

Book Recommender System Using LLMs and Vector Search

Semantic Search, LLM-Generated Descriptions, and Interactive
UI

Project Documentation

Author & Developer

Mohamed Amine Mammar El Hadj

LinkedIn: www.linkedin.com/in/mohamed-amine-mammar-el-hadj-715a41295

GitHub: <https://github.com/mimine47>

Email: mohamedamine.devtech@gmail.com

September 3, 2025

A comprehensive approach to building an intelligent book recommendation
system leveraging Large Language Models, Vector Embeddings, and
Modern UI Technologies

Contents

1	Author Information	2
1.1	Project Developer	2
1.2	Technical Expertise	2
2	Project Overview	2
2.1	Objective	2
2.2	End Goal	2
3	Motivation & Importance	3
4	Project Timeline & Workflow	3
5	Dataset Details	3
5.1	Source Files	3
5.2	Metadata Schema	4
6	Implementation Details	4
6.1	Book Description Generation	4
6.2	Metadata Merging	5
6.3	Text Chunking for Embeddings	5
6.4	Vector Embedding & Database	6
6.5	Thumbnail Fetching	7
6.6	Search Pipeline	7
6.7	Gradio Frontend Interface	8
7	Key Design Decisions	9
7.1	Model Selection	9
7.2	Technical Architecture	9
8	Performance Analysis	9
8.1	System Metrics	9
8.2	Quality Assessment	9
9	Future Improvements	10
9.1	Performance Enhancements	10
9.2	Feature Extensions	10
9.3	Scalability Considerations	10
10	Conclusion	10

1 Author Information

1.1 Project Developer

Mohamed Amine Mammar El Hadj

Deep Learning & Machine Learning Developer|Software Developer

- **LinkedIn:** www.linkedin.com/in/mohamed-amine-mammar-el-hadj-715a41295
- **GitHub:** <https://github.com/mimine47>
- **Email:** mohamedamine.devtech@gmail.com

1.2 Technical Expertise

This project demonstrates proficiency in:

- Large Language Models (LLMs) and Natural Language Processing
- Vector Databases and Semantic Search Technologies
- Machine Learning Pipeline Development
- Full-Stack Application Development
- API Integration and Data Processing
- Modern UI/UX Development with Gradio

2 Project Overview

2.1 Objective

This project aims to build a semantic book recommendation system leveraging:

- Large Language Models (Falcon-7B-Instruct) for **book description generation**
- Vector embeddings for **semantic similarity search**
- Chroma vector database for **persistent, efficient search**
- Gradio for an **interactive frontend** with thumbnails and metadata

2.2 End Goal

Enable users to enter natural language queries (e.g., *"books about magic"*) and receive highly relevant book recommendations with thumbnails and metadata.

3 Motivation & Importance

Traditional keyword-based book searches fail to capture semantic meaning, leading to suboptimal recommendations. This project addresses several key challenges:

- **Semantic Understanding:** Vector search ensures retrieval based on meaning, not just exact keywords
- **Enhanced Metadata:** LLMs generate engaging descriptions, enriching existing metadata
- **User Experience:** Interactive UI improves user engagement and usability
- **Scalability:** Vector databases provide efficient similarity search at scale

4 Project Timeline & Workflow

Table 1 presents the comprehensive project workflow with completion status.

Table 1: Project Timeline and Task Completion Status

Step	Task	Description	Status
1	Data Cleaning & Preparation	Consolidate books_cleaned.csv and description data	✓
2	Description Generation	Use Falcon-7B-Instruct to generate short descriptions	✓
3	Merge Metadata & Descriptions	Ensure all metadata aligns with generated descriptions	✓
4	Text Chunking	Split descriptions into chunks for embedding	✓
5	Vector Embedding	Encode chunks with all-MiniLM-L6-v2	✓
6	Vector DB Creation	Persist embeddings in Chroma	✓
7	Thumbnail Fetching	Use Google Books API to fetch book covers	✓
8	Search Pipeline Development	Combine vector search with metadata matching	✓
9	Frontend Interface	Build Gradio gallery UI for interactive queries	✓
10	Testing & Validation	Evaluate search results & refine pipeline	✓

5 Dataset Details

5.1 Source Files

- books_cleaned.csv – book metadata
- books_with_descriptions_cleaned_new.csv – LLM-generated descriptions

5.2 Metadata Schema

Table 2 describes the complete metadata schema used in the system.

Table 2: Book Metadata Schema

Column	Description
bookID	Internal unique identifier
title	Book title
authors	Author(s)
average_rating	Average user rating
isbn	ISBN-10
isbn13	ISBN-13
language_code	Language code
num_pages	Page count
ratings_count	Number of ratings
text_reviews_count	Number of reviews
publication_date	Date of publication
publisher	Publisher
isbn13_title	Concatenation of ISBN13 + title
description	LLM-generated short description
thumbnail	URL of book cover

6 Implementation Details

6.1 Book Description Generation

The description generation process utilizes the Falcon-7B-Instruct model to create engaging book descriptions.

```

1 from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
2 import pandas as pd
3
4 # Load dataset and prepare book list
5 books = pd.read_csv("books_cleaned.csv")
6 books_list = books["isbn13_title"].tolist()[1:2000]
7
8 # Initialize Falcon-7B-Instruct model
9 model_name = "tiiuae/falcon-7b-instruct"
10 tokenizer = AutoTokenizer.from_pretrained(model_name)
11 model = AutoModelForCausalLM.from_pretrained(
12     model_name,
13     device_map="auto",
14     torch_dtype=torch.bfloat16
15 )
16
17 pipe = pipeline(
18     "text-generation",
19     model=model,
20     tokenizer=tokenizer,
21     device_map="auto"
22 )
23

```

```

24 def generate_book_description(book_info, max_new_tokens=100):
25     prompt = f"Write a short, engaging description for the book: {
book_info}"
26     sequences = pipe(
27         prompt,
28         max_new_tokens=max_new_tokens,
29         do_sample=True,
30         top_k=10,
31         num_return_sequences=1,
32         pad_token_id=tokenizer.eos_token_id
33     )
34     return sequences[0]["generated_text"]
35
36 # Batch processing for memory optimization
37 results = []
38 batch_size = 500
39 for start in range(0, len(books_list), batch_size):
40     batch = books_list[start:start+batch_size]
41     for book in batch:
42         desc_text = generate_book_description(book, max_new_tokens=100)
43         results.append({
44             "isbn13_title": book,
45             "description": desc_text
46         })
47
48 # Save results
49 desc_df = pd.DataFrame(results)
50 desc_df.to_csv("books_with_descriptions.csv", index=False)

```

Listing 1: Book Description Generation Pipeline

Key Design Decisions:

- Batch processing optimizes GPU memory usage
- Sampling with `top_k=10` ensures diversity in descriptions
- Maximum token limit prevents overly long descriptions

6.2 Metadata Merging

The metadata merging process combines original book data with generated descriptions.

```

1 books = pd.read_csv("books_cleaned.csv")
2 descriptions = pd.read_csv("books_with_descriptions.csv")
3
4 # Merge datasets on common key
5 final_books_df = books.merge(descriptions, on="isbn13_title", how="
inner")
6 final_books_df.to_csv("final_books_df.csv", index=False)

```

Listing 2: Metadata Merging Process

6.3 Text Chunking for Embeddings

Text chunking optimizes embedding quality by ensuring appropriate context length.

```
1 from langchain.text_splitter import RecursiveCharacterTextSplitter
2
3 splitter = RecursiveCharacterTextSplitter(
4     chunk_size=200,
5     chunk_overlap=30
6 )
7
8 docs = []
9 for isbn, desc in zip(final_books_df["isbn13"], final_books_df["
    description"]):
10     chunks = splitter.split_text(desc)
11     for i, chunk in enumerate(chunks):
12         docs.append({
13             "isbn13": isbn,
14             "chunk_id": i,
15             "text": chunk
16         })
```

Listing 3: Text Chunking Implementation

Chunking Analysis:

- Most books generate 3–4 chunks
- Maximum observed chunks: 6 (rare occurrences)
- Overlap ensures context preservation across boundaries

6.4 Vector Embedding & Database

The vector database implementation uses Chroma for persistent storage and fast similarity search.

```
1 from langchain_community.embeddings import HuggingFaceEmbeddings
2 from langchain_community.vectorstores import Chroma
3 from langchain.schema import Document
4
5 # Initialize embedding model
6 embedding_fn = HuggingFaceEmbeddings(model_name="all-MiniLM-L6-v2")
7
8 # Prepare documents for vector storage
9 documents = [
10     Document(
11         page_content=doc["text"],
12         metadata={
13             "isbn13": doc["isbn13"],
14             "chunk_id": doc["chunk_id"]
15         }
16     ) for doc in docs
17 ]
18
19 # Create and persist vector database
20 vector_db = Chroma.from_documents(
21     documents=documents,
22     embedding=embedding_fn,
23     persist_directory="books_chroma_db"
24 )
```

```
25 vector_db.persist()
```

Listing 4: Vector Database Creation

Technology Choices:

- **all-MiniLM-L6-v2**: Fast, high-quality sentence embeddings
- **Chroma**: Provides persistence and efficient similarity search
- **Metadata storage**: Enables full book retrieval from chunk matches

6.5 Thumbnail Fetching

Book cover thumbnails enhance the user interface experience through visual representation.

```
1 import requests
2 import pandas as pd
3 import time
4
5 def fetch_thumbnail(isbn):
6     if pd.isna(isbn):
7         return "cover-not-found.jpg"
8
9     url = f"https://www.googleapis.com/books/v1/volumes?q=isbn:{isbn}"
10    try:
11        response = requests.get(url, timeout=10)
12        if response.status_code == 200:
13            data = response.json()
14            return data["items"][0]["volumeInfo"]["imageLinks"]["thumbnail"]
15        elif response.status_code == 429: # Rate limiting
16            time.sleep(30)
17            return fetch_thumbnail(isbn)
18        else:
19            return "cover-not-found.jpg"
20    except:
21        return "cover-not-found.jpg"
22
23 # Apply thumbnail fetching
24 books_df["thumbnail"] = books_df["isbn13"].apply(fetch_thumbnail)
25 books_df.to_csv("books_with_thumbnails.csv", index=False)
```

Listing 5: Thumbnail Fetching via Google Books API

6.6 Search Pipeline

The search pipeline combines vector similarity search with metadata matching for comprehensive results.

```
1 def search_books(query, db, books, k=5):
2     """
3     Search for books using semantic similarity
4
5     Args:
6         query: User search query
7         db: Chroma vector database
```



```

8         books: DataFrame containing book metadata
9         k: Number of results to return
10
11     Returns:
12         DataFrame of matched books with metadata
13     """
14     results = db.similarity_search(query, k=k)
15     matched_books = []
16
17     for r in results:
18         isbn = r.metadata.get("isbn13", None)
19         if isbn:
20             match = books[books["isbn13"] == int(isbn)]
21             if not match.empty:
22                 for _, row in match.iterrows():
23                     matched_books.append(row.to_dict())
24
25     return pd.DataFrame(matched_books)
26
27 # Example usage
28 df_matched = search_books("books about love", db, books, k=5)
29 print(df_matched)

```

Listing 6: Integrated Search Pipeline

6.7 Gradio Frontend Interface

The interactive frontend provides an intuitive gallery-based interface for book discovery.

```

1 import gradio as gr
2
3 def recommend_books(query):
4     """Generate book recommendations based on user query"""
5     df = search_books(query, db, books, k=10)
6     return [
7         (row['thumbnail'], f"{row['title']} by {row['authors']}")
8         for _, row in df.iterrows()
9     ]
10
11 # Create Gradio interface
12 with gr.Blocks() as demo:
13     gr.Markdown("# Book Recommender System")
14
15     with gr.Row():
16         user_query = gr.Textbox(
17             label="Enter a description or keyword:",
18             placeholder="e.g., 'fantasy novels with magic'"
19         )
20         submit_btn = gr.Button("Find Books")
21
22     output_gallery = gr.Gallery(
23         label="Recommended Books",
24         columns=3,
25         rows=2
26     )
27
28     submit_btn.click(
29         fn=recommend_books,

```

```
30     inputs=user_query ,
31     outputs=output_gallery
32 )
33
34 # Launch application
35 demo.launch(server_port=7860, server_name="0.0.0.0")
```

Listing 7: Gradio Interactive Interface

7 Key Design Decisions

7.1 Model Selection

- **Falcon-7B-Instruct**: Chosen for human-like description generation capabilities
- **all-MiniLM-L6-v2**: Balance between speed and embedding quality
- **Chroma**: Efficient vector database with persistence support

7.2 Technical Architecture

- **Chunking Strategy**: 200 characters with 30-character overlap optimizes embedding effectiveness
- **Batch Processing**: Memory-efficient description generation
- **API Integration**: Google Books API for enhanced visual experience
- **UI Framework**: Gradio enables rapid deployment with interactive galleries

8 Performance Analysis

8.1 System Metrics

- **Dataset Size**: 2,000 books processed
- **Average Chunks per Book**: 3–4 chunks
- **Embedding Dimensions**: 384 (all-MiniLM-L6-v2)
- **Search Response Time**: < 1 second for typical queries

8.2 Quality Assessment

The system demonstrates strong semantic understanding through several test cases:

- Query: "books about love" → Successfully identifies romance novels
- Query: "fantasy adventure" → Returns fantasy genre books with adventure themes
- Query: "historical fiction" → Accurately matches historical narrative books

9 Future Improvements

9.1 Performance Enhancements

- Implement parallel processing for description generation
- Add caching mechanisms for API calls
- Optimize embedding model selection for domain-specific tasks

9.2 Feature Extensions

- Add filtering capabilities (language, rating, publication date)
- Incorporate user feedback for recommendation refinement
- Implement collaborative filtering alongside content-based recommendations
- Develop mobile-responsive interface

9.3 Scalability Considerations

- Database partitioning for larger datasets
- Distributed embedding computation
- Load balancing for production deployment

10 Conclusion

This project successfully demonstrates the development of a robust, interactive, and intelligent book recommendation system. The integration of Large Language Models, semantic vector search, and modern UI technologies creates a comprehensive solution that addresses the limitations of traditional keyword-based search systems.

The system showcases end-to-end capabilities from raw data processing to deployable frontend application, with detailed documentation of the complete pipeline. The modular architecture ensures maintainability and extensibility for future enhancements.

Key achievements include:

- Successful integration of LLM-generated content with traditional metadata
- Implementation of efficient semantic search using vector embeddings
- Development of an intuitive, visual user interface
- Creation of a scalable, persistent vector database solution

This work provides a foundation for advanced recommendation systems and demonstrates the practical application of modern NLP technologies in information retrieval tasks.

About the Author

Mohamed Amine Mammar El Hadj is a Deep Learning & Machine Learning Developer and Software Developer focusing on NLP, vector search, and AI-powered applications. He has experience building end-to-end deep learning pipelines, including data pre-processing, embedding-based retrieval systems, LLM integration, and deployment-ready interfaces. This book recommender system showcases his practical skills in applying deep learning techniques to real-world problems.

For more projects and technical insights, connect with Mohamed Amine:

- **LinkedIn:** Mohamed Amine Mammar El Hadj
- **GitHub:** github.com/mimine47
- **Email:** mohamedamine.devtech@gmail.com

References

1. Falcon-7B-Instruct Model Documentation, Technology Innovation Institute
2. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, Reimers & Gurevych, 2019
3. LangChain: Building Applications with LLMs, Harrison Chase et al.
4. Chroma: Open-source embedding database documentation
5. Google Books API Documentation, Google Developers