

# Mitigating Token-Level Misalignment in Generative Information Retrieval via Reinforcement Learning from Human Feedback on

Assignee Research

June 15, 2026

## Abstract

Generative information retrieval (GenIR) is a promising neural retrieval paradigm that formulates document retrieval as a document identifier (docid) generation task, allowing for end-to-end optimization toward a unified global retrieval objective. However, existing GenIR models suffer from token-level misalignment, where models trained to predict the next token often fail to capture document-level relevance effectively. While reinforcement learning-based methods, such as reinforcement learning from relevance feedback (RLRF), aim to address this misalignment through reward modeling, they intro

## 1 Introduction

This paper examines: Lightweight and Direct Document Relevance Optimization for Generative Information Retrieval. Research question: To what extent does reinforcement learning from human feedback mitigate token-level misalignment in generative information retrieval systems evaluated on cross-lingual datasets?.

## 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 7.4/10.

## 3 Results

13 papers retrieved. 13 claims extracted; 10 independently verified. Quality review score: 7.4/10.

## 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
DDRO improves retrieval accuracy on MS-MARCO and Natural Questions benchmarks, outperforming multiple baselines.	×	0.10
DDRO maintains competitive performance with established baselines on broader metrics like R@10.	✓	0.18
An ablation study highlights the contributions of pairwise ranking optimization to the observed performance improvements	✓	0.24
DDRO is a pairwise ranking approach that aligns docid generation with document-level relevance judgments.	✓	0.23
DDRO unifies training objectives within a single framework, optimizing directly for document-level relevance.	✓	0.23
Experimental results demonstrate improvements in retrieval accuracy with DDRO.	×	0.14
DDRO eliminates the need for explicit reward model training and reinforcement learning fine-tuning.	✓	0.25
DDRO reduces computational overhead and improves optimization efficiency compared to GenRRL.	×	0.07
ROGER combines dense and generative retrieval by using dense retrievers as intermediaries to provide relevance signals.	✓	0.27
ROGER relies on external dense retrievers and does not directly optimize for document-level relevance within the generat	✓	0.32
DDRO eliminates dependency on external dense retrievers by incorporating pairwise ranking directly into the generative m	✓	0.25
The SFT phase serves as a pretraining step, aligning the model with initial relevance signals from training data.	✓	0.27
The proposed workflow comprises three key stages: (1) Construction of document identifiers (docids), (2) Supervised fine	✓	0.39

## References

- <http://arxiv.org/abs/2504.05181v2>
- <http://arxiv.org/abs/2407.14477v4>
- <http://arxiv.org/abs/2402.18571v3>