

Cross-Domain Generalization in RAG Systems: Sparse vs. Dense Retrieval Evaluation on TriviaQA and MusicQA

Assignee Research

June 12, 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

This paper examines: Evaluating Multi-Hop Reasoning in RAG Systems: A Comparison of LLM-Based Retriever Evaluation Strategies. Research question: What is the impact of sparse vs. dense retrieval combinations in RAG systems on cross-domain generalization when evaluated on datasets like TriviaQA and MusicQA?.

2 Methodology

Systematic literature search across multiple databases yielded 6 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.0/10.

3 Results

6 papers retrieved. 14 claims extracted; 10 independently verified. Quality review score: 7.0/10.

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 |
|---|----------|------------|
| CARE consistently outperforms existing methods for evaluating multi-hop reasoning in RAG systems. | ✓ | 0.30 |
| 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 indirect evaluation approach derived from the eRAG method [24] involves comparing generated answers with ground truth | ✓ | 0.19 |
| The direct evaluation method based on the ARES framework [23] involves determining if a context is crucial to answering | ✓ | 0.23 |
| The CARE method involves determining if a context is crucial to answering a question with a ground truth answer, given a | ✓ | 0.17 |
| In the indirect method, an LLM attempts to answer the query using only a single context document, and if the answer is e | ✓ | 0.26 |
| CARE outperformed other approaches across all models except for the LLaMa 3.1-8b model. | × | 0.09 |
| The LLaMa 3.1-8b model experienced a significant decline in overall performance, with substantial drops in both F1-Score | × | 0.11 |
| The indirect approach led to a significant improvement in F1-Score for the small LLaMa model. | ✓ | 0.30 |
| The direct approach resulted in a decline in F1-Score for the reasoning model o4-mini. | ✓ | 0.33 |
| For CARE, the reasoning model o4-mini exhibited a decrease in accuracy, F1-Score, and recall compared to GPT-4.1, while | ✓ | 0.30 |
| Among smaller models, we observed a shift from precision to recall. | × | 0.06 |
| CARE demonstrated stable performance when context size and parameters were varied. | × | 0.03 |

References

- <http://arxiv.org/abs/2404.07220v2>
- <http://arxiv.org/abs/2604.18234v1>
- <http://arxiv.org/abs/2510.25518v1>