH	77.56	75.12	76.99	76.43	77.27	77.97	77.06	77.70
	GPQA	32.83	36.36	35.35	36.87	37.37	37.88	37.37	41.41
	HellaSwag	88.88	87.15	87.65	87.50	87.60	87.06	88.78	88.87
	MuSR	52.28	51.48	51.59	49.64	51.23	50.41	57.07	58.10
	KOR-Bench	51.84	56.00	54.64	56.24	54.72	56.72	53.84	55.60
Math	GSM8K	90.22	89.31	89.92	89.99	90.45	90.98	91.05	91.96
	MATH-500	76.00	70.60	76.80	74.40	76.00	77.60	77.60	78.40
Knowledge	CMLLU	81.70	82.26	81.90	82.42	82.18	82.75	84.79	84.39
	MMLU-Pro	57.49	58.05	57.94	57.99	57.96	58.20	60.00	59.98
	MMMLU-Lite	43.64	40.97	40.02	41.54	39.59	42.30	48.40	49.00
Average		56.49	57.82	57.92	58.09	58.34	60.49	62.60	<b>63.18</b>

Table 8: Reward model comparison in RLHF training. The policy model is initialized from Llama-3.1-8B-Instruct [26] (i.e., baseline). Sky8B and Sky27B denote Skywork-Reward-8B and Skywork-Reward-8B [61], respectively. Int7B and Int20B denote InternLM2-Reward-7B and InternLM2-Reward-20B [11], respectively. WPM72B denotes WorldPM-72B-UltraFeedback [108].

Type	Dataset	Policy model: Llama-3.1-8B-Instruct							
		Baseline	Int7B	Sky8B	Int20B	Sky27B	WPM72B	POLAR-1.8B	POLAR-7B
General Task	AlignBench	4.66	4.89	5.07	4.90	5.18	4.94	4.89	5.25
	AlpacaEval	23.73	48.57	50.01	51.68	50.56	41.16	50.19	72.30
	Arena-Hard	42.31	45.26	39.84	37.27	39.22	38.78	52.60	61.95
	WildBench	-0.89	19.88	20.81	30.70	19.54	13.66	23.69	37.12
	MT-Bench	8.15	8.24	8.27	8.03	8.33	8.28	8.42	8.46
Instruct Following	FollowBench	89.70	91.00	91.20	90.90	91.20	90.30	92.10	92.90
	FoFo	37.00	30.60	32.40	34.80	32.20	37.40	38.30	45.70
Coding	HumanEval	71.34	54.27	64.63	59.15	63.41	70.12	70.72	70.73
	MBPP	70.04	64.20	70.43	57.98	69.26	71.60	71.98	73.54
General Reasoning	DROP	79.46	79.33	81.42	74.43	82.53	81.60	84.48	84.14
	BBH	54.26	49.37	55.57	58.25	61.79	60.50	64.34	70.02
	GPQA	29.80	32.32	31.88	31.31	27.78	28.28	32.83	39.90
	HellaSwag	46.90	63.56	53.06	52.60	51.03	53.99	60.92	67.84
	MuSR	60.73	54.22	51.15	53.76	56.18	56.86	55.70	50.73
	KOR-Bench	46.56	46.08	48.16	44.08	48.16	47.20	46.86	46.48
Math	GSM8K	83.40	83.02	85.06	83.09	82.49	85.52	85.82	88.40
	MATH-500	51.80	47.40	47.40	45.80	50.00	53.60	56.40	56.00
Knowledge	CMLLU	54.78	52.81	54.94	52.86	55.63	54.96	55.29	53.10
	MMLU-Pro	49.57	48.11	50.14	47.53	49.84	50.68	51.57	52.66
	MMMLU-Lite	43.83	38.12	43.01	34.87	44.43	43.45	47.04	49.32
Average		47.36	48.06	49.22	47.70	49.44	49.64	52.71	<b>56.33</b>

Table 9: Reward model comparison in RLHF training. The policy model is initialized from Qwen2.5-7B-Instruct [117] (i.e., baseline). Sky8B and Sky27B denote Skywork-Reward-8B and Skywork-Reward-8B [61], respectively. Int7B and Int20B denote InternLM2-Reward-7B and InternLM2-Reward-20B [11], respectively. WPM72B denotes WorldPM-72B-UltraFeedback [108].

Type	Dataset	Policy model: Qwen2.5-7B-Instruct							
		Baseline	Int7B	Sky8B	Int20B	Sky27B	WPM72B	POLAR-1.8B	POLAR-7B
General Task	AlignBench	6.18	6.22	6.27	6.25	6.30	6.48	6.42	6.49
	AlpacaEval	37.02	58.39	58.63	62.81	61.61	62.18	60.13	61.86
	Arena-Hard	56.73	51.88	54.10	53.81	59.64	61.52	67.11	71.38
	WildBench	24.30	35.10	34.78	33.99	35.40	39.99	36.63	40.30
	MT-Bench	8.38	8.37	8.40	8.40	8.43	8.43	8.51	8.49
Instruct Following	FollowBench	88.00	89.00	89.20	89.40	89.90	90.70	90.50	90.30
	FoFo	44.10	39.10	46.80	47.40	49.00	50.40	52.20	52.00
Coding	HumanEval	84.15	78.66	79.88	76.83	81.71	81.41	84.15	87.20
	MBPP	74.32	66.93	73.54	71.28	74.71	73.54	76.65	75.10
General Reasoning	DROP	78.35	78.09	81.10	78.44	80.25	80.79	82.67	83.22
	BBH	56.92	62.73	66.83	65.15	67.50	66.50	68.67	67.14
	GPQA	32.32	33.33	33.84	33.36	34.34	38.64	38.38	37.37
	HellaSwag	67.47	74.63	74.73	65.48	71.27	76.24	84.91	81.25
	MuSR	46.53	41.79	44.11	40.84	45.31	46.22	51.04	50.64
	KOR-Bench	41.36	48.32	47.92	49.20	46.08	48.48	49.44	48.16
Math	GSM8K	91.13	88.70	90.45	89.84	91.05	91.40	91.58	92.19
	MATH-500	75.80	72.00	75.80	75.40	76.10	75.30	76.80	77.20
Knowledge	CMLLU	77.56	73.80	75.53	74.97	76.94	76.76	78.07	77.92
	MMLU-Pro	55.33	52.30	55.54	54.89	55.34	55.48	55.97	56.65
	MMMLU-Lite	53.04	39.61	43.28	45.22	46.01	46.19	47.13	53.15
Average		54.95	54.95	57.04	56.15	57.84	58.83	60.35	<b>60.90</b>

Table 10: Reward model comparison in RLHF training. The policy model is initialized from Qwen2.5-32B-Instruct [117] (i.e., baseline). Sky8B and Sky27B denote Skywork-Reward-8B and Skywork-Reward-8B [61], respectively. Int7B and Int20B denote InternLM2-Reward-7B and InternLM2-Reward-20B [11], respectively. WPM72B denotes WorldPM-72B-UltraFeedback [108].

Type	Dataset	Policy model: Qwen2.5-32B-Instruct							
		Baseline	Int7B	Sky8B	Int20B	Sky27B	WPM72B	POLAR-1.8B	POLAR-7B
General Task	AlignBench	6.76	6.81	6.90	6.98	6.86	7.05	7.02	7.12
	AlpacaEval	37.27	56.19	58.42	60.75	60.71	65.66	61.92	78.88
	Arena-Hard	77.75	70.98	72.03	76.27	78.59	81.34	84.20	89.52
	WildBench	25.09	38.03	36.33	37.39	37.49	40.34	39.50	45.82
	MT-Bench	8.47	8.49	8.54	8.52	8.51	8.54	8.54	8.57
Instruct Following	FollowBench	91.30	94.10	94.00	94.00	94.40	95.00	94.50	94.20
	FoFo	59.70	57.30	65.20	54.90	65.80	62.30	66.00	66.80
Coding	HumanEval	90.24	84.15	90.24	89.02	90.85	90.00	90.89	91.46
	MBPP	82.88	82.10	78.60	82.49	77.04	83.58	84.05	86.38
General Reasoning	DROP	89.24	87.68	89.77	89.24	89.67	89.50	90.52	90.31
	BBH	79.29	79.34	82.50	82.02	83.12	81.43	81.51	82.49
	GPQA	43.43	43.94	46.46	44.44	47.98	47.47	47.47	52.02
	HellaSwag	82.40	83.92	83.62	76.78	85.77	80.24	87.98	87.99
	MuSR	66.70	65.35	64.17	64.56	65.19	64.87	67.64	67.49
	KOR-Bench	57.36	58.08	59.92	58.88	59.84	58.74	58.24	58.72
Math	GSM8K	95.30	95.00	95.38	95.45	91.96	95.56	96.06	95.98
	MATH-500	83.40	79.40	83.20	83.24	81.60	84.24	84.00	84.80
Knowledge	CMLLU	84.41	81.70	82.64	82.86	83.06	83.73	86.02	85.52
	MMLU-Pro	68.38	67.56	69.88	69.82	70.23	69.15	69.94	70.42
	MMMLU-Lite	60.43	45.71	53.79	47.36	54.16	54.26	65.06	64.83
Average		64.49	64.29	66.08	65.25	66.64	67.15	68.55	<b>70.47</b>

Table 11: Ablation study of POLAR on preference evaluation. **w/o PT** refers to the reward model that is fine-tuned solely with human criteria, without our pre-training phase. **w/o PT & Ref** denotes the reward model trained via the traditional non-pre-trained method without reference trajectories. Since strong performance in preference evaluation does not guarantee effective reward signals, we conduct an extra ablation study on RLHF training in Table 2 and Table 12, which provides a stronger validation.

Tasks	POLAR-1.8B			POLAR-7B		
	Origin	w/o PT	w/o PT & Ref	Origin	w/o PT	w/o PT & Ref
<b>Harmlessness</b>	74.2	73.9	63.8	77.7	78.4	74.0
<b>NLP Tasks</b>	68.5	69.9	65.4	71.3	71.3	70.9
<b>Multilingual</b>	80.8	84.6	56.0	80.0	84.0	72.0
<b>Chat</b>	75.1	70.6	65.5	79.1	71.8	63.8
<b>Brainstorming</b>	74.8	72.3	67.9	72.3	79.9	76.1
<b>Role Playing</b>	70.6	60.3	70.6	67.7	70.6	63.2
<b>STEM</b>	79.8	75.0	59.5	81.0	77.4	54.8
<b>Instruct Following</b>	73.1	76.9	50.0	65.4	61.5	61.5
<b>Closed &amp; Open QA</b>	80.3	73.2	54.9	80.3	76.1	64.8
<b>Creative Writing</b>	85.5	81.6	57.9	85.5	86.8	65.8
<b>Coding</b>	65.8	57.9	54.0	68.4	64.5	59.2
<b>Reasoning</b>	69.2	69.1	46.2	69.2	76.9	46.2
<b>Average</b>	74.0	72.6	63.3	76.3	76.5	70.8

Table 12: Ablation study of POLAR in RLHF training. **w/o PT** denotes the reward model fine-tuned solely on human criteria, without any pre-training phase. **w/o PT & Ref** represents the reward model trained via the traditional non-pre-trained method without reference trajectories.

Type	Dataset	Policy Model: InternLM3-8B-Instruct						
		Baseline	Reward Model Size: 1.8B			Reward Model Size: 7B		
			POLAR	w/o PT	w/o PT & Ref	POLAR	w/o PT	w/o PT & Ref
<b>General Task</b>	<b>AlignBench</b>	6.13	6.54	6.39	6.19	6.53	6.39	6.44
	<b>AlpacaEval</b>	51.80	68.57	53.54	58.63	68.70	57.89	53.42
	<b>Arena-Hard</b>	40.79	63.47	49.84	42.71	63.34	54.99	51.68
	<b>WildBench</b>	13.67	40.34	33.79	22.64	39.75	28.06	35.49
	<b>MT-Bench</b>	7.97	8.57	8.28	8.17	8.43	8.35	8.33
<b>Instruct Following</b>	<b>FollowBench</b>	82.20	91.10	84.30	81.80	90.40	89.30	84.00
	<b>FoFo</b>	43.10	54.30	51.80	45.50	56.10	47.20	52.60
<b>Coding</b>	<b>HumanEval</b>	81.10	82.93	82.32	81.71	84.15	81.71	81.10
	<b>MBPP</b>	67.70	73.54	71.60	69.26	75.10	76.26	72.37
<b>General Reasoning</b>	<b>DROP</b>	82.82	86.64	84.22	83.31	85.63	86.86	83.89
	<b>BBH</b>	77.56	77.06	77.76	77.40	77.70	76.41	77.74
	<b>GPQA</b>	32.83	37.37	35.86	36.87	41.41	39.90	38.38
	<b>HellaSwag</b>	88.88	88.78	88.46	88.04	88.87	89.01	87.21
	<b>MuSR</b>	52.28	57.07	54.53	49.90	58.10	57.71	53.32
	<b>KOR-Bench</b>	51.84	53.84	50.56	53.84	55.60	55.28	56.40
<b>Math</b>	<b>GSM8K</b>	90.22	91.05	89.99	90.30	91.96	90.60	90.90
	<b>MATH-500</b>	76.00	77.60	76.60	74.70	78.40	77.00	75.20
<b>Knowledge</b>	<b>CMMLU</b>	81.70	84.79	82.70	81.83	84.39	83.98	82.60
	<b>MMLU-Pro</b>	57.49	60.00	58.28	57.34	59.98	58.92	58.13
	<b>MMMLU-Lite</b>	43.64	48.40	48.15	41.92	49.00	49.28	45.46
<b>Average</b>		56.49	62.60	59.45	57.60	63.18	60.76	59.73

Table 13: Ablation study of POLAR on RFT vs. SFT. **RFT** denotes the reinforcement fine-tuning using our proposed reward models. **SFT** denotes the straightforward supervised fine-tuning. Two training processes employ the same prompt-reference data. Results indicate that RFT with a reward model is significantly more effective than SFT. The substantial improvement in policy model performance cannot be attributed solely to the training data (i.e., prompt-reference data) used in RFT.

Policy Model		InternLM3-8B-Instruct			Qwen2.5-7B-Instruct			Qwen2.5-32B-Instruct		
	Dataset	Baseline	RFT <sub>POLAR-7B</sub>	SFT	Baseline	RFT <sub>POLAR-7B</sub>	SFT	Baseline	RFT <sub>POLAR-7B</sub>	SFT
General Task	AlignBench	6.13	6.53	5.84	6.18	6.49	5.62	6.76	7.12	6.52
	AlpacaEval	51.80	68.70	57.39	37.02	61.86	56.40	37.27	78.88	75.53
	Arena-Hard	40.79	63.34	49.74	56.73	71.38	59.43	77.75	89.52	80.17
	WildBench	13.67	39.75	5.77	24.30	40.30	2.77	25.09	45.82	26.67
	MT-Bench	7.97	8.43	7.61	8.38	8.49	7.60	8.47	8.57	8.05
Instruct Following	FollowBench	82.20	90.40	83.60	88.00	90.30	88.60	91.30	94.20	92.30
	FoFo	43.10	56.10	51.49	44.10	52.00	45.12	59.60	66.80	64.47
Coding	HumanEval	81.10	84.15	78.05	84.15	87.20	76.22	90.24	91.46	86.59
	MBPP	67.70	75.10	72.37	74.32	75.10	69.65	82.88	86.38	79.77
General Reasoning	DROP	82.82	85.63	87.12	78.35	83.22	84.71	89.24	90.31	90.54
	BBH	77.56	77.70	68.22	56.92	67.14	59.76	79.29	82.49	74.99
	GPQA	32.83	41.41	33.33	32.32	37.37	27.78	43.43	52.02	42.93
	HellaSwag	88.88	88.87	88.96	67.47	81.25	84.16	82.40	87.99	90.73
	MuSR	52.28	58.10	53.45	46.53	50.64	46.01	66.70	67.49	62.87
	KOR-Bench	51.84	55.60	42.48	41.36	48.16	44.16	57.36	58.72	45.12
Math	GSM8K	90.22	91.96	89.76	91.13	92.19	88.63	95.30	95.98	95.45
	MATH-500	76.00	78.40	71.00	75.80	77.20	71.60	83.40	84.80	80.40
Knowledge	CMMLU	81.70	84.39	82.03	77.56	77.92	74.23	84.41	85.52	83.62
	MMLU-Pro	57.49	59.98	54.74	55.33	56.65	52.42	68.38	70.42	66.37
	MMMLU-Lite	43.64	49.00	45.90	53.04	53.15	48.42	60.43	64.83	63.28
Average		56.49	63.18	56.44	54.95	60.90	54.66	64.49	70.47	65.82