

Impact of Multitask Fine-Tuning on Qwen2.5 Alignment Capabilities via AlpacaEval and Chatbot Arena Preference Scores

Assignee Research

June 11, 2026

Abstract

Reward models (RMs) play a crucial role in Reinforcement Learning from Human Feedback by serving as proxies for human preferences in aligning large language models. However, they suffer from various biases which could lead to reward hacking. In this paper, we identify a model preference bias in RMs, where they systematically assign disproportionately high scores to responses from certain policy models, leading to unfair judgments. To mitigate this bias, we propose a calibration method named CHatbot Arena calibrated Reward Modeling (CHARM) that leverages Elo scores from the Chatbot Arena to con

1 Introduction

This paper examines: CHARM: Calibrating Reward Models With Chatbot Arena Scores. Research question: What is the effect of multitask fine-tuning on the alignment capabilities of Qwen2.5, as evaluated by preference scores on AlpacaEval and Chatbot Arena?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

3 Results

15 papers retrieved. 22 claims extracted; 19 independently verified. Quality review score: 7.2/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
CHARM provides a simple, effective, and broadly applicable approach to building more reliable and fair reward models.	✓	0.33
The code for CHARM is available at https://github.com/HexagonStar/CHARM .	×	0.15
Benchmarks such as MT-Bench, Alpaca-Eval, and Arena-hard employ LLM-as-a-judge systems to assess model response quality	✓	0.21
ChatBot Arena employs a crowdsourced, pairwise comparison system where users select preferred responses from two anonymo	✓	0.20
ChatBot Arena aggregates human preference data to compute a dynamic Elo rating for each model.	✓	0.17
Park et al. (2024) identified six distinct types of bias in evaluation models.	✓	0.28
Li et al. (2025) found that judge models may develop bias favoring content generated by themselves or closely related LL	✓	0.27
Dubois et al. (2024) proposed a regression-based method to mitigate length bias.	✓	0.30
Huang et al. (2025) introduced a post hoc calibration technique for reward models.	✓	0.22
Chatbot Arena represents a universal distribution of real-world prompts.	✓	0.16
The AlpacaEval dataset consists of 805 carefully curated questions.	✓	0.17
The AlpacaEval dataset exhibits a 98% Spearman correlation with Chatbot Arena.	×	0.15
The study evaluated five popular reward models on a diverse set of policy models with varying Elo scores.	✓	0.21
Some models on AlpacaEval did not participate in the Chatbot Arena, resulting in unavailable Elo scores for them.	✓	0.19
Reward model scores correlate positively with human preferences.	×	0.14
Calibrated RMs achieve improved evaluation accuracy on RM-Bench.	✓	0.24
Calibrated RMs achieve improved evaluation accuracy on the Chat-Hard domain of Reward-Bench.	✓	0.26
Calibrated RMs exhibit a stronger correlation with human preferences by producing scores more closely aligned with Elo r	✓	0.32
Calibrated RMs improve downstream post-training performance.	✓	0.19
In the Elo rating system implementation described, a model receives a score of 1 for a win, 0.5 for a tie, and 0 for a loss.	✓	0.18

References

- <http://arxiv.org/abs/2504.10045v2>
- <http://arxiv.org/abs/2402.18571v3>
- <http://arxiv.org/abs/2402.11690v1>