

Scalability of DPO and RLHF in Large Multimodal Reasoning Benchmarks

Assignee Research

June 1, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the scalability of DPO compare to RLHF in terms of training throughput when applied to large multimodal reasoning benchmarks like MMBench or SEED-Bench. This paper studies the alignment process of generative models with Reinforcement Learning from Human Feedback (RLHF). We first identify the primary challenges of existing popular methods like offline PPO and offline DPO as lacking in strategical exploration of the environment. 8 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.1/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Iterative Preference Learning from Human Feedback: Bridging Theory and Practice for RLHF under KL-Constraint. Research question: How does the scalability of DPO compare to RLHF in terms of training throughput when applied to large multimodal reasoning benchmarks like MMBench or SEED-Bench?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.1/10.

3 Results

13 papers retrieved. 8 claims extracted; 0 independently verified. Quality review score: 3.1/10.

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
Under Assumption 1, with probability at least $1-\delta$, the suboptimality gap $J(\pi) - J(\pi^*)$ for Algorithm 1 Option I is bounded	×	0.01
Under Assumption 1, with probability at least $1-\delta$, the suboptimality gap $J(\pi) - J(\pi^*)$ for Algorithm 1 Option II includes	×	0.02
By Jensen’s inequality, the uncertainty bonus bound for Option I ($\mathbb{E}[\text{sim}_0(x, \pi)] - \nu$) is less than or equal to the bound	×	0.03
If the reference vector ν is set to $\mathbb{E}[\text{sim}_0(x, \pi^{\text{ref}})]$, the resulting policy from Option I is theoretically guaranteed to	×	0.02
In best-of-n sampling, n independent responses are sampled by policy π_{1_t} for each prompt, and the response with the highest	×	0.03
In best-of-n sampling, the KL divergence between the initial policy π_{1_t} and the resulting policy π_{2_t} is upper bounded	×	0.04
The LLaMA2 project adjusts the sampling temperature of policy π_{1_t} to induce policy π_{2_t} .	×	0.02
Offline learning in the context of RLHF is defined as learning without further querying human feedback.	×	0.13

References

- <http://arxiv.org/abs/2407.14477v4>
- <http://arxiv.org/abs/2407.04973v1>
- <http://arxiv.org/abs/2312.11456v4>