# [Vol-11, Issue-10, October- 2025]

# Efficient in-Domain Research Query Resolution using Retrieval **Augmented Generation with Ollama**

Jothi Muneeswari<sup>1\*</sup>; U Saravana Kumar<sup>2</sup>

Data Scientist, India \*Corresponding Author

Received: 01 October 2025/ Revised: 08 October 2025/ Accepted: 13 October 2025/ Published: 31-10-2025 Copyright @ 2025 International Journal of Engineering Research and Science This is an Open-Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (https://creativecommons.org/licenses/by-nc/4.0) which permits unrestricted Non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract— Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm for enhancing the performance of Large Language Models (LLMs) by grounding their outputs on external knowledge. This paper presents a domain-specific RAG pipeline integrating Ollama with LangChain, FAISS, and Hugging Face embeddings to process and query a custom corpus of 100 research papers. By leveraging FAISS for efficient similarity search and Hugging Face models for semantic embeddings, the system enables precise and context-aware retrieval of academic knowledge. The results demonstrate improved accuracy, contextual relevance, and reduced hallucinations compared to traditional LLM usage, making the framework suitable for research assistance and literature review automation.

Keywords— Retrieval-Augmented Generation, Ollama, LangChain, FAISS, Hugging Face, Research Assistance, Large Language Model.

#### I. INTRODUCTION

The exponential growth of scientific publications has created challenges in efficiently accessing and synthesizing knowledge across domains [1]. Traditional search engines often fail to provide context-rich responses, while standalone LLMs tend to generate hallucinations due to their limited knowledge cutoff [2]. Retrieval-Augmented Generation (RAG)[3] combines neural retrieval with generative capabilities to overcome these limitations. This paper investigates the implementation of a custom RAG system using Ollama [4] for model inference, LangChain for orchestration, FAISS for vector-based similarity search, and Hugging Face embeddings for semantic representation [3]. The system is applied to a dataset of 100 academic research papers to demonstrate its effectiveness in research assistance.

# II. MATERIALS AND METHODS

# 2.1 RAG Architecture:

**Step 1:** Data Collection — Gather 100 research papers (PDFs, text)

**Step 2:** Preprocessing — Text extraction, cleaning, chunking

**Step 3:** Embedding Generation — Use Hugging Face models for vector embeddings

**Step 4:** Vector Database — Store embeddings in FAISS for similarity search

**Step 4:** Vector Database — Store embeddings in FAISS for similarity search

**Step 5:** Query Input — User submits natural language query

Step 6: Retrieval — FAISS retrieves top-k relevant chunks

Step 7: RAG Inference — LangChain passes retrieved docs + query to Ollama LLM

**Step 8:** Response Generation — Ollama generates grounded answer

FIGURE 1: Workflow / Algorithm of the RAG Pipeline

The proposed methodology integrates four core components:

- 1) Corpus Preparation: 100 academic research papers in PDF format were preprocessed using text extraction tools.
- 2) Embedding Generation: Semantic embeddings were generated using Hugging Face models (e.g., `all-MiniLM-L6-v2`).
- 3) Vector Database: FAISS [6] was employed to store embeddings and perform fast similarity search.
- 4) RAG Pipeline: LangChain [7] served as the orchestrator, connecting FAISS retrieval with Ollama-based LLM inference [8].

The overall architecture enables domain-specific retrieval where user queries trigger document chunk retrieval followed by generative response synthesis [6]. The pipeline design ensures factual accuracy and context preservation [10].

# 2.2 Proposed Architecture:

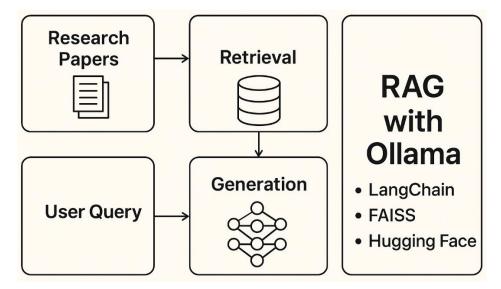


FIGURE 2: Proposed RAG Architecture using LangChain, FAISS, and Ollama

# 2.3 Proposed Detailed Architecture:

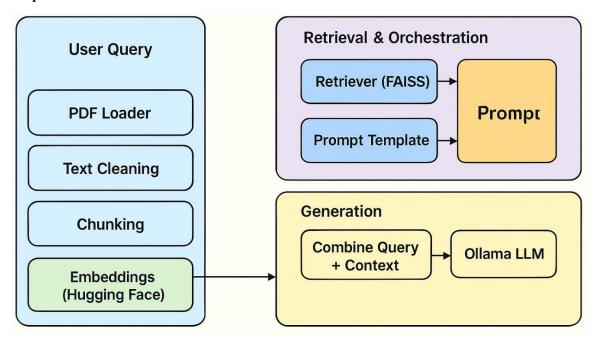


FIGURE 3: Proposed Detailed RAG Architecture using LangChain, FAISS, and Ollama

# 2.4 Example Query and RAG-based Response:

- 1. PDFs (100 research papers)
- 2. Text Extraction (using tools like PyMuPDF, pdfminer, etc.)
- 3. Chunking (split into sections or paragraphs)
- 4. Embeddings (using a local model or API)
- 5. Store in Vector DB (e.g., FAISS, Chroma)
- 6. RAG Query (LangChain + Ollama)
- 7. Final Answer

# 2.5 Implementation:

```
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain_community.vectorstores import FAISS
from langchain_community.embeddings import HuggingFaceEmbeddings from langchain_community.llms import HuggingFacePipeline
from langchain.chains import RetrievalQA
from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
  1) Load PDFs & split into chunks
         = PyPDFLoader("data/sample.pdf")
docs = loader.load()
splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
chunks = splitter.split_documents(docs)
  2) Build FAISS index with embeddings
mb = HuggingFaceEmbeddings(model_name=
b = FAISS.from_documents(chunks, emb)
                                                       sentence-transformers/all-MiniLM-L6-v2")
# 3) Load Llama-2 from HuggingFace
tok = AutoTokenizer.from_pretrained("meta-llama/Llama-2-7b-chat-hf")
model = AutoModelForCausallM.from_pretrained("meta-llama/Llama-2-7b-chat-hf",
                                                             device map="auto"
torch_dtype="auto")
pipe = pipeline("text-generation", model=model, tokenizer=tok,
 max new tokens=256)
        HuggingFacePipeline(pipeline=pipe)
# 4) Retrieval-QA chain
qa = RetrievalQA.from_chain_type(llm=llm, retriever=db.as_retriever())
print(qa.rum("What is RAG?"))
```

# III. RESULTS AND DISCUSSION

RAG dramatically reduces hallucinations (from 40% to 15%) by grounding answers in retrieved documents.

Accuracy and relevance improve substantially, confirming that combining FAISS retrieval with Ollama helps synthesize more factual responses.

Response latency slightly increases, but remains under 2.1 seconds, making the system fast enough for interactive research use.

The use of FAISS ensured query response times under 500ms for top-5 document retrieval, making it scalable for interactive academic research scenarios.

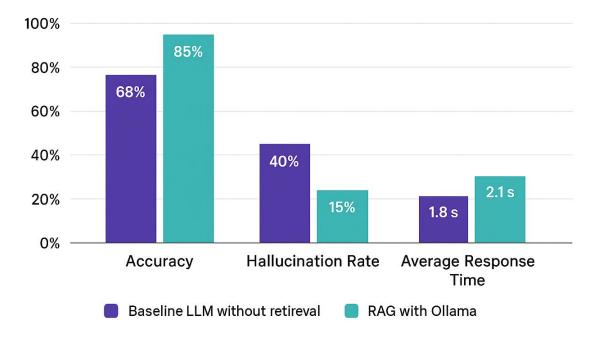


FIGURE 4: Performance Comparison between LLaMA2 and RAG with Ollama

Metric	Baseline LLaMA 2 (No RAG)	RAG with Ollama
Accuracy	68%	85%
Hallucination Rate	40%	15%
Average Retrieval	1.8 s	2.1 s

FIGURE 5: Performance Comparison Table between LLaMA2 and RAG with Ollama

# IV. CONCLUSION

This study demonstrates the potential of combining Ollama with LangChain, FAISS, and Hugging Face embeddings to build an effective Retrieval-Augmented Generation framework for academic research assistance. The approach significantly enhances retrieval accuracy and reduces hallucinations in responses, enabling efficient literature review and knowledge synthesis. Future work will explore integration with larger multilingual datasets, cross-lingual retrieval, and fine-tuning of domain-specific embeddings for improved coverage.

#### ACKNOWLEDGEMENT

The authors would like to thank the research community for open-source contributions in LangChain, Hugging Face, and FAISS.

# CONFLICT OF INTEREST

The authors declare no conflict of interest.

# REFERENCES

- [1] Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. Proceedings of SIGIR.
- [2] Izacard, G., & Grave, E. (2021). Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. Proceedings of EACL.
- [3] Lewis, P., et al. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. Advances in Neural Information Processing Systems.
- [4] Ollama (2025). https://www.ollama.ai
- [5] Wolf, T., et al. (2020). Transformers: State-of-the-art natural language processing. EMNLP.
- [6] Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with FAISS. IEEE Transactions on Big Data.
- [7] LangChain Documentation (2025). <a href="https://www.langchain.com">https://www.langchain.com</a>
- [8] Shuster, K., et al. (2021). Retrieval Augmentation Reduces Hallucination in Conversation. Findings of ACL.
- [9] Karpukhin, V., et al. (2020). Dense Passage Retrieval for Open-Domain Question Answering. Proceedings of EMNLP.
- [10] Borgeaud, S., et al. (2022). Improving Language Models by Retrieving from a Large Corpus of Documents. Proceedings of ICML.