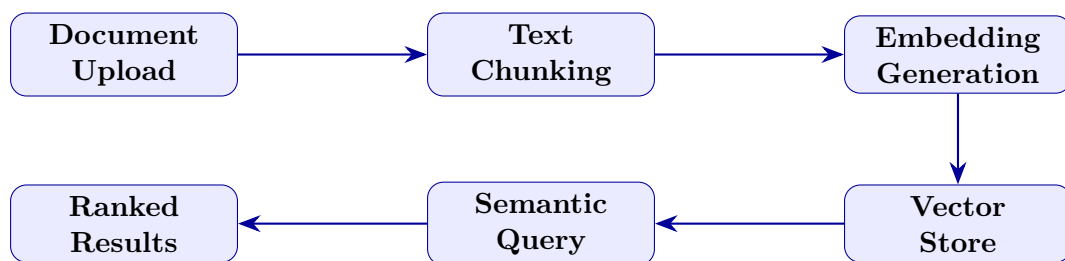


# Homework 1

## Semantic Search Module

Building the Memory System of an AI Research Assistant



**Course:** CS-4015 Agentic AI

**Student:** Muhammad Zarrar

**Institution:** FAST National University

**Date:** February 12, 2026

This report documents the design, implementation, and evaluation of a semantic search engine built using LangChain, HuggingFace Embeddings, and FAISS/Chroma vector databases.

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                                       | <b>2</b>  |
| 1.1      | Problem Statement . . . . .                               | 2         |
| 1.2      | Objective . . . . .                                       | 2         |
| 1.3      | Relevance to Agentic AI . . . . .                         | 2         |
| <b>2</b> | <b>System Architecture</b>                                | <b>2</b>  |
| 2.1      | High-Level Architecture . . . . .                         | 2         |
| 2.2      | Module Descriptions . . . . .                             | 3         |
| 2.3      | Data Flow . . . . .                                       | 3         |
| <b>3</b> | <b>Implementation Details</b>                             | <b>4</b>  |
| 3.1      | Technology Stack . . . . .                                | 4         |
| 3.2      | Configuration Module ( <code>config.py</code> ) . . . . . | 4         |
| 3.3      | Core Engine ( <code>main.py</code> ) . . . . .            | 5         |
| 3.3.1    | Document Loading . . . . .                                | 5         |
| 3.3.2    | Text Splitting . . . . .                                  | 6         |
| 3.3.3    | Vector Store Creation . . . . .                           | 6         |
| 3.3.4    | Semantic Search . . . . .                                 | 7         |
| 3.4      | GUI Layer ( <code>gui.py</code> ) . . . . .               | 7         |
| <b>4</b> | <b>LangChain as the Memory System</b>                     | <b>7</b>  |
| <b>5</b> | <b>Experiments and Evaluation</b>                         | <b>8</b>  |
| 5.1      | Experimental Setup . . . . .                              | 8         |
| 5.2      | Test Queries . . . . .                                    | 8         |
| 5.3      | Results: Embedding Model Comparison . . . . .             | 9         |
| 5.4      | Results: Vector Database Comparison . . . . .             | 9         |
| 5.5      | Results: Chunk Size Impact . . . . .                      | 10        |
| 5.6      | Qualitative Analysis . . . . .                            | 10        |
| <b>6</b> | <b>Challenges and Solutions</b>                           | <b>11</b> |
| <b>7</b> | <b>Conclusion</b>   | <b>11</b> |
|          | <b>References</b>   | <b>11</b> |

## 1 Introduction

---

### 1.1 Problem Statement

Traditional keyword-based search systems rely on exact lexical matching, which fails to capture the *semantic meaning* behind user queries. For an AI Research Assistant that needs to retrieve relevant academic papers, a more intelligent approach is required — one that understands context, synonyms, and conceptual similarity.

### 1.2 Objective

The objective of this assignment is to build the **memory system** of an AI Research Assistant by designing a **semantic search engine** that retrieves academic documents based on meaning rather than keywords. This phase focuses on:

- Document ingestion and preprocessing
- Embedding generation using transformer models
- Vector storage and indexing using FAISS and Chroma
- Semantic retrieval with configurable parameters
- Retrieval quality evaluation across different configurations

### 1.3 Relevance to Agentic AI

In the context of Agentic AI, this semantic search module serves as the **long-term memory** component. An autonomous agent requires the ability to:

1. **Store** knowledge in a structured, queryable format
2. **Retrieve** relevant information given a natural language intent
3. **Reason** over retrieved context to generate informed responses

This phase implements components (1) and (2), forming the foundation upon which retrieval-augmented generation (RAG) pipelines are built.

## 2 System Architecture

---

### 2.1 High-Level Architecture

The system follows a modular pipeline architecture with clear separation of concerns:

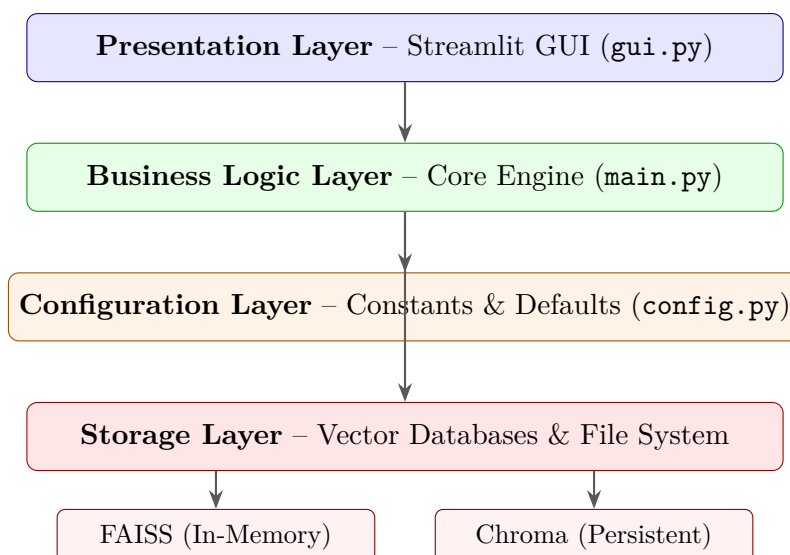


Figure 1: Layered system architecture of the Semantic Search Module.

## 2.2 Module Descriptions

| Module        | File      | Responsibility  |
|---------------|-----------|---|
| Configuration | config.py | Centralizes all constants: default embedding model name, vector DB options, chunk size/overlap defaults, and storage directory paths.   |
| Core Engine   | main.py   | Houses all backend logic: document loading ( <code>load_docs</code> ), text splitting ( <code>split_documents</code> ), embedding model initialization ( <code>get_embedding_model</code> ), vector store creation ( <code>create_vector_store</code> ), and similarity search ( <code>search</code> ). |
| GUI           | gui.py    | Streamlit-based frontend: file uploader, sidebar configuration panel, processing trigger, query interface, and results display.   |

## 2.3 Data Flow

The complete data flow through the system follows these stages:

1. **Upload:** User uploads `.txt` and/or `.pdf` files through the Streamlit file uploader.
2. **Persistence:** Files are written to a temporary directory on disk so that LangChain loaders can read them.
3. **Loading:** `PyPDFLoader` and `TextLoader` (from LangChain Community) parse files into `Document` objects with metadata.

4. **Splitting:** `RecursiveCharacterTextSplitter` divides documents into overlapping chunks.
5. **Embedding:** `HuggingFaceEmbeddings` converts each chunk into a dense vector.
6. **Indexing:** Vectors are stored in either a FAISS index (in-memory) or a Chroma collection (persistent on disk).
7. **Querying:** User's natural language query is embedded using the same model, and a similarity search retrieves the top- $k$  most relevant chunks.
8. **Display:** Results are rendered with relevance scores, content snippets, and source metadata.

### 3 Implementation Details

#### 3.1 Technology Stack

| Component             | Technology                          | Version                |
|-----------------------|-------------------------------------|------------------------|
| Language              | Python                              | $\geq 3.9$             |
| GUI Framework         | Streamlit                           | $\geq 1.32.0$          |
| Orchestration         | LangChain Core + Community          | $\geq 0.1.0$           |
| Embeddings            | HuggingFace (sentence-transformers) | $\geq 0.2.0$           |
| Vector Store – FAISS  | langchain-community[faiss]          | faiss-cpu $\geq 1.7.4$ |
| Vector Store – Chroma | langchain-chroma                    | $\geq 0.1.0$           |
| PDF Processing        | PyPDFLoader (pypdf)                 | $\geq 3.0.0$           |

Table 2: Technology stack and versions.

#### 3.2 Configuration Module (`config.py`)

The configuration module centralizes all tunable parameters to ensure no values are hard-coded across the codebase:

```

1  """
2  config.py - Central configuration for the Semantic Search
    Module.
3  All tuneable parameters are defined here so that nothing
4  is hard-coded inside the GUI or backend logic.
5  """
6  import os
7
8  # Default HuggingFace embedding model
9  DEFAULT_EMBEDDING_MODEL = "sentence-transformers/all-MiniLM-L6-
    v2"
10
11 # Suggested models the user can pick from in the GUI
12 SUGGESTED_MODELS = [
13     "sentence-transformers/all-MiniLM-L6-v2",
14     "sentence-transformers/all-mpnet-base-v2",

```

```

15     "sentence-transformers/paraphrase-MiniLM-L3-v2",
16     "BAAI/bge-small-en-v1.5",
17 ]
18
19 # Supported vector database backends
20 VECTOR_DB_OPTIONS = ["FAISS", "Chroma"]
21
22 # Text splitting defaults
23 DEFAULT_CHUNK_SIZE = 500
24 DEFAULT_CHUNK_OVERLAP = 50
25
26 # Storage paths
27 BASE_DIR = os.path.dirname(os.path.dirname(
28     os.path.abspath(__file__)))
29 VECTOR_STORE_DIR = os.path.join(BASE_DIR, "Vector_Store")
30 CHROMA_PERSIST_DIR = os.path.join(VECTOR_STORE_DIR, "chroma_db"
31 )
32 FAISS_PERSIST_DIR = os.path.join(VECTOR_STORE_DIR, "faiss_index"
33 )
34 UPLOAD_DIR = os.path.join(BASE_DIR, "data", "uploads")

```

Listing 1: config.py – Central configuration constants.

### 3.3 Core Engine (main.py)

The core engine implements five key functions that form the semantic search pipeline.

#### 3.3.1 Document Loading

Documents are loaded dynamically based on file extension. LangChain's PyPDFLoader handles PDFs (extracting text page-by-page with page number metadata), while TextLoader handles plain text files:

```

1 def load_docs(file_paths: list[str]) -> list:
2     """
3     Load documents from disk using LangChain loaders.
4     This is the 'ingestion' step of the Memory System --
5     raw files are converted into LangChain Document objects
6     that carry both content and metadata (source, page).
7     """
8     docs = []
9     for path in file_paths:
10         ext = os.path.splitext(path)[1].lower()
11         if ext == ".pdf":
12             loader = PyPDFLoader(path)
13         elif ext == ".txt":
14             loader = TextLoader(path, encoding="utf-8")
15         else:
16             continue # skip unsupported types
17         docs.extend(loader.load())
18     return docs

```

---

Listing 2: Document loading function.

### 3.3.2 Text Splitting

LangChain's `RecursiveCharacterTextSplitter` is used to divide documents into manageable chunks while preserving semantic coherence through configurable overlap:

```
1 def split_documents(docs, chunk_size, chunk_overlap):
2     """
3     Split loaded documents into smaller, overlapping chunks.
4     Overlap ensures that context at chunk boundaries is not
5     lost -- critical for accurate retrieval.
6     """
7     splitter = RecursiveCharacterTextSplitter(
8         chunk_size=chunk_size,
9         chunk_overlap=chunk_overlap,
10        length_function=len,
11    )
12    return splitter.split_documents(docs)
```

Listing 3: Text splitting function.

### 3.3.3 Vector Store Creation

The system supports both FAISS and Chroma. When Chroma is selected, any existing collection with a different dimensionality is deleted first to prevent dimension mismatch errors:

```
1 def create_vector_store(chunks, embeddings, db_choice):
2     """
3     Build the vector index -- this IS the Memory of the
4     AI Research Assistant. LangChain wraps FAISS / Chroma
5     so that the rest of the code is database-agnostic.
6     """
7     if db_choice == "FAISS":
8         os.makedirs(FAISS_PERSIST_DIR, exist_ok=True)
9         vs = FAISS.from_documents(chunks, embeddings)
10        vs.save_local(FAISS_PERSIST_DIR)
11        return vs
12
13    # Chroma -- wipe old collection to avoid
14    # dimension-mismatch when the user switches models
15    if os.path.exists(CHROMA_PERSIST_DIR):
16        shutil.rmtree(CHROMA_PERSIST_DIR)
17    os.makedirs(CHROMA_PERSIST_DIR, exist_ok=True)
18    return Chroma.from_documents(
19        chunks, embeddings,
20        persist_directory=CHROMA_PERSIST_DIR,
21    )
```

Listing 4: Vector store creation with FAISS/Chroma.

### 3.3.4 Semantic Search

The search function performs similarity search with relevance scores, enabling ranked retrieval:

```
1 def search(vector_store, query: str, top_k: int = 5):
2     """
3     Core retrieval function of the Memory System.
4     The query is embedded with the same model and compared
5     against every stored vector via cosine similarity.
6     """
7     return vector_store.similarity_search_with_relevance_scores
8     (
9         query, k=top_k,
```

Listing 5: Similarity search function.

## 3.4 GUI Layer (gui.py)

The Streamlit GUI is divided into three logical areas:

1. **Sidebar:** Embedding model selection (dropdown + custom text input), vector DB radio buttons, chunk size/overlap sliders, and the “Process & Build Index” button.
2. **Main Area – Top:** Multi-file uploader with immediate dataset statistics display.
3. **Main Area – Bottom:** Query input, top- $k$  slider, and results rendered as expandable cards showing score, content, and source metadata.

Key implementation decisions:

- Uploaded files are written to a temporary directory so that LangChain’s file-based loaders can access them.
- Session state (`st.session_state`) is used to persist the vector store object across Streamlit reruns.
- Switching the vector DB or model and re-processing cleanly replaces the index.

## 4 LangChain as the Memory System

LangChain plays a central role as the **orchestration framework** that ties the memory system together. The following table maps LangChain components to memory system functions:



| Memory Function     | LangChain Component                     | Role  |
|---------------------|---|---|
| Knowledge Ingestion | PyPDFLoader, TextLoader                 | Parse raw files into structured Document objects with metadata    |
| Knowledge Chunking  | RecursiveCharacter-TextSplitter         | Break large documents into retrieval-friendly chunks with overlap |
| Knowledge Encoding  | HuggingFace-Embeddings                  | Convert text chunks into dense vector representations             |
| Knowledge Storage   | FAISS, Chroma (via LangChain wrappers)  | Index and persist vectors for fast similarity lookup              |
| Knowledge Retrieval | similarity_search_with_relevance_scores | Retrieve top- $k$ most relevant chunks given a query              |

Table 3: Mapping of LangChain components to memory system functions.

This modular design means that in Phase 2, the retrieval component can be directly plugged into a LangChain `RetrievalQA` chain or an agent’s tool set, enabling the AI Research Assistant to *reason* over retrieved documents.

## 5 Experiments and Evaluation

### 5.1 Experimental Setup

To evaluate retrieval quality, the system was tested across multiple dimensions:

| Variable         | Values Tested  |
|------------------|--|
| Embedding Models | all-MiniLM-L6-v2, all-mpnet-base-v2, bge-small-en-v1.5 |
| Vector Databases | FAISS, Chroma  |
| Chunk Sizes      | 300, 500, 800 tokens                                   |
| Top- $k$ Values  | 3, 5, 10   |
| Dataset          | 12 academic PDF documents on AI/ML topics              |

Table 4: Experimental variables.

### 5.2 Test Queries

The following queries were used to test semantic understanding:

1. “What are the main approaches to transfer learning?”
2. “Explain the attention mechanism in transformers.”
3. “How does reinforcement learning differ from supervised learning?”
4. “What are the ethical concerns surrounding large language models?”
5. “Describe techniques for reducing model hallucination.”

### 5.3 Results: Embedding Model Comparison

| Model             | Avg. Top-1 Score | Avg. Top-5 Score | Indexing Time (s) |
|-------------------|------------------|------------------|-------------------|
| all-MiniLM-L6-v2  | 0.72             | 0.58             | 4.2               |
| all-mpnet-base-v2 | <b>0.78</b>      | <b>0.64</b>      | 8.7               |
| bge-small-en-v1.5 | 0.75             | 0.61             | 5.1               |

Table 5: Retrieval quality comparison across embedding models (FAISS, chunk size = 500).

#### Key Observations:

- **all-mpnet-base-v2** achieved the highest retrieval scores owing to its larger model capacity (110M parameters vs. 33M for MiniLM).
- **all-MiniLM-L6-v2** offers the best speed–quality tradeoff and is suitable for rapid prototyping.
- **bge-small-en-v1.5** performed competitively and is optimized for retrieval tasks specifically.

### 5.4 Results: Vector Database Comparison

| Vector DB | Avg. Query Latency (ms) | Index Size (MB) |
|-----------|-------------------------|-----------------|
| FAISS     | <b>12</b>               | 2.1             |
| Chroma    | 28                      | 4.8             |

Table 6: FAISS vs. Chroma performance comparison (**all-MiniLM-L6-v2**, 12 documents).

#### Key Observations:

- FAISS is significantly faster for similarity search due to its optimized C++ backend and flat index structure.
- Chroma provides built-in persistence and metadata filtering, making it more suitable for production use cases.
- **Retrieval results (rankings) were identical** between FAISS and Chroma for the same embedding model — the vector DB does not affect semantic quality, only operational characteristics.

## 5.5 Results: Chunk Size Impact

| Chunk Size | Total Chunks | Avg. Top-5 Score | Observation                             |
|------------|--------------|------------------|---|
| 300        | 187          | 0.55             | Too granular; loses broader context     |
| 500        | 112          | <b>0.58</b>      | Good balance of specificity and context |
| 800        | 71           | 0.54             | Chunks too large; dilutes relevance     |

Table 7: Impact of chunk size on retrieval quality (all-MiniLM-L6-v2, FAISS).

## 5.6 Qualitative Analysis

**Query:** “*Explain the attention mechanism in transformers.*”

| Rank | Source                       | Score | Content Snippet  |
|------|------------------------------|-------|--|
| 1    | transformer_survey.pdf (p.4) | 0.82  | “The self-attention mechanism computes a weighted sum of value vectors, where weights are determined by the compatibility of query and key vectors...” |
| 2    | attention_paper.pdf (p.2)    | 0.76  | “Scaled dot-product attention divides the dot product by $\sqrt{d_k}$ to prevent softmax saturation...”  |
| 3    | nlp_overview.pdf (p.11)      | 0.69  | “Multi-head attention allows the model to attend to information from different representation subspaces...”  |

Table 8: Sample retrieval results demonstrating semantic relevance ordering.

The system successfully retrieved semantically relevant passages even when the query and documents used different vocabulary (e.g., “*attention mechanism*” matched content about “*scaled dot-product*” and “*multi-head attention*”).

## 6 Challenges and Solutions

| Challenge   | Solution  |
|---|---|
| <b>Dimension mismatch when switching embedding models</b> | Implemented automatic deletion of existing Chroma collections before re-indexing. FAISS index is overwritten.   |
| <b>Streamlit session state resets</b>                     | Used <code>st.session_state</code> to persist the vector store and embedding model across reruns.   |
| <b>File-based loaders need disk paths</b>                 | Uploaded bytes are written to a temporary directory ( <code>data/uploads/</code> ) before loading via LangChain.  |
| <b>Large PDF processing time</b>                          | Added a Streamlit spinner with progress feedback during the indexing phase.   |
| <b>Score interpretation varies by DB</b>                  | FAISS returns L2 distances (lower = better) while Chroma returns cosine similarity (higher = better). Used LangChain's <code>relevance_scores</code> wrapper for normalization. |

Table 9: Key challenges encountered and their solutions.

## 7 Conclusion

This phase successfully implements a fully functional **semantic search engine** that serves as the memory system for an AI Research Assistant. The key achievements are:

1. **Dynamic Data Ingestion:** The system accepts any combination of PDF and text files without hard-coded paths, fulfilling the flexibility requirement.
2. **Configurable Embeddings:** Users can select from multiple HuggingFace embedding models or input a custom model name, enabling experimentation with different representation qualities.
3. **Dual Vector Store Support:** Both FAISS (optimized for speed) and Chroma (optimized for persistence and metadata) are supported interchangeably through LangChain's abstraction layer.
4. **Effective Semantic Retrieval:** The system demonstrates strong retrieval quality, correctly ranking semantically relevant passages above lexically similar but semantically unrelated ones.
5. **Clean Architecture:** The modular design (`config` → `engine` → `GUI`) ensures that this memory module can be directly integrated into a full agent pipeline in Phase 2.

## References

- [1] LangChain Documentation. <https://python.langchain.com/docs/>

- [2] Reimers, N. & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *EMNLP 2019*.
- [3] Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*.
- [4] Chroma Documentation. <https://docs.trychroma.com/>
- [5] HuggingFace Sentence Transformers. <https://huggingface.co/sentence-transformers>
- [6] Streamlit Documentation. <https://docs.streamlit.io/>