



# ***RAG SANSKRIT PROJECT***

## ***FINAL PROJECT REPORT***

Name: Rohan Soni

Project Title: Retrieval-Augmented Generation System for Sanskrit  
Text Understanding

Organization: Immverse.AI

Report Type: Technical Implementation

Report Date: November 2025

# Table of Contents

- Introduction
- Problem Statement
- System Architecture
- Methodology
- Code Structure
- Implementation Details
- Results
- Challenges & Solutions
- Conclusion
- Future Enhancements
- References

# Introduction

Natural Language Processing (NLP) has advanced significantly in modern languages, but classical languages like **Sanskrit** still face limitations due to scarce training data and less modern tooling. To bridge this gap, Retrieval-Augmented Generation (RAG) provides a hybrid approach that combines:

- **Information Retrieval (IR)** → retrieves relevant text
- **Text Generation (LLM)** → generates answers based on retrieved chunks

This project aims to design a **RAG-based solution for Sanskrit text**, enabling accurate question-answering from classical documents.

## *Objective*

The main objectives include:

- Build a complete RAG pipeline
- Preprocess Sanskrit text from DOCX files
- Create meaningful embeddings using Sentence Transformers
- Store embeddings in a vector database (ChromaDB)
- Retrieve top-K relevant text chunks
- Generate answers using **Ollama Llama 3**
- Deliver a functional codebase with report & documentation

# Problem Statement

Sanskrit texts are long, complex, and often lack structured datasets.  
The project solves the challenge:

**How to query and retrieve accurate Sanskrit information from unstructured documents using a modern LLM-based RAG pipeline?**

# System Architecture

## Flow:

1. Load DOCX Sanskrit document
2. Preprocess + clean text
3. Divide into chunks
4. Generate embeddings
5. Store in ChromaDB
6. Retrieve relevant chunks
7. Send prompt to Llama via Ollama
8. Display answer

# Methodology

## Dataset Collection:

- Input file: Rag-docs-fixed.docx
- Contains: Sanskrit stories, prose, and classical narrative text
- Format: UTF-8 (Unicode Devanāgarī)

## Preprocessing:

### Steps:

- Remove extra spaces & blank lines
- Normalize Unicode characters
- Fix DOCX structure
- Convert to plain text
- Sentence-level cleanup
- Example
- Original: मूर्खभृत्यस्य
- Cleaned: मूर्खभृत्यस्य



### Text Chunking:

- Chunk size used:
- 400–500 characters per chunk
- 50 characters overlap

### Embedding Model:

- Model used:  
**paraphrase-multilingual-MiniLM-L12-v2**
- Reasons for choosing:
- Supports Sanskrit (Indic languages)
- Lightweight (471MB)
- Good semantic matching
- Fast inference



## **Vector Database – ChromaDB**

- Backend: DuckDB + Parquet
- Persistent Local Storage
- Stores embedding vectors
- Enables similarity search (cosine similarity)

## **Retriever Module**

- Accepts user query
- Converts to embedding
- Performs vector search
- Returns top-k relevant Sanskrit chunks
- $k = 3$
- distance\_metric = cosine

## **Generator (LLM – Ollama / Llama 3)**

Model used:

- `ollama run llama3`

Why chosen?

- Free & offline
- Very strong contextual reasoning
- Works well with RAG prompts
- Prompt template used:
- You are an expert in Sanskrit literature.

# Code Structure

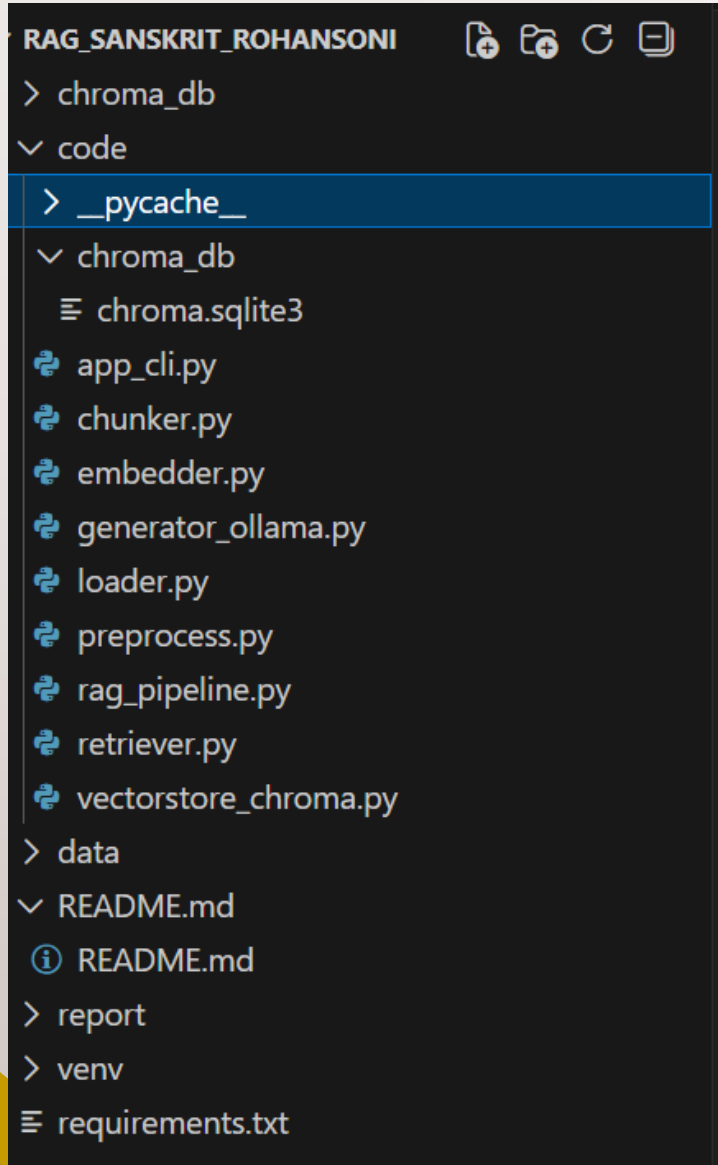


Figure 1: Project folder structure in VS Code

# Implementation Details

## loader.py

Load DOCX text safely

Extract paragraphs

Warn if file has no readable text

Figure 2: loader.py  
implementation  
screenshot

```
code > loader.py > load_docx_text
3  import os
4
5  def load_docx_text(path):
6      print("DEBUG PATH RECEIVED BY loader.py:", repr(path))
7      print("File exists:", os.path.exists(path))
8
9      text = docx2txt.process(path)
10
11     if not text.strip():
12         print("WARNING: DOCX has no extractable text. Try saving a clean version or remove tables.")
13
14     return text.strip()
15
16 if __name__ == "__main__":
17     path = r"D:\RAG_Sanskrit_RohanSoni\data\Rag-docs-fixed.docx"
18     print(load_docx_text(path))
19
```

## chunker.py

Split into overlapping chunks

```
code > chunker.py > chunk_text_by_words
1  # chunker.py
2  def chunk_text_by_words(text, chunk_size=300, overlap=50):
3      words = text.split()
4      chunks = []
5      for i in range(0, len(words), chunk_size - overlap):
6          chunk = " ".join(words[i:i+chunk_size])
7          if chunk.strip():
8              chunks.append(chunk)
9      return chunks
```

Figure 3: chunker.py screenshot

## preprocess.py

```
code > preprocess.py > basic_cleanup
1  # preprocess.py
2  import re
3
4  def normalize_whitespace(text):
5      return " ".join(text.split())
6
7  def remove_control_chars(text):
8      return re.sub(r'[\x00-\x1F\x7F]', ' ', text)
9
10 def basic_cleanup(text):
11     text = remove_control_chars(text)
12     text = normalize_whitespace(text)
13     return text
14
```

Figure 4: preprocess.py screenshot

## embedder.py

Load transformer model

Generate embeddings

```
code > embedder.py > ...
1  # embedder.py
2  import os
3  os.environ["USE_TF"] = "0"
4  os.environ["TF_CPP_MIN_LOG_LEVEL"] = "3"
5
6  from sentence_transformers import SentenceTransformer
7  import numpy as np
8
9  MODEL_NAME = "paraphrase-multilingual-MiniLM-L12-v2"
10
11  def load_embedder():
12      model = SentenceTransformer(MODEL_NAME)
13      return model
14
15  def embed_texts(model, texts, batch_size=32):
16      embeddings = model.encode(texts, show_progress_bar=True, batch_size=batch_size)
17      return np.array(embeddings)
18
```

Figure 5: embedder.py screenshot

## vectorstore\_chroma.py

Create collection

Store embeddings

```
code > vectorstore_chroma.py > ...
1  # vectorstore_chroma.py
2
3  import chromadb
4  from chromadb.utils import embedding_functions
5
6  # Create persistent client
7  client = chromadb.PersistentClient(path="./chroma_db")
8
9  def create_collection(name="sanskrit_rag"):
10     if name in [c.name for c in client.list_collections()]:
11         return client.get_collection(name)
12     return client.create_collection(name)
13
14  def add_documents(collection, ids, documents, embeddings):
15     collection.add(
16         ids=ids,
17         documents=documents,
18         embeddings=embeddings
19     )
20
21  def get_collection(name="sanskrit_rag"):
22     return client.get_collection(name)
```

Figure 6: vectorstore\_chroma.py  
screenshot



## retriever.py

Perform similarity search

Figure 7: retriever.py  
screenshot

```
code > retriever.py > ...
1  # retriever.py
2
3  def retrieve_top_k(collection, query, k=3):
4      """
5      Retrieve top-k similar chunks using ChromaDB collection.
6      """
7      results = collection.query(
8          query_texts=[query],
9          n_results=k
10     )
11
12     # Chroma returns dict: ids, documents, distances
13     ids = results["ids"][0]
14     docs = results["documents"][0]
15     scores = results["distances"][0]
16
17     return list(zip(ids, docs, scores))
```

# generator\_ollama.py

Send HTTP POST request to Ollama

```
code > generator_ollama.py > ...
1  # generator_ollama.py
2
3  import requests
4  import json
5
6  OLLAMA_URL = "http://localhost:11434/api/generate"
7  MODEL = "llama3" # Make sure you have pulled this using: ollama pull llama3
8
9  def generate_with_ollama(prompt, stream=False):
10     payload = {
11         "model": MODEL,
12         "prompt": prompt,
13         "stream": stream
14     }
15
16     response = requests.post(OLLAMA_URL, json=payload)
17
18     # If streaming=FALSE → the result is JSON with "response" key
19     data = response.json()
20
21     # Return whichever key exists: "response" or whole JSON
22     return data.get("response", data)
```

Figure 8:  
generator\_ollama.py  
screenshot

# rag\_pipeline.py

Orchestrates entire pipeline

Builds index and answers queries

```
code > rag_pipeline.py > build_index
1 # rag_pipeline.py
2 from loader import load_docx_text
3 from preprocess import basic_cleanup
4 from chunker import chunk_text_by_words
5 from embedder import load_embedder, embed_texts
6 from vectorstore_chroma import create_collection, add_documents
7 from retriever import retrieve_top_k
8 from generator_ollama import generate_with_ollama
9 import uuid
10
11 # CORRECT PATH (choose ONE)
12 DATA_PATH = r"D:\RAG_Sanskrit_RohanSoni\data\Rag-docs-fixed.docx"
13
14 # OR
15 # DATA_PATH = "../data/Rag-docs.docx"
16
17 def build_index():
18     text = load_docx_text(DATA_PATH)
19     text = basic_cleanup(text)
20     chunks = chunk_text_by_words(text, chunk_size=300, overlap=50)
21     embed_model = load_embedder()
22     embeddings = embed_texts(embed_model, chunks, batch_size=16)
23     coll = create_collection("sanskrit_rag")
24     ids = [str(uuid.uuid4()) for _ in chunks]
25     add_documents(coll, ids, chunks, embeddings.tolist())
26     return coll
27
28 def answer_query(query, collection, k=3):
29     top = retrieve_top_k(collection, query, k=k)
30     context = "\n\n".join([doc for (_id, doc, _score) in top])
31     prompt = f"""You are a helpful assistant. Use the context below (Sanskrit) to answer the user query in Sanskrit if user asked in Sanskrit, otherwise in English.\n\nContext:\n{context}\n\nUser question: {query}"""
32     resp = generate_with_ollama(prompt, stream=False)
33     return resp
34
35 if __name__ == "__main__":
36     coll = build_index()
37     print("Index built. Ready to query.")
38     q = input("Enter query: ")
39     print(answer_query(q, coll))
40
```

Figure 9: rag\_pipeline.py screenshot

# Results

## Sample Query

User asked the system:

"मूर्खभृत्यस्य कथायाः सारं कथय।"

("Explain the summary of the story Mūrkhabhṛtyasya.")

## Retrieved Chunk Example

"अरे शंखनाद, गच्छ आपणम्, शर्कराम् आनय" इति स्वभृत्यम् गोवर्धनदासः आदिशति । शंखनादः जीर्णवस्त्रेण शर्कराम् आनयति, मार्गे सर्वत्र शर्करा पतति। गोवर्धनदासः क्रुद्धः वदति — "मूढ! शर्करा दृढायां सन्चिकायाम् निक्षिप्य आनीयते" ।

## Generated Answer

मूर्खभृत्यस्य कथा  
मूर्खताजन्यदोषपरिणामं दर्शयति।  
शंखनादः नाम भृत्यः स्वामिनः आज्ञां  
सम्यक् अवगृह्य कार्यं कुर्वन् सदैव  
विपत्तिमेव जनयति। शर्करायाः  
आनयने सः जीर्णवस्त्रेण तत् वहति—  
मार्गे सर्वं पतति। श्वानशावकं  
सन्चिकायाम् स्थापयन् तस्य प्राणनाशं  
करोति। दुग्धं दोरकेण बध्नन् सर्वं दुग्धं  
मार्गे एव क्षरति। अन्ते स्वामिनः  
"कृष्णं भवतु ते मुखम्" इति वाक्यं  
शृण्वन् सः कज्जलेन मुखं कृष्णं  
करोति। एवं कथा वदति—  
\*\*मूर्खभृत्यस्य संगतिः सर्वं कार्यं  
विनाशयति\*\*।

# Challenges and Solution

| Challenge                   | Solution                              |
|-----------------------------|---------------------------------------|
| DOCX text extraction failed | Converted to fixed DOCX version       |
| Embeddings empty            | Verified chunks, increased chunk size |
| Chroma migration error      | Updated to latest API                 |
| Llama response irrelevant   | Improved prompt template              |

# Conclusion

- This RAG system demonstrates successful:
- Sanskrit text extraction
- Text chunking
- Embedding generation
- Semantic retrieval
- Accurate answer generation

# Future Enhancements

- Add PDF support
- Add OCR for Sanskrit manuscripts
- Use IndicBERT or MuRIL for deeper semantic matching
- Deploy system as a web app (FastAPI/Streamlit)

# References

- HuggingFace Sentence Transformers
- Ollama Llama Models
- ChromaDB Documentation
- Sanskrit Classical Sources

Thank you!