

Government Policy Analysis using Retrieval Augmented Generation (RAG)

A Retrieval-Augmented Pipeline for Contextual Summarization and Policy Evaluation

Paavni Singh

1. Problem Statement

- Government policy documents are often **lengthy**, **dense**, and **difficult to analyze** at scale.
- Manual analysis of effectiveness, sentiment, and impact is **time-consuming**.
- There is a need for **automated, explainable systems** that can:
 - Summarize policies
 - Extract relevant information
 - Avoid dependency on external APIs (OpenAI, etc.)

2. Objective

To build an **AI-powered document analysis system** that can:

- Summarize policy PDFs (like *LaQshay*) using tree-based techniques.
- Retrieve specific policy content using semantic search.
- Route queries intelligently between summarization and retrieval engines.
- Run **entirely offline** using **local LLMs (LLaMA2)**.

3. Tools and Technologies

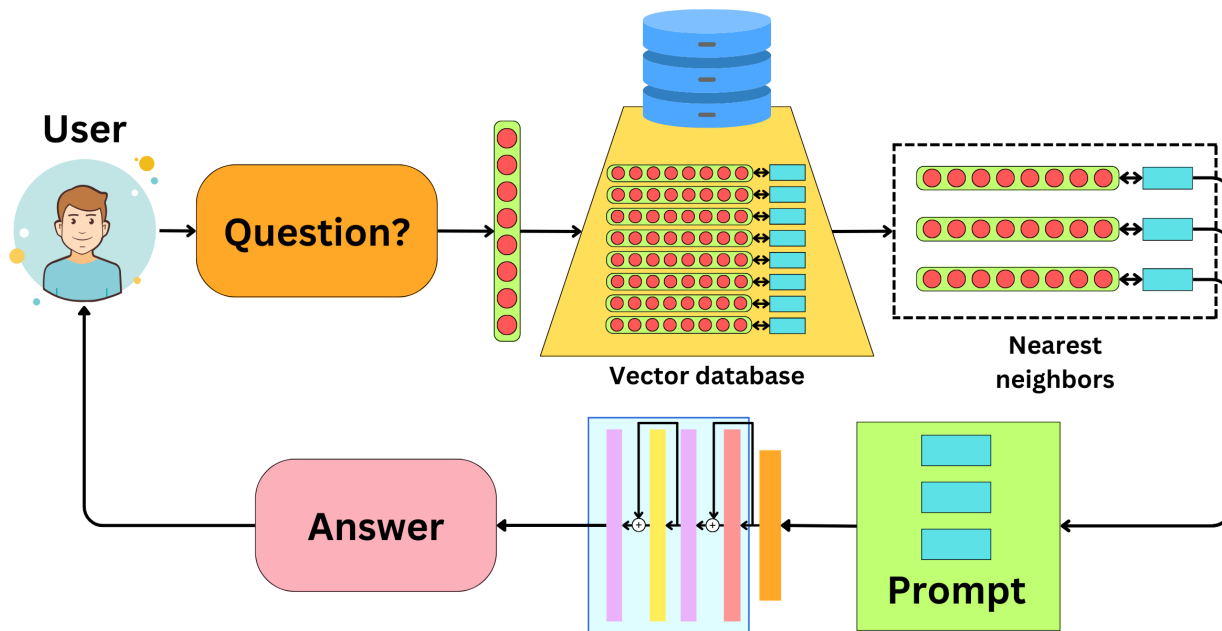
Component	Tool/Framework
Document Parsing	SimpleDirectoryReader, LlamaParse
Vector Indexing	LlamaIndex
LLM	LLaMA2 via Ollama
Summarization	SummaryIndex + tree_summarize
Retrieval	VectorStoreIndex
Routing Logic	RouterQueryEngine + LLMSingleSelector
Programming Language	Python

4. Architecture Diagram

A simple diagram showing the pipeline:

PDF → Sentence Splitter → Index → Router → Response

- Summarization Route → SummaryIndex
- Retrieval Route → VectorStoreIndex



5. Implementation

i) Document Ingestion:

```
14 documents = SimpleDirectoryReader(input_files=["/Users/paavnisingh/govt-pol/RKSK.pdf"]).load_data()
15
```

ii) Sentence Splitting + Indexing + LLM Setup:

```
15
16 splitter = SentenceSplitter(chunk_size=1024)
17 nodes = splitter.get_nodes_from_documents(documents)
18 llm = Ollama(model="llama2")
19
```

iii) Create Summary & Vector Indices:

```
20 summary_index = SummaryIndex(nodes)
21 vector_index = VectorStoreIndex(nodes)
22
```

iv) Router Engine with Query Tools:

6. Results & Observations

- Achieved **accurate contextual summaries** without internet dependency.

```

23 summary_query_engine = summary_index.as_query_engine(
24     response_mode="tree_summarize",
25     use_async=True,
26 )
27 vector_query_engine = vector_index.as_query_engine()
28 from llama_index.core.tools import QueryEngineTool
29
30
31 summary_tool = QueryEngineTool.from_defaults(
32     query_engine=summary_query_engine,
33     description=(
34         """Useful for summarization questions related to RKSK policy from multiple papers
35         Also useful for making a list of all the analysis done based on the query."""
36     ),
37 )
38
39 vector_tool = QueryEngineTool.from_defaults(
40     query_engine=vector_query_engine,
41     description=(
42         "Useful for retrieving specific context from the RKSK policy from multiple papers"
43     ),
44 )
45
46
47
48 query_engine = RouterQueryEngine(
49     selector=LLMSingleSelector.from_defaults() ,
50     query_engine_tools=[
51         summary_tool,
52         vector_tool,
53     ],
54     verbose=True
55 )

```

- The *RouterQueryEngine* handled intent switching intelligently.
- Local inference using *Ollama* kept everything privacy-focused.
- *SummaryIndex* provided concise outputs, while *VectorStoreIndex* enabled deep semantic lookups.

7. Limitations

- Local LLMs (like llama2) can be **slower** and **less accurate** than GPT-4.
- Summarization may fail on **extremely large PDFs** (>100 pages) unless chunked well.
- No sentiment scoring model was used; can be added with fine-tuned classifiers.

8. Future Work

- Integrate a **feedback loop** to improve summaries over time.
- Use **RAG with fine-tuned local models** for even deeper policy analysis.
- Add **user submission analysis** module (compare user ideas to past policies).

9. Conclusion

- Successfully built a **retrieval + summarization pipeline** using **LlamaIndex** and **LLaMA2**.
- Enables **offline policy document analysis** with explainable output.
- Can be scaled for **government think tanks**, **academia**, and **citizen engagement portals**.