# Implementation of an Agentic RAG Framework for Comparative Analysis of Classical Islamic Corpora

Hossein Shafaei

Independent Researcher & AI Systems Developer

Mashhad, Iran

Email: Shafaeihossein0@gmail.com

GitHub: https://github.com/shafaeihossin-hub/Classical-Shia-AI-RAG

*Abstract*—This research introduces a localized Retrieval-Augmented Generation (RAG) system specialized for theological discourse. Using an Agentic Query Expansion mechanism with DeepSeek-R1, we developed a system capable of cross-referencing the Holy Quran, Al-Kafi, Nahj al-Balagha, and Musnad Ahmad ibn Hanbal. The architecture ensures data privacy via Docker and Qdrant.

## I. INTRODUCTION

Navigating classical theological texts requires high precision. Traditional LLMs often suffer from hallucination. This project implements a private RAG pipeline to ensure strict adherence to verified local contexts.

## II. SYSTEM ARCHITECTURE

### A. Infrastructure and Storage

The vector storage is handled by Qdrant, running within a Docker container to ensure environment isolation.
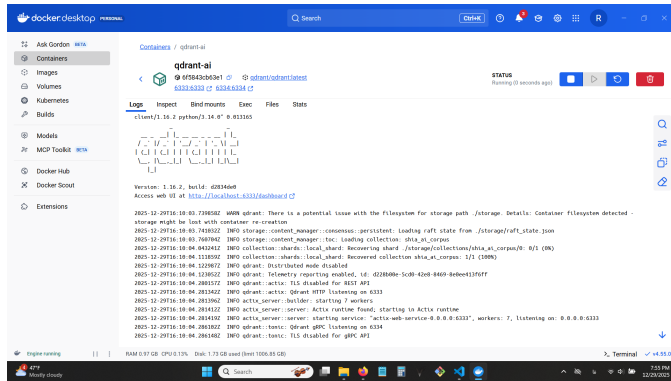


Fig. 1. Containerized Qdrant-AI environment in Docker.

### B. Agentic Search Logic

The system uses an "Agentic" approach, generating additional search keywords to improve retrieval hit rate before querying the database.

## III. IMPLEMENTATION DETAILS

Listing 1 demonstrates the core logic for query expansion and the RAG loop.

```python
def get_smart_queries(user_question):
    prompt = f"Generate 2 short search keywords in Persian related to: {
        user_question}."
    res = ollama.generate(model=MODEL_NAME, prompt=prompt)
    keywords = res['response'].strip().split(',')
    return [user_question] + [k.strip() for k in keywords]

def shia_ai_rag_query(user_question: str):
    queries = get_smart_queries(user_question)
    combined_context = ""
    for q in queries:
        combined_context += query_database(q) + "\n---\n"
    return generate_response(full_prompt)
```
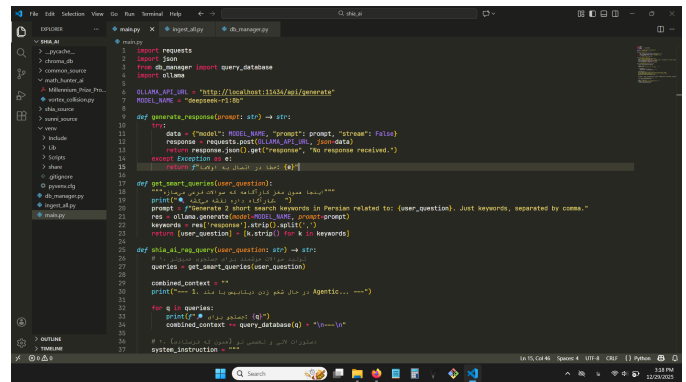
Listing 1. Agentic Query Expansion Logic



Fig. 2. Implementation of the RAG pipeline in VS Code.

## IV. PROJECT MANAGEMENT

The project is open-sourced under the MIT License and hosted on GitHub.

## V. CONCLUSION

This study validates that Small Language Models (SLMs) can achieve professional-grade performance in specialized fields when paired with intelligent retrieval mechanisms.
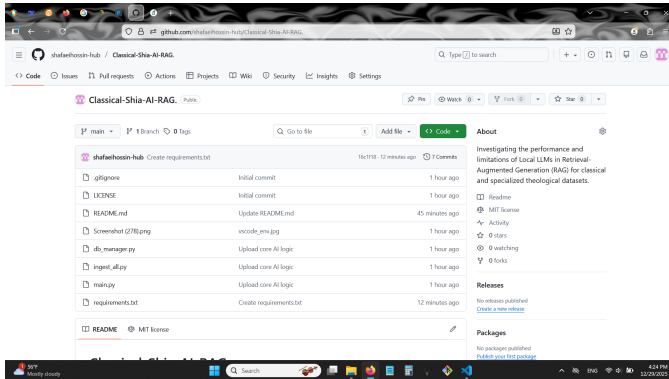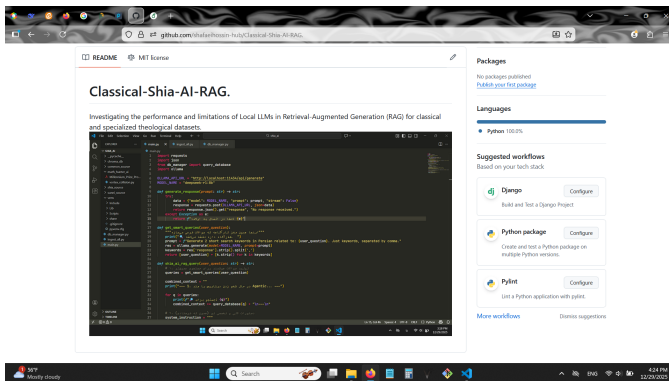
Fig. 3. Project structure and version control on GitHub.



Fig. 4. Project Documentation and README details.