

SOVEREIGN: Evaluating Multi-Hop Reasoning in RAG Systems: A Comparison of LLM-Based Retrieval

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 27, 2026

Abstract

Retrieval-augmented generation (RAG) enhances large language models (LLMs) with external knowledge to answer questions more accurately. However, research on evaluating RAG systems-particularly the retriever component-remains limited, as most existing work focuses on single-context retrieval rather than multi-hop queries, where individual contexts may appear irrelevant in isolation but are essential when combined. In this research, we use the HotPotQA, MuSiQue, and SQuAD datasets to simulate a RAG system and compare three LLM-as-judge evaluation strategies, including our proposed Context-Awar

1 Introduction

Analysis of: Evaluating Multi-Hop Reasoning in RAG Systems: A Comparison of LLM-Based Retriever Evaluation Strategies. Research goal: Does the accuracy gain from extending context windows to 128K tokens saturate beyond a certain retrieval step count (e.g., 3 steps) for multi-hop reasoning on HotPotQA, and how does this trade-off vary across model scales (7B vs 70B parameters)?.

2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

3 Results

11 papers retrieved. 9 claims extracted, 9 verified. Tribunal: 7.8/10 → APPROVE (revision_round=0). Policy: AUTO_APPROVE.

4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

5 Extracted Claims

Claim	Verified	Confidence
Retrieval-augmented generation (RAG) enhances large language models (LLMs) with external knowledge to answer questions in	✓	0.30
Research on evaluating RAG systems, particularly the retriever component, remains limited, as most existing work focuses	✓	0.43
In this research, the HotPotQA, MuSiQue, and SQuAD datasets are used to simulate a RAG system.	✓	0.17
Three LLM-as-judge evaluation strategies are compared, including the proposed Context-Aware Retriever Evaluation (CARE).	✓	0.31
Experiments with LLMs from OpenAI, Meta, and Google demonstrate that CARE consistently outperforms existing methods for	✓	0.42
The performance gains of CARE are most pronounced in models with larger parameter counts and longer context windows.	✓	0.24
Single-hop queries show minimal sensitivity to context-aware evaluation.	✓	0.29
The results highlight the critical role of context-aware evaluation in improving the reliability and accuracy of retrieval	✓	0.39
The complete data of the experiments is provided at https://github.com/lorenzbrehme/CARE .	✓	0.21

References

- <http://arxiv.org/abs/2605.02173v1>
- <http://arxiv.org/abs/2604.18234v1>
- <https://www.semanticscholar.org/paper/7a39d0ccffe0c30b98a3a4b2c1bb82bb0e8c6026>