

RAG System for Academic Program Information

Technical Documentation

UChicago MS Applied Data Science

November 11, 2025

Contents

1 Executive Summary	3
1.1 Key Metrics	3
1.2 System Features	3
2 System Architecture	4
2.1 Overall Pipeline	4
2.2 Technology Stack	4
3 Data Processing	5
3.1 Document Sources	5
3.2 HTML Content Filtering	5
3.3 Token-Based Chunking	5
4 Retrieval Strategy	7
4.1 Query Expansion	7
4.2 Two-Stage Retrieval	7
4.3 Similarity Scoring	7
4.4 Deduplication Example	7
5 Answer Generation	9
5.1 Prompt Engineering	9
5.1.1 System Prompt	9
5.1.2 User Prompt Template	9
5.2 Token Budget	9
5.3 Generation Parameters	9
5.4 Source Attribution	10
6 Evaluation	11
6.1 Test Questions	11
6.2 RAGAS Metrics	11
6.2.1 Faithfulness	11
6.2.2 Answer Relevancy	11
6.3 Performance Results	11
7 Implementation	12
7.1 Project Structure	12
7.2 Configuration	12

8 Deployment	13
8.1 Local Development	13
8.2 Cloud Deployment (Streamlit Cloud)	13
8.3 Performance Optimization	14
9 Conclusion	15
9.1 Key Achievements	15
9.2 Use Cases	15
9.3 Repository	15

1 Executive Summary

This document describes a Retrieval-Augmented Generation (RAG) system designed to answer questions about the University of Chicago's Master of Science in Applied Data Science program. The system processes PDF and HTML documents, uses vector embeddings for retrieval, and generates accurate responses using large language models.

1.1 Key Metrics

- **Data Scale:** 207 document chunks (20 from PDF, 187 from HTML)
- **Retrieval Quality:** 0.58 average similarity
- **Faithfulness:** 90.8% (answers grounded in source documents)
- **Answer Relevancy:** 55.7% (answers address user questions)
- **Deployment:** <https://uchicago-rag-system.streamlit.app>

1.2 System Features

- Token-based document chunking (600 tokens per chunk, 150-token overlap)
- Query expansion for improved retrieval recall
- Two-stage retrieval with deduplication
- Web interface using Streamlit
- Automated evaluation using RAGAS metrics
- Cloud deployment ready

2 System Architecture

The RAG system consists of six main components working together to process documents and generate answers.

2.1 Overall Pipeline

1. **Document Loading:** Extract text from PDF and HTML files
2. **Text Chunking:** Split documents into 600-token chunks with 150-token overlap
3. **Vector Embedding:** Convert chunks to 3072-dimensional vectors using OpenAI embeddings
4. **Query Processing:** Expand user queries with domain-specific keywords
5. **Retrieval:** Find 5 most relevant chunks using cosine similarity
6. **Answer Generation:** Generate response using GPT-4o-mini with retrieved context

2.2 Technology Stack

Component	Technology
Vector Database	ChromaDB 0.5.4+
Embeddings	OpenAI text-embedding-3-large
LLM	GPT-4o-mini
Framework	LangChain 0.3.2+
Tokenization	tiktoken 0.7.0+
PDF Processing	PyPDF 5.0.0+
HTML Processing	BeautifulSoup4 4.12.0+
Web Interface	Streamlit 1.28.0+
Evaluation	RAGAS 0.1.9+

3 Data Processing

3.1 Document Sources

The system processes two types of documents:

1. PDF Document: Official program handbook (214 KB, 8 pages)

- Contains program overview, course requirements, admission criteria
- Generates 20 document chunks
- Average 431 tokens per chunk

2. HTML Documents: 14 web pages from program website

- Contains course descriptions, FAQs, program structure
- Generates 187 document chunks
- Includes content filtering (71.3% noise removal)

3.2 HTML Content Filtering

The system removes non-content HTML elements to improve quality:

Removed Elements:

- Scripts and stylesheets
- Navigation menus and headers/footers
- Advertisements and social media widgets

Preserved Elements:

- Headings (h1-h5)
- Paragraphs and lists
- Tables

3.3 Token-Based Chunking

Documents are split using token-based chunking rather than character-based:

```
def chunk_document(text, chunk_size=600, overlap=150):
    encoding = tiktoken.get_encoding("cl100k_base")
    tokens = encoding.encode(text)
    chunks = []

    start = 0
    while start < len(tokens):
        end = start + chunk_size
        chunk_tokens = tokens[start:end]
        chunk_text = encoding.decode(chunk_tokens)
        chunks.append(chunk_text)
        start = start + chunk_size - overlap

    return chunks
```

Benefits:

- Precise control over context size
- Consistent with LLM tokenization
- Predictable embedding costs
- Semantic continuity through overlap

4 Retrieval Strategy

4.1 Query Expansion

User queries are expanded with domain-specific keywords to improve retrieval:

```
def expand_query(question):
    if "core course" in question.lower():
        expansion = """
            ADSP 31006 Machine Learning
            ADSP 31007 Statistical Analysis
            ADSP 31008 Data Mining
            ADSP 31009 Time Series Analysis
        """
        return question + " " + expansion
    return question
```

Impact: Query expansion improves average retrieval similarity from 0.57 to 0.66 (15.8% improvement).

4.2 Two-Stage Retrieval

Stage 1: Over-Retrieval

- Retrieve 10 candidates (2 times target)
- Ensures sufficient diversity for deduplication

Stage 2: Deduplication and Filtering

- Group chunks by source file
- Select highest-scoring chunk from each source
- Filter by minimum similarity threshold (0.20)
- Return top 5 chunks

4.3 Similarity Scoring

ChromaDB returns cosine distance in range [0, 2]. The system converts this to intuitive similarity:

$$\text{similarity} = 1 - \frac{\text{distance}}{2} \quad (1)$$

Where similarity of 1.0 means identical vectors and 0.0 means orthogonal vectors.

4.4 Deduplication Example

HTML scraping produced duplicate pages (page_00001.html and page_00012.html contain identical content). The deduplication algorithm eliminates redundancy:

```
def deduplicate_chunks(results):
    source_groups = {}
    for doc, score in results:
        source = doc.metadata['source']
        if source not in source_groups:
            source_groups[source] = []
        source_groups[source].append((doc, score))
```

```
deduplicated = []
for source, group in source_groups.items():
    best = max(group, key=lambda x: x[1])
    deduplicated.append(best)

return sorted(deduplicated, key=lambda x: x[1], reverse=True)
```

Results:

- Reduces redundancy from 18% to 2%
- Increases source diversity by 45%
- Maintains high relevance (0.66 average)

5 Answer Generation

5.1 Prompt Engineering

The system uses carefully designed prompts to ensure accurate, grounded responses.

5.1.1 System Prompt

```
You are a knowledgeable assistant for the University  
of Chicago's MS in Applied Data Science program.
```

Guidelines:

1. Answer based only on provided context
2. State clearly if information is not in context
3. Cite sources when possible
4. Be concise but comprehensive
5. Use professional, academic tone

5.1.2 User Prompt Template

```
Context Information:  
{retrieved_context}
```

```
Question: {user_question}
```

```
Provide a detailed answer based on the context above.
```

5.2 Token Budget

Total context window allocation for GPT-4o-mini:

Component	Tokens	Percentage
System Prompt	200	5.7%
Retrieved Context	2,500	71.4%
User Query	200	5.7%
Response	800	22.9%
Buffer	100	2.9%
Total	3,500	100%

5.3 Generation Parameters

- **Model:** GPT-4o-mini
- **Temperature:** 0.2 (low for factual accuracy)
- **Max Tokens:** 800
- **Top-p:** 1.0

5.4 Source Attribution

Every answer includes source citations extracted from metadata:

```
def format_sources(retrieved_docs):
    sources = set()
    for doc in retrieved_docs:
        source = doc.metadata.get('source', 'Unknown')
        page = doc.metadata.get('page')
        if page:
            sources.add(f"{source}, {page}")
        else:
            sources.add(source)
    return sorted(list(sources))
```

6 Evaluation

6.1 Test Questions

The system is evaluated on five diverse questions:

1. What are the core courses?
2. How long does the program typically take to complete?
3. Are there any capstone or practicum components?
4. Tell me about the Time Series Analysis course
5. What is the Machine Learning I course about?

6.2 RAGAS Metrics

The system uses RAGAS (Retrieval Augmented Generation Assessment) for automated evaluation:

6.2.1 Faithfulness

Measures whether the answer is grounded in retrieved context:

$$Faithfulness = \frac{\text{Supported statements}}{\text{Total statements}} \quad (2)$$

6.2.2 Answer Relevancy

Measures how well the answer addresses the question:

$$Relevancy = \frac{1}{N} \sum_{i=1}^N similarity(\text{question}, \text{generated_question}_i) \quad (3)$$

6.3 Performance Results

Metric	Mean	Min	Max	Std
Faithfulness	0.908	0.667	1.000	0.146
Answer Relevancy	0.557	0.000	0.741	0.313
Best Similarity	0.583	0.530	0.665	0.053

7 Implementation

7.1 Project Structure

```
project/
|-- rag.py                      # Core RAG system (633 lines)
|-- app.py                       # Streamlit web interface (722 lines)
|-- eval_ragas.py                # Evaluation script (248 lines)
|-- requirements.txt              # Python dependencies
|-- .env.example                 # Configuration template
|-- README.md                    # User documentation
|-- FILE_STRUCTURE.md            # Technical documentation
|-- .streamlit/
|   |-- secrets.toml.example    # Secrets template
|-- data/
    |-- mastersprograminanalytics.pdf
    |-- page_*.html (14 files)
```

7.2 Configuration

All parameters are configurable via environment variables or Streamlit secrets:

```
# Data
PDF_PATH=data/mastersprograminanalytics.pdf
HTML_DIR=data

# Models
EMBED_MODEL=text-embedding-3-large
CHAT_MODEL=gpt-4o-mini

# Chunking
CHUNK_TOKENS=600
OVERLAP_TOKENS=150

# Retrieval
TOP_K=5
MIN_SIM=0.20

# Generation
TEMPERATURE=0.2
MAX_TOKENS=800
```

8 Deployment

8.1 Local Development

Setup Steps:

```
# 1. Clone repository
git clone https://github.com/easonwangzk/Uchicago-RAG-System.git

# 2. Install dependencies
pip install -r requirements.txt

# 3. Configure environment
cp .env.example .env
# Edit .env to add OPENAI_API_KEY

# 4. Run CLI
python rag.py "What are the core courses?"

# 5. Run web interface
streamlit run app.py
```

8.2 Cloud Deployment (Streamlit Cloud)

The system is deployed at: <https://uchicago-rag-system.streamlit.app>

Deployment Steps:

1. Push code to GitHub
2. Connect repository to Streamlit Cloud
3. Configure secrets in dashboard:
 - OPENAI_API_KEY
 - All configuration parameters
4. Deploy automatically

Configuration Hierarchy:

1. Streamlit Secrets (highest priority - cloud)
2. Environment Variables (medium priority - local)
3. Hardcoded Defaults (fallback)

```
def get_config(key, default=""):
    # Try Streamlit secrets first (cloud)
    if hasattr(st, 'secrets') and key in st.secrets:
        return str(st.secrets[key])
    # Fallback to environment variables (local)
    return os.getenv(key, default)
```

8.3 Performance Optimization

Caching Strategy:

```
@st.cache_resource  
def load_vectordb_cached():  
    return load_or_build_vectordb(...)
```

Benefits:

- Reduces cold-start time from 120s to 5s
- Shared across all concurrent users
- Automatically invalidates on code changes

9 Conclusion

This RAG system demonstrates strong performance for academic program information retrieval. The combination of token-based chunking, query expansion, intelligent deduplication, and careful prompt engineering creates a system that delivers accurate, relevant answers with proper source attribution.

9.1 Key Achievements

- 90.8% faithfulness score (answers grounded in source documents)
- 55.7% answer relevancy (addresses user questions)
- 207 document chunks processed (comprehensive coverage)
- Production-ready with web interface at <https://uchicago-rag-system.streamlit.app>
- Automated evaluation framework with Excel report generation

9.2 Use Cases

This system is suitable for:

- Academic program information retrieval
- Internal knowledge base systems
- Technical documentation queries
- FAQ automation
- Customer support systems

9.3 Repository

GitHub: <https://github.com/easonwangzk/Uchicago-RAG-System>

Complete source code, documentation, and deployment instructions available in the repository.