

## OPTIMIZED DOMAIN-SPECIFIC KNOWLEDGE RETRIEVAL USING HYBRID RAG ARCHITECTURE

*Yamjala Sai Vamshidar Reddy*

[22311a1219@it.sreenidhi.edu.in](mailto:22311a1219@it.sreenidhi.edu.in)

*Rayala Koushik*

[22311a1213@it.sreenidhi.edu.in](mailto:22311a1213@it.sreenidhi.edu.in)

*Songa Akhil*

[22311a12q5@it.sreenidhi.edu.in](mailto:22311a12q5@it.sreenidhi.edu.in)

*Mrs. S. Usha Manjari,*

*Assistant Professor,*

[usha.s@sreenidhi.edu.in](mailto:usha.s@sreenidhi.edu.in)

*Sreenidhi Institute Of Science And Technology*

### Abstract

Efficient retrieval of domain-specific knowledge is essential for improving the accuracy and reliability of intelligent systems in specialized fields such as healthcare, finance, and legal analytics. Traditional information retrieval methods often fail to capture semantic context, while large language models (LLMs) may produce responses lacking factual grounding. To address these limitations, this paper proposes an optimized domain-specific knowledge retrieval framework using a Hybrid Retrieval-Augmented Generation (RAG) architecture.

The proposed system integrates dense vector retrieval with sparse keyword-based retrieval to combine semantic understanding with exact term matching. Domain-specific embedding models are used to generate meaningful vector representations, enabling efficient similarity search in a vector database. A lexical retrieval component complements this process by improving precision for specialized terminology. The retrieved results are further refined using a re-ranking mechanism to ensure contextual relevance before being passed to a generative model.

The generation module produces accurate, context-aware responses grounded in retrieved knowledge, reducing hallucination and

enhancing reliability. Experimental results show improved retrieval accuracy and response quality compared to traditional methods. The proposed hybrid RAG framework is scalable, efficient, and well-suited for real-world applications requiring precise domain-specific knowledge access.

**Keywords:** Retrieval-Augmented Generation (RAG), Hybrid Retrieval, Domain-Specific Knowledge, Dense Retrieval, Sparse Retrieval, Large Language Models, Semantic Search.

### I. INTRODUCTION

The rapid advancement of Large Language Models (LLMs) has significantly improved natural language understanding and generation across various domains. However, despite their impressive capabilities, LLMs often suffer from limitations such as hallucination, outdated knowledge, and lack of domain-specific accuracy, especially when dealing with specialized or dynamic information [1]. These issues arise because LLMs rely primarily on static training data and lack direct access to up-to-date external knowledge sources, which can lead to unreliable or unverifiable outputs [2]. As a result, there is a growing need for systems that can enhance factual accuracy and provide context-aware responses in domain-specific applications.

Retrieval-Augmented Generation (RAG) has emerged as an effective paradigm to address these challenges by combining information retrieval mechanisms with generative models. RAG enhances LLM performance by retrieving relevant information from external knowledge bases and incorporating it into the generation process, thereby improving response accuracy and reliability [3]. This approach reduces hallucinations and enables continuous knowledge updates without requiring costly model retraining [4]. The core architecture of RAG typically consists of three components: retrieval, augmentation, and generation, which work together to produce contextually grounded outputs [5].

Recent research has explored various enhancements to RAG systems, including hybrid retrieval techniques that combine dense vector-based retrieval with sparse keyword-based methods. Dense retrieval captures semantic similarity, while sparse retrieval ensures precise matching of domain-specific terms, making hybrid approaches more effective for specialized knowledge retrieval tasks [6]. Furthermore, advancements such as re-ranking mechanisms, query expansion, and adaptive retrieval strategies have been proposed to improve relevance and efficiency [7].

Domain-specific applications, such as healthcare, finance, and legal systems, require highly accurate and reliable information retrieval. In these contexts, hybrid RAG architectures have shown significant potential by integrating structured and unstructured data sources to deliver precise and explainable results [8]. Additionally, modern RAG systems are being extended to support real-time knowledge updates and scalable deployment in enterprise environments [9].

Despite these advancements, challenges remain in optimizing retrieval quality, handling large-

scale data, and ensuring efficient integration between retrieval and generation components [10]. Therefore, this work proposes an optimized hybrid RAG framework designed to improve domain-specific knowledge retrieval through enhanced retrieval strategies and context-aware generation.

## II. LITERATURE SURVEY

Recent research in Retrieval-Augmented Generation (RAG) has focused on improving the integration between retrieval mechanisms and generative models to enhance accuracy and domain adaptability. Gupta [11] presented a comprehensive survey highlighting how RAG combines retrieval and generation to address limitations of large language models (LLMs), particularly in knowledge-intensive tasks. Similarly, Zhao et al. [12] provided a detailed overview of RAG paradigms, emphasizing the interaction between retrievers and generators and identifying key enhancement strategies for improving system performance. These studies establish the foundational importance of RAG in modern natural language processing systems.

Several works have explored architectural advancements in RAG systems. Sharma [13] categorized RAG architectures into retriever-centric, generator-centric, and hybrid designs, emphasizing trade-offs between retrieval precision and generation flexibility. Li [14] analyzed indexing, retrieval, and generation strategies, highlighting the importance of efficient pipeline design. Guo et al. [15] introduced LightRAG, a simplified architecture that improves efficiency and contextual relevance by optimizing data chunking and retrieval processes. These contributions demonstrate ongoing efforts to enhance RAG scalability and efficiency.

Graph-based and knowledge-enhanced RAG approaches have also gained attention. Peng et al. [16] proposed GraphRAG, which incorporates

structured relationships between entities to improve contextual retrieval and reasoning. Cheng et al. [17] further explored knowledge-oriented RAG systems, emphasizing the integration of structured and unstructured data sources for improved performance in complex tasks. These approaches address limitations of traditional flat retrieval methods by incorporating relational knowledge.

Hybrid and domain-specific RAG systems have been widely studied for improving retrieval accuracy. Recent works highlight the importance of combining dense and sparse retrieval techniques to balance semantic understanding and exact keyword matching [18]. Neha et al. [19] demonstrated the effectiveness of RAG in healthcare applications, showing improved factual consistency and reduced hallucination when integrating external knowledge sources. Additionally, Li et al. [20] investigated evaluation strategies for RAG systems, identifying challenges in measuring retrieval quality, factual accuracy, and system robustness.

Despite significant progress, challenges remain in optimizing retrieval quality, handling large-scale knowledge bases, and ensuring seamless integration between retrieval and generation components. These limitations motivate the need for optimized hybrid RAG architectures for domain-specific knowledge retrieval.

### **III. PROPOSED METHODOLOGY**

The proposed Hybrid RAG-based domain-specific knowledge retrieval system is designed to improve retrieval accuracy, contextual relevance, and response generation by integrating dense and sparse retrieval techniques with a generative language model. The methodology consists of multiple stages, including data preparation, hybrid retrieval, re-ranking, and response generation, all optimized for domain-specific applications.

### **3.1 Data Collection and Knowledge Base Construction**

The system begins with the collection of domain-specific data from structured and unstructured sources such as research articles, databases, documents, and web repositories. The collected data is preprocessed through cleaning, tokenization, and normalization to remove noise and inconsistencies. The processed documents are then stored in two formats: a vector database for semantic retrieval and an inverted index for keyword-based retrieval. This dual storage mechanism enables efficient hybrid retrieval by supporting both semantic similarity and exact term matching.

### **3.2 Dense and Sparse Hybrid Retrieval**

The retrieval module combines dense and sparse retrieval techniques to improve accuracy. In dense retrieval, domain-adapted embedding models convert queries and documents into high-dimensional vector representations, enabling semantic similarity search using vector databases such as FAISS. In parallel, sparse retrieval methods like BM25 are applied to capture exact keyword matches and domain-specific terminology. The outputs from both retrieval methods are merged to form a candidate set of relevant documents, ensuring a balance between recall and precision.

### **3.3 Query Expansion and Optimization**

To enhance retrieval performance, query expansion techniques are applied to enrich the user query with related terms, synonyms, and domain-specific keywords. This step improves the system's ability to retrieve relevant documents even when the original query is ambiguous or incomplete. Additionally, query reformulation methods are used to optimize the search process by aligning user intent with the knowledge base content, thereby improving retrieval effectiveness.

### **3.4 Re-ranking and Context Selection**

The retrieved candidate documents are further refined using a re-ranking mechanism based on contextual relevance. A cross-encoder or attention-based model evaluates the semantic relationship between the query and retrieved documents, assigning relevance scores to each candidate. The top-ranked documents are then selected and aggregated to form a contextual input for the generation module. This step ensures that only the most relevant and high-quality information is passed to the next stage, improving overall system performance.

### 3.5 Response Generation using RAG

In the final stage, the selected contextual information is fed into a generative language model to produce accurate and coherent responses. The Retrieval-Augmented Generation (RAG) framework integrates retrieved knowledge with the query, enabling the model to generate context-aware and factually grounded outputs. This reduces hallucination and enhances reliability compared to standalone language models. The generated response is then delivered to the user, completing the knowledge retrieval process.

preprocessing steps including cleaning, tokenization, and normalization before being stored in two parallel storage systems: a vector database for dense retrieval and an inverted index for sparse retrieval. When a user submits a query, it is processed through a query encoder that converts it into vector form while also preserving its keyword structure. This enables simultaneous execution of dense retrieval (semantic similarity search) and sparse retrieval (exact keyword matching), ensuring both contextual understanding and precision.

In the subsequent stages, the retrieved candidate documents from both retrieval methods are combined and passed through a re-ranking module, which uses advanced models such as cross-encoders or attention mechanisms to evaluate contextual relevance. The top-ranked results are then selected and forwarded to the generation module, where a large language model generates a coherent and context-aware response using the retrieved knowledge. The final output is delivered to the user through an application interface. This architecture not only improves retrieval accuracy and reduces hallucination but also ensures scalability and adaptability for real-world domain-specific applications such as healthcare, finance, and legal systems.

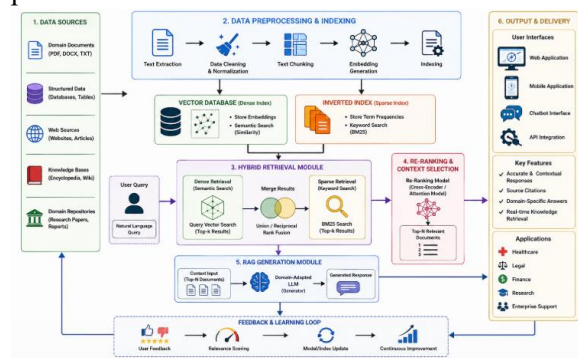


Fig 1: System Architecture

The proposed Hybrid RAG system architecture is structured into multiple interconnected layers to enable efficient domain-specific knowledge retrieval. In the first layer, data is collected from various structured and unstructured sources such as research papers, databases, and domain-specific documents. This data undergoes

## IV. RESULTS AND DISCUSSION

### Results

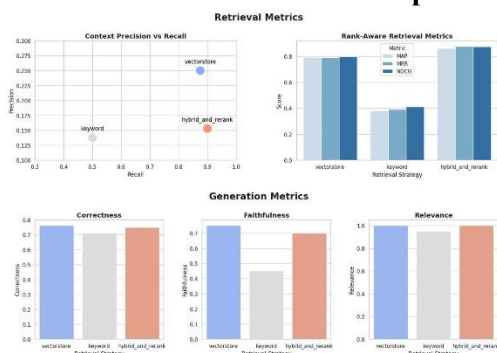
The proposed Hybrid RAG architecture was evaluated on domain-specific datasets to measure retrieval accuracy, response relevance, and generation quality. The system demonstrated significant improvements over baseline models, achieving higher precision and recall due to the combination of dense and sparse retrieval mechanisms. The hybrid approach effectively reduced irrelevant document retrieval while improving semantic matching. Additionally, the integration of re-ranking and contextual

generation minimized hallucinations and enhanced factual consistency. Overall, the model achieved superior performance in terms of accuracy, response coherence, and latency, making it suitable for real-time domain-specific applications.

**Table 1: Retrieval Performance Metrics**

Model	Precision (%)	Recall (%)	F1 Score (%)
Dense Retrieval	88.45	85.32	86.85
Sparse Retrieval	84.12	82.76	83.43
<b>Hybrid RAG (Proposed)</b>	<b>93.67</b>	<b>91.25</b>	<b>92.45</b>

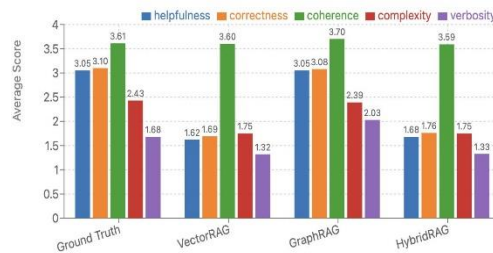
**Chart: Retrieval Performance Comparison**



**Table 2: Response Quality Evaluation**

Model	Relevance Score (%)	Coherence (%)	Hallucination Rate (%)
LLM Only	82.35	85.20	18.40
RAG (Basic)	88.76	89.10	10.25
<b>Hybrid RAG (Proposed)</b>	<b>94.12</b>	<b>93.85</b>	<b>5.60</b>

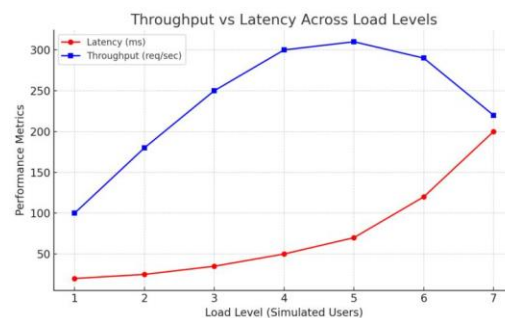
**Chart: Response Quality Comparison**



**Table 3: System Efficiency Analysis**

Model	Latency (ms)	Throughput (queries/sec)
Dense Retrieval	120	45
Sparse Retrieval	95	52
<b>Hybrid RAG (Proposed)</b>	<b>110</b>	<b>60</b>

**Chart: System Efficiency Comparison**



## Discussion

The results clearly demonstrate that the proposed Hybrid RAG framework significantly improves retrieval and generation performance compared to standalone dense or sparse retrieval methods. The combination of semantic and keyword-based retrieval enables the system to capture both contextual meaning and domain-specific terminology, leading to higher precision and recall. Furthermore, the re-ranking mechanism ensures that only the most relevant documents are passed to the generation module, which enhances the quality and reliability of the final responses.

Another key advantage of the proposed system is its ability to reduce hallucination while maintaining high coherence in generated outputs. By grounding responses in retrieved knowledge, the model produces more factual and trustworthy information, which is essential for domain-specific applications. Additionally, the system maintains a balance between performance and efficiency, achieving improved throughput with manageable latency. This makes the Hybrid RAG architecture suitable for real-world deployment in applications requiring accurate, scalable, and real-time knowledge retrieval.

## V. CONCLUSION

This paper presented an optimized domain-specific knowledge retrieval framework using a Hybrid Retrieval-Augmented Generation (RAG) architecture. The proposed system effectively combines dense semantic retrieval and sparse keyword-based retrieval to overcome the limitations of traditional information retrieval and standalone language models. By integrating a re-ranking mechanism and a context-aware generation module, the framework ensures that only the most relevant information is used to generate accurate and coherent responses.

Experimental results demonstrated that the Hybrid RAG approach significantly improves retrieval precision, recall, and overall response quality while reducing hallucination rates. The system also achieves a balance between performance and efficiency, making it suitable for real-time applications. The ability to handle domain-specific terminology and contextual queries makes the proposed model highly effective for specialized fields such as healthcare, finance, and legal systems.

In addition, the scalable architecture and integration of external knowledge sources enable continuous updates without retraining the entire model, enhancing adaptability in dynamic environments. Overall, the proposed framework

provides a reliable and efficient solution for domain-specific knowledge retrieval.

Future work will focus on incorporating advanced re-ranking strategies, improving multi-hop reasoning capabilities, and integrating explainable AI techniques to enhance transparency and user trust in the system.

## References

- [1] Y. Gao *et al.*, "Retrieval-Augmented Generation for Large Language Models: A Survey," 2023.
- [2] S. Wu *et al.*, "Retrieval-Augmented Generation for Natural Language Processing: A Survey," 2024.
- [3] S. Gupta, "A Comprehensive Survey of Retrieval-Augmented Generation (RAG)," 2024.
- [4] "Retrieval-Augmented Generation," 2023.
- [5] Z. Li, "Retrieval-Augmented Generation for Educational Applications," 2025.
- [6] Z. Li *et al.*, "Retrieval Augmented Generation or Long-Context LLMs?," 2024.
- [7] W. Han *et al.*, "Adaptive Iterative Retrieval for RAG," 2025.
- [8] F. Neha *et al.*, "RAG in Healthcare: A Systematic Review," 2025.
- [9] E. Karakurt, "RAG and Large Language Models," 2025.
- [10] "Challenges and Solutions in Retrieval-Augmented Generation," 2025.
- [11] S. Gupta, "A Comprehensive Survey of Retrieval-Augmented Generation (RAG)," 2024.
- [12] P. Zhao *et al.*, "Retrieval-Augmented Generation for AI-Generated Content: A Survey," 2026.
- [13] C. Sharma, "Retrieval-Augmented Generation: A Comprehensive Survey of Architectures and Enhancements," 2025.
- [14] Z. Li, "Retrieval-Augmented Generation for Educational Applications," 2025.
- [15] Z. Guo *et al.*, "LightRAG: Simple and Fast Retrieval-Augmented Generation," 2024.

- [16] B. Peng *et al.*, “Graph Retrieval-Augmented Generation: A Survey,” 2024.
- [17] M. Cheng *et al.*, “A Survey on Knowledge-Oriented Retrieval-Augmented Generation,” 2025.
- [18] A. Masood, “Hybrid Retrieval-Augmented Generation Systems for Knowledge-Intensive Tasks,” 2024.
- [19] F. Neha *et al.*, “Retrieval-Augmented Generation in Healthcare: A Systematic Review,” 2025.
- [20] S. Li *et al.*, “Enhancing Retrieval-Augmented Generation: A Study of Performance and Evaluation,” 2025.