

RAG Chatbot System Technical Documentation

Retrieval-Augmented Generation Pipeline Implementation Report

Project: RAG Chatbot System

Technology Stack: Python, Groq API, FAISS, Streamlit

Generated: August 13, 2025

Version: 1.0

1. System Overview

This RAG (Retrieval-Augmented Generation) chatbot system combines document retrieval with large language model generation to provide accurate, contextual responses. The system uses the Groq API for fast inference with Meta's LLaMA 3.1 models, FAISS for efficient vector similarity search, and Streamlit for the user interface. **Key Components:**

- Document Processing Pipeline
- Vector Database with Semantic Search
- RAG Pipeline with Context Retrieval
- Groq API Integration for Response Generation
- Interactive Streamlit Web Interface

2. System Architecture

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

Component	Technology	Purpose
Document Processor	LangChain, PyPDF, spaCy	Load, clean, and chunk documents
Vector Store	FAISS, SentenceTransformers	Store and search document embeddings
RAG Pipeline	Custom Python Classes	Orchestrate retrieval and generation
LLM Interface	Groq API (LLaMA 3.1)	Generate responses from context
Web Interface	Streamlit	Provide user-friendly chat interface

3. Document Processing & Chunking

3.1 Document Structure

The system supports multiple document formats including PDF, DOCX, TXT, and HTML files. Documents are loaded using LangChain's document loaders and processed through a cleaning pipeline. **Supported Formats:**

- PDF files (via PyPDFLoader)
- Microsoft Word documents (via Docx2txtLoader)
- Plain text files (via TextLoader)
- HTML files (via BeautifulSoup parsing)

3.2 Chunking Logic

Documents are split into semantic chunks using a hybrid approach: **Primary Strategy:**

SpacyTextSplitter with linguistic boundaries

Fallback Strategy: RecursiveCharacterTextSplitter with hierarchical separators

Chunk Size: 300 characters (configurable)

Overlap: 50 characters to maintain context continuity

Separators: [paragraph, sentence, clause, word, character]

```
# Chunking Configuration
chunk_size = 300
chunk_overlap = 50
separators = ["\n\n", "\n", ". ", "! ", "? ", " ", ""]

# Text cleaning pipeline
text = re.sub(r'\s+', ' ', text)
# Normalize whitespace
text = re.sub(r'^\w\s.\!?\,\:\;\-\(\)\[\]\{\}\|\'\\"', '', text)
text =
```

```
re.sub(r'([.!?])\1+', r'\1', text) # Remove repeated punctuation
```

4. Embedding Model & Vector Database

4.1 Embedding Model

Model: sentence-transformers/all-MiniLM-L6-v2

Embedding Dimension: 384

Max Sequence Length: 512 tokens

Normalization: L2 normalization enabled

Performance: ~14k sentences/second on CPU This model was chosen for its excellent balance between performance and quality, providing high-quality embeddings while maintaining fast inference speed suitable for real-time applications.

4.2 Vector Database (FAISS)

Technology: Facebook AI Similarity Search (FAISS)

Index Type: Flat L2 (exact search for small-medium datasets)

Distance Metric: Cosine similarity

Storage: Persistent disk storage with JSON metadata

Search Method: k-NN with configurable similarity threshold FAISS provides efficient similarity search with sub-millisecond query times for datasets up to 100k documents. The system uses LangChain's FAISS wrapper for seamless integration.

5. Prompt Engineering & Generation

5.1 Prompt Template

The system uses a carefully engineered prompt template to ensure high-quality responses:

SYSTEM PROMPT: You are an expert AI assistant that provides detailed, accurate, and well-structured answers. You have access to relevant information that you should use to answer questions comprehensively. Instructions for high-quality responses: 1. Provide COMPREHENSIVE answers that fully address the user's question 2. Structure your response with clear sections, bullet points, or numbered lists 3. Include SPECIFIC details and examples when relevant 4. Use a professional, informative tone while being accessible 5. Do not reference "context documents" or numbered sources 6. Integrate information seamlessly as if it's your own knowledge Relevant Information: {context} Please provide a detailed and natural answer to the following question:

5.2 Generation Parameters

Model: llama-3.1-8b-instant (via Groq API)

Temperature: 0.1 (low for factual accuracy)

Max Tokens: 2048

Top-K Retrieval: 5 most relevant chunks

Score Threshold: 0.1 (minimum similarity)

Context Window: ~8k tokens effective context

6. Example Queries & Responses

6.1 Successful Query Example

Query	What are the main features of artificial intelligence?
Context Score	0.87
Chunks Retrieved	5
Response Time	1.2s

Response:

Artificial intelligence encompasses several key features that enable machines to simulate human-like cognitive abilities. The main features include: **Machine Learning**: The ability to learn patterns from data without explicit programming, enabling systems to improve performance over time. **Natural Language Processing**: Understanding and generating human language, allowing AI to communicate effectively with users through text and speech. **Computer Vision**: Processing and interpreting visual information from images and videos to recognize objects, faces, and scenes. **Decision Making**: Analyzing complex scenarios and making informed decisions based on available data and predefined criteria. **Automation**: Performing tasks automatically without human intervention, increasing efficiency and reducing manual effort.

Analysis:

SUCCESS: High relevance score (0.87), comprehensive response covering all major aspects, well-structured format.

6.2 Moderate Success Example

Query	How does machine learning work in practice?
Context Score	0.64
Chunks Retrieved	4
Response Time	0.9s

Response:

Machine learning works through a systematic process of training algorithms on data to make predictions or decisions. The practical implementation involves several steps: **Data Collection**: Gathering relevant datasets that represent the problem domain. **Data Preprocessing**: Cleaning and preparing data by handling missing values, normalizing features, and splitting into training/testing sets. **Model Selection**: Choosing appropriate algorithms based on the problem type (supervised, unsupervised, or reinforcement learning). **Training Process**: The algorithm learns patterns by adjusting internal parameters based on the training data.

Analysis:

MODERATE: Good technical accuracy but somewhat generic response. Could benefit from more domain-specific examples.

6.3 Failure Case Example

Query	What is the latest research in quantum computing?
Context Score	0.23
Chunks Retrieved	2
Response Time	0.7s

Response:

I don't have enough specific information about the latest research in quantum computing to provide a comprehensive answer. Based on the available information, quantum computing is an emerging field that uses quantum mechanical phenomena to process information differently than classical computers. For the most current research developments, I'd recommend checking recent publications from quantum computing research institutions or academic journals.

Analysis:

FAILURE: Low relevance score (0.23) indicates insufficient context. System correctly identifies limitation rather than hallucinating.

7. System Limitations & Challenges

7.1 Hallucination Mitigation

Challenge: LLMs can generate plausible-sounding but factually incorrect information. **Mitigation Strategies:**

- Low temperature setting (0.1) to reduce creativity and maintain factual accuracy
 - Explicit instructions to avoid referencing non-existent sources
 - Confidence thresholding - system admits uncertainty when context relevance is low
 - Prompt engineering to encourage acknowledging information limitations
 - Context-grounded generation - responses must be based on retrieved documents
- Monitoring:** Manual review of responses and user feedback collection for quality assessment.

7.2 Performance Considerations

Response Latency:

- Average: 1-2 seconds (acceptable for most use cases)
 - Factors: Document retrieval (~100ms) + LLM generation (~800-1500ms)
 - Optimization: Groq API provides faster inference than standard cloud APIs
- Scalability Constraints:**
- FAISS flat index scales to ~100k documents before requiring approximate search
 - Embedding computation is CPU-bound (GPU acceleration available)
 - Concurrent users limited by Groq API rate limits
- Context Window Limitations:**
- Maximum 5 chunks retrieved per query
 - Chunk size optimized for balance between context and coherence
 - Long documents may lose important context if poorly chunked

7.3 Quality Challenges

Domain Specificity:

- Generic embedding model may not capture domain-specific terminology well
 - Recommendation: Fine-tune embeddings on domain-specific data
 - Current model performs well on general knowledge but may struggle with technical jargon
- Document**

Quality Dependency:

- System quality is directly tied to source document quality
 - Outdated or incorrect source documents propagate errors
 - Inconsistent document formatting affects chunking quality
- Multi-hop Reasoning:**
- Limited ability to synthesize information across multiple document sections
 - Complex queries requiring reasoning across disparate facts may fail
 - Current context window cannot handle very broad topics effectively

8. Recommendations & Future Improvements

Immediate Improvements:

- Implement query preprocessing for better retrieval (query expansion, spell correction)
- Add document relevance scoring to filter low-quality sources
- Implement response caching for frequently asked questions
- Add user feedback collection for continuous learning

Medium-term Enhancements:

- Upgrade to larger context window models (Llama-3.1-70b-versatile)
- Implement hybrid search combining keyword and semantic search
- Add document summarization for very long documents
- Implement conversation memory for multi-turn interactions

Long-term Strategic Improvements:

- Fine-tune embedding models on domain-specific data
- Implement retrieval quality scoring and automatic reranking
- Add multi-modal support (images, tables, charts)
- Develop automated quality assessment and monitoring pipeline