

# Direct Preference Optimization and RLHF in LLM Fine-Tuning: Sample Efficiency and Convergence on SQuTR with Noisy Inputs

Assignee Research

May 31, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does Direct Preference Optimization (DPO) compare to RLHF in terms of sample efficiency and convergence speed when fine-tuning LLMs on the SQuTR benchmark with noisy inputs. Aligning language models with human preferences through reinforcement learning from human feedback is crucial for their safe and effective deployment. The human preference is typically represented through comparison where one response is chosen over another for a given prompt. 11 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Data-Centric Human Preference with Rationales for Direct Preference Alignment. Research question: How does Direct Preference Optimization (DPO) compare to RLHF in terms of sample efficiency and convergence speed when fine-tuning LLMs on the SQuTR benchmark with noisy inputs?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

### **3 Results**

12 papers retrieved. 11 claims extracted; 0 independently verified. Quality review score: 3.7/10.

### **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
The study evaluates the impact of rationales on direct preference learning through multiple experiments.	×	0.14
The study uses three preference datasets: Orca DPO Pairs, UltraFeedback, and Anthropic Helpful and Harmless.	×	0.06
Each dataset has 512 fixed samples as the test set for winrate evaluations.	×	0.03
The models investigated include Mistral-7B-v0.1, Mistral-7B-Instruct-v0.2, Zephyr-7B-Beta, and Llama3-8B-Instruct.	×	0.01
GPT-4o is used as a judge to evaluate the responses generated by the models and to retrieve the winrate scores.	×	0.02
The study integrates rationales into preference learning frameworks such as DPO, ORPO, and SimPO.	×	0.07
RDPO shows better performance and 3x annotation saving compared to DPO.	×	0.02
The winrate of Mistral-7B-Instruct-v0.2 with RDPO is 27.55 compared to 19.52 with DPO.	×	0.00
The winrate of Llama-3.1-8B-Instruct with RDPO is 27.55 compared to 22.42 with DPO.	×	0.00
The study presents a demonstration of extending the direct preference optimization (DPO) algorithm to incorporate ration	×	0.10
The study analyzes the possible impact of rationales through the perspective of information theory.	×	0.04

## References

- <http://arxiv.org/abs/2408.07888v2>
- <http://arxiv.org/abs/2402.07314v3>
- <http://arxiv.org/abs/2407.14477v4>