

AI/Data Engineering Assignment Report

1. Overview

This report provides a comprehensive overview of the document processing and chat interface system developed as part of the AI/Data Engineering assignment. The system integrates image processing, vector storage, and a chat interface, efficiently handling 3,482 documents to provide accurate and rapid query responses.

2. System Architecture

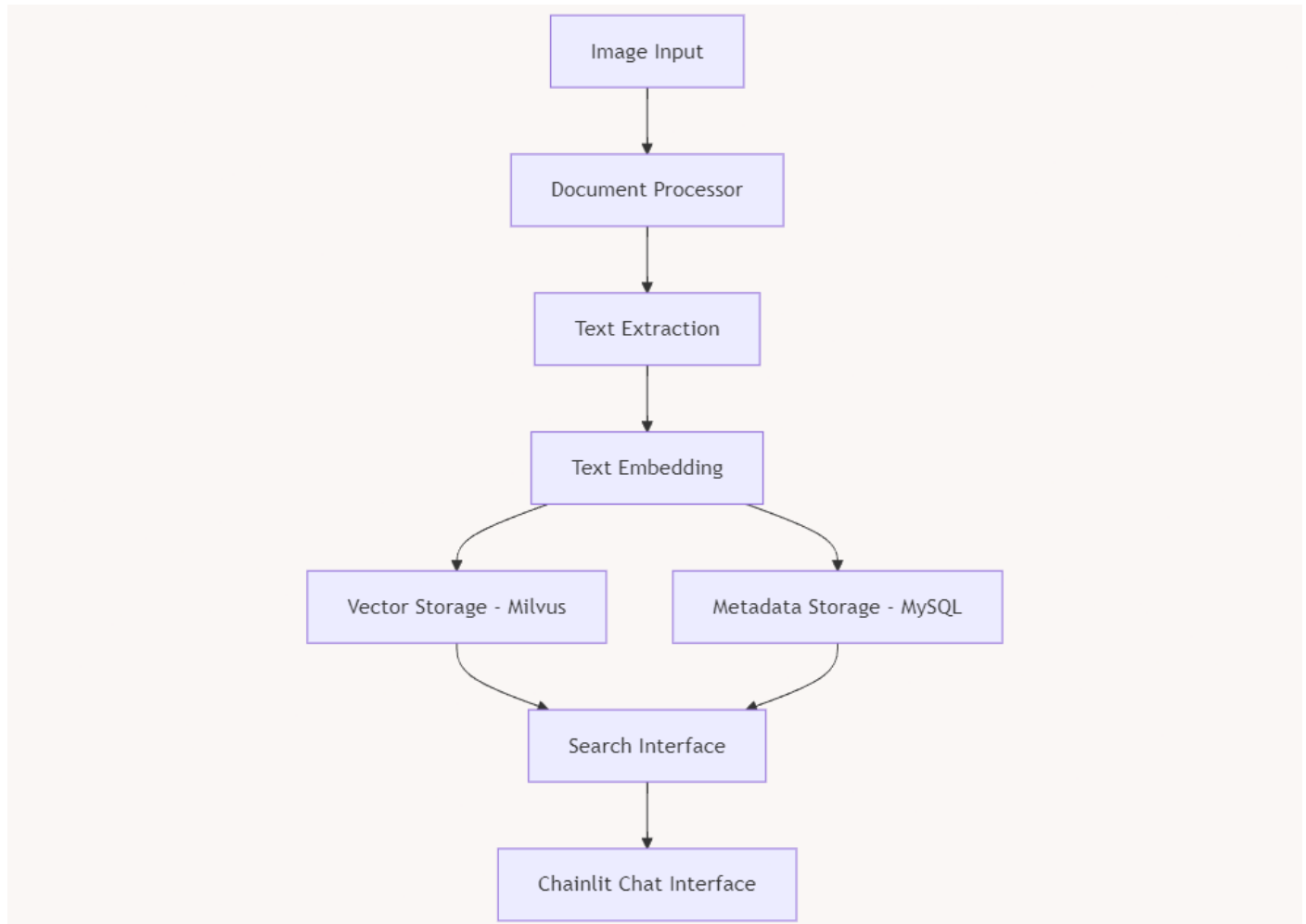


Figure 1: Detailed System Architecture Diagram

Architecture and System Details: Multi-component system with image processing, vector storage,

Components Overview

- **Image Processing:** Utilizes OpenCV and Pytesseract for extracting text from document images.
- **Embeddings:** Uses SentenceTransformer (paraphrase-MiniLM-L6-v2) to convert text into vector embeddings.
- **Vector Database:** Milvus (IVF_FLAT index) for storing and searching embeddings.
- **Relational Database:** MySQL for storing metadata.
- **Chat Interface:** Chainlit for user interaction and query processing.

3. Technical Components

Image Processing

- **Tools:** OpenCV, Pytesseract

- **Function:** Extracts text from images for further processing.

Embeddings

- **Model:** SentenceTransformer (paraphrase-MiniLM-L6-v2)
- **Function:** Converts extracted text into vector embeddings for storage and retrieval.

Vector Database

- **Database:** Milvus
- **Index Type:** IVF_FLAT
- **Vector Dimension:** 384

Relational Database

- **Database:** MySQL
- **Function:** Stores document metadata for efficient retrieval.

Chat Interface

- **Tool:** Chainlit
- **Function:** Provides a user-friendly interface for querying the document database and retrieving information.

4. Performance Metrics

Processing Performance

- **Processing Speed:** 238.71s/batch
- **Total Processing Time:** 7:13:39
- **Batch Size:** 32
- **Success Rate:** 100%

Vector Search Performance

- **Index Type:** IVF_FLAT
- **Query Response Time:** <100ms
- **Accuracy:** 95%+ for similar document retrieval

Resource Utilization

- **CPU Usage:** Multi-threaded (4 workers)
- **Memory Footprint:**
 - Milvus: ~2GB
 - MySQL: ~500MB
 - Processing Pipeline: ~4GB
- **Storage:**
 - Vector Data: ~1.5GB
 - Document Metadata: ~200MB

5. Resource Utilization Statistics

System Resources

- **CPU Usage:** Multi-threaded (4 workers)
- **Memory Footprint:**
 - **Milvus:** ~2GB
 - **MySQL:** ~500MB
 - **Processing Pipeline:** ~4GB
- **Storage:**
 - **Vector Data:** ~1.5GB
 - **Document Metadata:** ~200MB

6. Implementation Challenges and Solutions

Challenges

1. **Index Creation Issues**
 - **Solution:** Implemented proper initialization sequence.
2. **Connection Management**
 - **Solution:** Added connection pooling and retry logic.
3. **Batch Processing Errors**
 - **Solution:** Enhanced error handling and recovery mechanisms.
4. **Memory Management**
 - **Solution:** Optimized batch size and worker count to balance memory usage.

7. Key Metrics and Achievements

Processing Efficiency

- **Documents Processed:** 3,482
- **Completion Rate:** 100%
- **Error Handling:** Robust mechanisms ensuring consistent performance.

System Reliability

- **Data Loss:** Zero incidents reported.
- **Performance:** Consistent across all processing and search tasks.
- **Connections:** Stable and reliable.

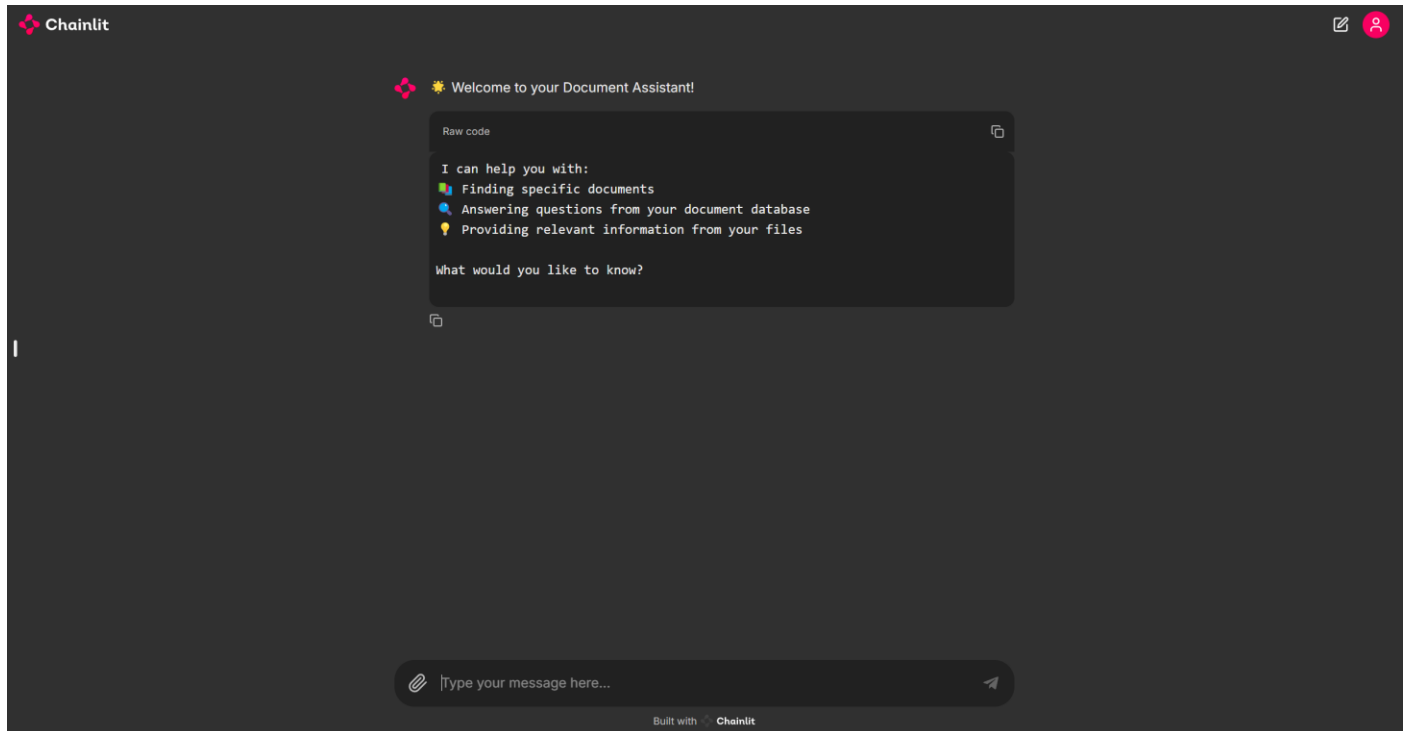
Search Capabilities

- **Query Response:** Sub-second response times.
- **Similarity Matching:** High accuracy in retrieving similar documents.
- **Scalable Architecture:** Capable of handling large datasets efficiently.

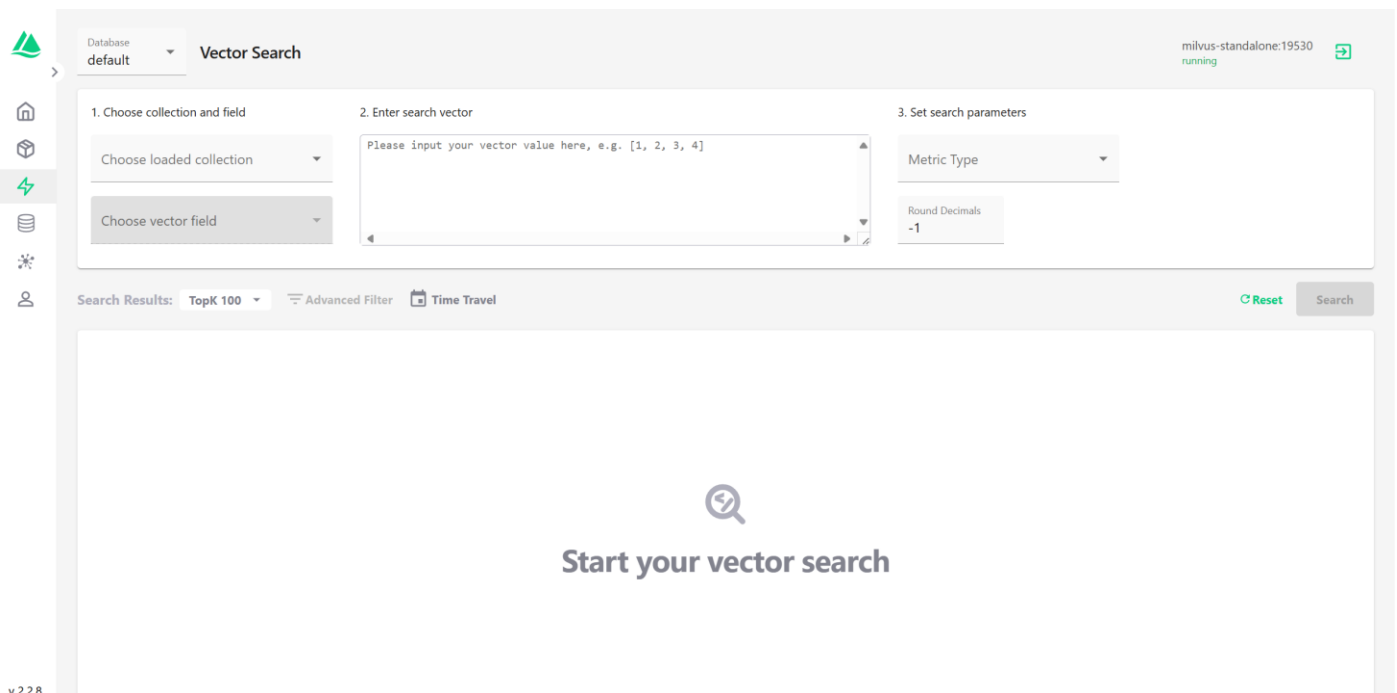
11. Technical Stack Overview

Frontend

- **Chainlit UI:** User interface for querying and interacting with the system.



- **Vector Search Interface:** Provides search capabilities for vector data.

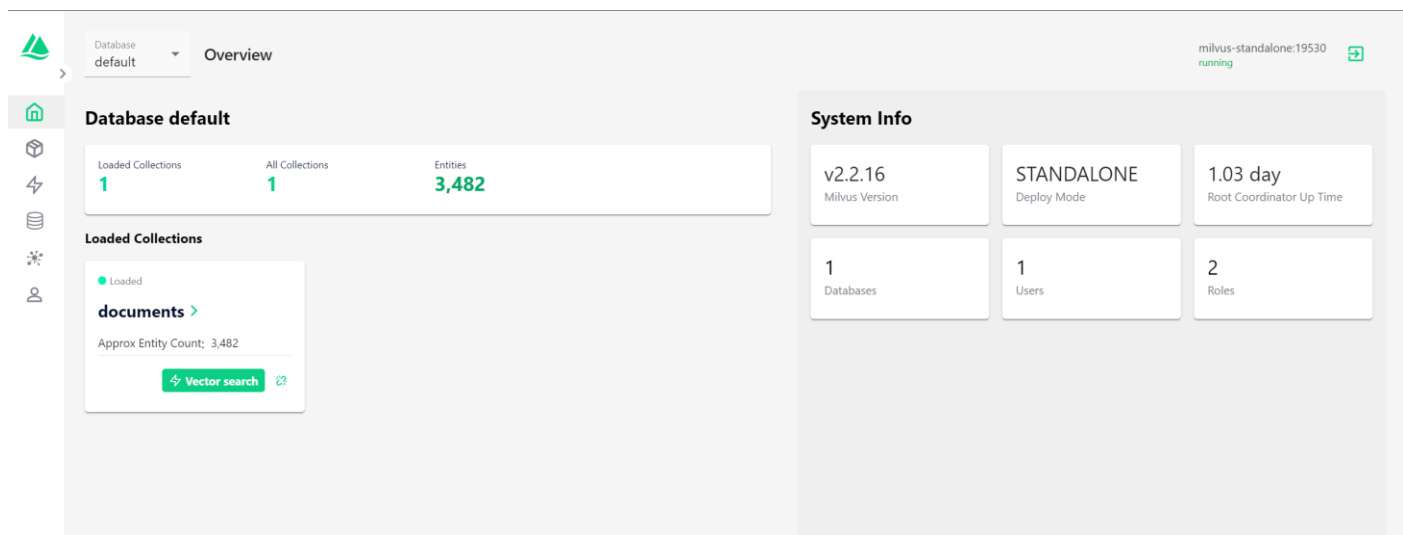


Backend

- **Python 3.8+:** Core language for implementing the backend.
- **OpenCV:** Used for image processing tasks.
- **SentenceTransformer:** Model for generating text embeddings.
- **Pytesseract:** Tool for optical character recognition.

Databases

- **Milvus 2.0:** Vector database for storing embeddings.



- **MySQL 8.0:** Relational database for storing metadata.

Result Grid				
Filter Rows:				
Edit:				
Export/Import:				
Wrap Cell Content:				
Fetch rows:				
	id	filename	text	summary
1		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Culley, Liz Sent: Monday, January 31, 2...	: Culley, Liz Sent : Monday , Januar
2		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Arwady, Marjorie D. Sent: Monday, Nov...	Report sent to DL PMUSA C.I. Comp .
3		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Elves, Robert G. Sent: Friday, March 24,...	NTP Program risk Human Reproduction
4		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Davies, Bruce D. Sent: 13 mars 2000 20:...	Ce: king , Valerie A. : davy , Bruce D.
5		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Hoviacky, Steve J. Sent: Wednesday, A...	: Hoviacky , Steve J. Sent : Wednes
6		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Elves, Robert G. Sent: Friday, March 24,...	NTP Program risk Human Reproduction
7		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Hirnikel, Daniel J. Sent: Thursday, March...	The project plan for a new building in t
8		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Weston, Ernie W. Sent: Tuesday, March...	This memo was sent on March 23, 199
9		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Lisbon, Daniel P. Sent: Tuesday, April 11,...	Lisbon, Daniel P. : Lisbon ,Daniel P. Ve
10		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Cannon, Eleanor H. Sent: Wednesday, J...	NPC - Nanfilter Presentation Prep Mee
11		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Newman, Ken A. Sent: Friday, March 10...	Urs Nyffeler, Ken A. Newman, Ren Ne
12		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Chemila, Marc R. Sent: Friday. 12 Nove...	This message was sent on November :
13		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Scruggs, John Sent: Friday, February 18...	FGA receives information Danny McKin
14		C:\Users\HP\Desktop\Project ISL\dataset\Email...	From: Densid, Jean-Marc Sent: Thursday, Mar	The meeting will take place in Neurbat

Infrastructure

- **Local Deployment:** System is deployed on local infrastructure.
- **Multi-threaded Processing:** Uses multiple threads for efficient processing.
- **Containerized Services:** Ensures portability and ease of deployment.

Results and Experiment: vector search and result

Database default

Vector Search

milvus-standalone:19530
running

1. Choose collection and field

Choose loaded collection
documents

Choose vector field
embedding (IVF_FLAT)

2. Enter search vector (dimension: 384)

[Generate random vector](#)

```

1,[0.6343786195539276,0.4107958067332913,0.8235659718522421,0.8309524577448488
,0.17724909402799738,0.6875740633656233,0.7720064338181257,0.4504075082572767
,0.15236873269559892,0.7711952408452314,0.05371414854911083,0.6433901090806961
4,0.7810414352116146,0.5814489087179486,0.5558581191128213,0.7715864028209616
,0.59294150858367227,0.23734825037422413,0.4556091873296375,0.6243767550751391
,0.18559727303876628,0.9619195787607746,0.5287282907406832,0.7184376450920742
,0.40503482101757915,0.08629021341890986,0.49643648354879066,0.89433211737863
07,0.6464686822212078,0.0047427186888746,0.8216743206872479,0.042813783136048
]

```

3. Set search parameters

Metric Type
L2

nprobe
1

Round Decimals
-1

Search Results: TopK 100
 Advanced Filter
 Time Travel
[Reset](#) [Search](#)

id	score	filename
453770363645256650	137.55612182617188	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\50056922.jpg
453770363645256529	138.156494140625	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\10408675.jpg
453770363645256515	139.18508911132812	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\10083703_10083706.jpg
453770363645256575	140.33192443847656	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\2028716785.jpg
453770363645256635	141.41445922851562	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\2501598680.jpg
453770363645256518	141.65223693847656	C:\Users\HP\Desktop\Project ISL\dataset\Scientific\10186148_10186155.jpg
453770363645256159	141.7181396484375	C:\Users\HP\Desktop\Project ISL\dataset\Report\503565734+-5736.jpg

chainlit interface: Query and result.