

Adversarial Robustness in Retrieval-Augmented Generation for Quranic Studies Across Open-Source LLM Scales and Architectural

Assignee Research

June 12, 2026

Abstract

Accurate and contextually faithful responses are critical when applying large language models (LLMs) to sensitive and domain-specific tasks, such as answering queries related to quranic studies. General-purpose LLMs often struggle with hallucinations, where generated responses deviate from authoritative sources, raising concerns about their reliability in religious contexts. This challenge highlights the need for systems that can integrate domain-specific knowledge while maintaining response accuracy, relevance, and faithfulness. In this study, we investigate 13 open-source LLMs categorized in

1 Introduction

This paper examines: Investigating Retrieval-Augmented Generation in Quranic Studies: A Study of 13 Open-Source Large Language Models. Research question: How does adversarial robustness in retrieval-augmented generation (measured by accuracy drop under perturbed or misleading prompts) vary across different open-source LLMs (7B vs. 70B) when applied to domain-specific tasks like Quranic studies, and what architectural modifications (e.g., attention mechanisms) mitigate this?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.5/10.

3 Results

16 papers retrieved. 9 claims extracted; 7 independently verified. Quality review score: 7.5/10.

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 investigates 13 open-source Large Language Models in the context of Quranic studies.	✓	0.17
The system employs a Retrieval-Augmented Generation (RAG) architecture combining retrieval-based and generative methods.	✓	0.21
The system performs semantic similarity search over a vectorized dataset obtained from Qur'anic surah descriptions.	✓	0.25
Generated responses include references to original dataset entries, such as surah descriptions or specific virtues.	✓	0.23
Human evaluators assessed response quality based on three dimensions: Context Relevance, Answer Faithfulness, and Answer	✓	0.16
Context Relevance is calculated using the precision@k metric.	✓	0.16
The evaluation platform logged and stored data for research purposes after responses were evaluated, scored, and comment	✓	0.19
The dataset selection criteria included authenticity, descriptive richness, clarity and accessibility, and relevance.	×	0.13
The dataset source was reviewed to confirm compliance with recognized Islamic scholarship and the absence of speculative	×	0.02

References

- <http://arxiv.org/abs/2503.16581v1>
- <http://arxiv.org/abs/2502.00306v2>
- <http://arxiv.org/abs/2410.05451v3>