

# Retrieval-Augmented Generation Performance Across Sparse Dense and Hybrid Retrieval Methods

Assignee Research

June 9, 2026

## **Abstract**

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does the retrieval-augmented generation performance compare between sparse (BM25), dense, and hybrid retrieval methods when evaluated on a benchmark of religious text question-answering tasks. 6 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## **1 Introduction**

This paper examines: Almanac — Retrieval-Augmented Language Models for Clinical Medicine. Research question: How does the retrieval-augmented generation performance compare between sparse (BM25), dense, and hybrid retrieval methods when evaluated on a benchmark of religious text question-answering tasks using metrics such as exact match accuracy and BLEU score?.

## **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 5.0/10.

## **3 Results**

16 papers retrieved. 6 claims extracted; 6 independently verified. Quality review score: 5.0/10.

## 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
Large language models (LLMs) have shown impressive zero-shot capabilities in natural language tasks such as summarization	✓	0.28
LLMs have promising applications in clinical medicine but are limited by their tendency to generate incorrect and sometimes	✓	0.26
Almanac is an LLM framework augmented with retrieval capabilities from curated medical resources for medical guideline a	✓	0.31
A panel of eight board-certified clinicians and two health care practitioners evaluated Almanac and compared it with sta	✓	0.40
Almanac showed a significant improvement in performance compared with the standard LLMs across axes of factuality, compl	✓	0.34
The study was funded by the National Institutes of Health, National Heart, Lung, and Blood Institute.	✓	0.25

## References

- <https://doi.org/10.15575/join.v11i1.1733>
- <https://doi.org/10.1145/3637528.3672065>
- <https://doi.org/10.1056/aioa2300068>