

Document Q&A System with RAG - Complete Hackathon Project

Project Overview

This comprehensive Document Question & Answer system uses **Retrieval-Augmented Generation (RAG)** to enable intelligent querying of uploaded documents. Built specifically for hackathons, it demonstrates modern AI/ML concepts while remaining accessible and extensible.

▮ Key Features

- **Multi-Format Document Support:** PDF, TXT, Markdown, HTML
- **Intelligent Text Processing:** Advanced chunking with configurable overlap
- **Vector-Based Retrieval:** TF-IDF embeddings with cosine similarity search
- **Natural Language Q&A:** Rule-based answer generation with LLM integration support
- **Interactive Web Interface:** Built with Streamlit for user-friendly experience
- **Source Citations:** Transparent answer provenance with confidence scores
- **Export Functionality:** Download Q&A sessions and chat history

▮ System Architecture

The system follows a modular RAG pipeline architecture:

```
Document Upload → Text Extraction → Chunking → Vectorization → Storage
                                                                ↓
User Interface ← Answer Generation ← Context Retrieval ← Query Processing
```

Core Components

1. **DocumentProcessor** (`document_processor.py`)
 - Multi-format text extraction (PDF, TXT, MD, HTML)
 - Intelligent text chunking with overlap
 - Text preprocessing and cleaning
2. **VectorStore** (`vector_store.py`)
 - TF-IDF based text vectorization
 - Cosine similarity search
 - In-memory document management
3. **RAGPipeline** (`rag_pipeline.py`)

- Question type detection and routing
- Context-aware answer generation
- Confidence scoring and source ranking

4. **Main Application** (main_app.py)

- Streamlit web interface
- Session state management
- File upload and validation

▮ Quick Start Guide

Prerequisites

- Python 3.8 or higher
- pip package manager

Installation Steps

1. Download the Project

```
# Extract the provided files to a directory
cd rag_qa_system
```

2. Install Dependencies

```
pip install -r requirements.txt
```

3. Run the Application

```
streamlit run main_app.py
```

4. Access the Interface

- Open browser to <http://localhost:8501>
- Upload documents using the sidebar
- Start asking questions!

Sample Usage Workflow

1. **Upload Documents:** Use sidebar to upload PDF, TXT, MD, or HTML files
2. **Process Documents:** Click "Process Documents" to chunk and vectorize
3. **Ask Questions:** Enter natural language questions in the main interface
4. **Review Answers:** View generated answers with source citations
5. **Export Results:** Download Q&A sessions as Markdown files

▮ Technical Implementation

Document Processing Pipeline

The system handles various document formats through specialized extraction methods:

PDF Processing:

- Uses `pypdf` or `PyPDF2` for text extraction
- Handles multi-page documents gracefully
- Supports both selectable and scanned text (with OCR recommendations)

Text Chunking Strategy:

- Configurable chunk sizes (200-800 tokens)
- Semantic-aware splitting at sentence boundaries
- Overlapping chunks to preserve context
- Metadata tracking for source attribution

Vector Storage:

- TF-IDF vectorization with stop word filtering
- Cosine similarity for retrieval ranking
- In-memory storage with efficient search algorithms

Question Answering Pipeline

The RAG pipeline employs sophisticated question processing:

Question Type Detection:

- Definition queries: "What is X?"
- Comparison queries: "Difference between X and Y?"
- List queries: "Types of X?"
- Process queries: "How does X work?"
- Explanation queries: "Why does X happen?"

Answer Generation:

- Rule-based extraction for reliable responses
- Context synthesis from multiple sources
- Confidence scoring based on retrieval quality
- Source citation with relevance scores

▮ User Interface Design

Web Application Features

Main Interface:

- Clean, intuitive design with clear navigation
- Real-time processing feedback
- Responsive layout for desktop and mobile

Document Management:

- Drag-and-drop file upload
- Document list with metadata
- Individual document deletion
- Processing progress indicators

Q&A Interface:

- Natural language question input
- Streaming answer display
- Expandable source citations
- Chat history management

Configuration Panel:

- Adjustable chunk size and overlap
- Retrieval parameter tuning
- Performance monitoring

▮ Advanced Features

LLM Integration Support

The system supports integration with various Language Models:

Supported Providers:

- OpenAI GPT (GPT-3.5, GPT-4)
- Anthropic Claude
- Hugging Face models
- Local models via Ollama

API Integration Example:

```
# OpenAI Integration
from api_integration import OpenAIIntegration
```

```
llm = OpenAIIntegration(api_key="your-api-key")
enhanced_pipeline = EnhancedRAGPipeline(
    vector_store=vector_store,
    llm_provider="openai",
    api_key="your-api-key"
)
```

Performance Optimization

Memory Management:

- Efficient vector storage algorithms
- Garbage collection for large documents
- Memory usage monitoring

Caching Layer:

- Query result caching
- Embedding caching for repeated documents
- Session state optimization

Quality Assessment

Answer Quality Metrics:

- Source diversity scoring
- Confidence level calculation
- Response completeness analysis
- User feedback integration

▮ Testing and Validation

Comprehensive Test Suite

The project includes extensive testing capabilities:

Unit Tests:

- Document processor validation
- Vector store functionality
- RAG pipeline components
- Utility function testing

Integration Tests:

- End-to-end workflow testing

- Multi-document scenarios
- Question type coverage
- Performance benchmarking

Sample Test Scenarios:

1. Basic Q&A with single document
2. Multi-document information synthesis
3. Edge cases and error handling
4. Performance under load

Sample Documents and Questions

Provided Test Data:

- Machine Learning basics (Markdown)
- Deep Learning guide (Markdown)
- RAG Systems overview (Text)

Test Question Categories:

- Definition queries
- Comparison questions
- Process explanations
- Technical details
- Cross-document synthesis

Deployment Options

Local Development

- Streamlit development server
- Virtual environment setup
- Hot reloading for development

Cloud Deployment

- **Streamlit Cloud:** Free hosting with GitHub integration
- **Heroku:** Container-based deployment
- **AWS EC2:** Full control deployment
- **Docker:** Containerized deployment

Production Considerations

- Database integration (PostgreSQL + pgvector)
- Load balancing and scaling
- Security enhancements
- Monitoring and logging

▮ Future Enhancements

Planned Features

- **Advanced Embeddings:** Transformer-based semantic search
- **Multi-language Support:** International document processing
- **Real-time Collaboration:** Shared document spaces
- **Analytics Dashboard:** Usage statistics and insights
- **Voice Interface:** Speech-to-text question input

Integration Opportunities

- **Knowledge Graphs:** Entity relationship mapping
- **OCR Enhancement:** Advanced document image processing
- **Workflow Automation:** API endpoints for programmatic access
- **Enterprise Features:** User management and access control

▮ Hackathon Tips

Presentation Strategy

1. **Demo Flow:** Start with document upload → show processing → ask compelling questions
2. **Technical Depth:** Explain RAG concept and implementation choices
3. **Scalability:** Discuss production deployment options
4. **Innovation:** Highlight unique features and optimizations

Key Selling Points

- **Complete Implementation:** Full RAG pipeline from scratch
- **User Experience:** Intuitive interface with real-time feedback
- **Extensibility:** Modular architecture for easy enhancement
- **Production Ready:** Deployment guides and testing suite

Common Questions and Answers

- **Q:** "How does this compare to ChatGPT?"
- **A:** "Our system provides transparent source citations and can be trained on private documents without data privacy concerns."
- **Q:** "What makes your chunking strategy special?"
- **A:** "We use semantic-aware splitting with configurable overlap, preserving context while optimizing retrieval accuracy."
- **Q:** "How scalable is this solution?"
- **A:** "The modular architecture supports various vector databases and can integrate with enterprise-grade infrastructure."

▮ Learning Resources

Key Concepts Demonstrated

- **Retrieval-Augmented Generation (RAG)**
- **Vector similarity search**
- **Text preprocessing and chunking**
- **Information retrieval systems**
- **Natural language processing**

Technologies Showcased

- **Python:** Core programming language
- **Streamlit:** Web application framework
- **NumPy:** Numerical computing
- **Natural Language Processing:** Text processing
- **Machine Learning:** Vector operations and similarity

▮ Contributing and Extension

Adding New Document Types

1. Extend `DocumentProcessor.supported_formats`
2. Implement extraction method
3. Update file validation
4. Add test cases

Customizing Answer Generation

1. Modify question type detection
2. Add new answer templates
3. Implement custom ranking algorithms
4. Test with diverse question types

Performance Optimization

1. Implement advanced vector databases
2. Add caching layers
3. Optimize memory usage
4. Profile and benchmark

▮ Project Structure

```
rag_qa_system/  
├── main_app.py           # Main Streamlit application  
├── document_processor.py # Document parsing and chunking  
├── vector_store.py       # Vector storage and similarity search  
├── rag_pipeline.py       # RAG pipeline and answer generation  
├── utils.py              # Utility functions  
├── api_integration.py     # LLM API integrations  
├── testing_examples.py    # Test cases and sample data  
├── deployment_guide.md   # Deployment instructions  
├── requirements.txt       # Python dependencies  
└── README.md             # Project documentation
```

▮ Conclusion

This Document Q&A System with RAG represents a comprehensive implementation of modern information retrieval and natural language processing techniques. It demonstrates practical AI/ML applications while maintaining accessibility for learning and extension.

The system is designed to impress judges with its technical depth, user experience, and practical applicability, while providing a solid foundation for further development and real-world deployment.

Perfect for hackathons, this project showcases:

- Deep understanding of RAG architectures
- Practical implementation skills
- User-centered design principles
- Production-ready development practices
- Clear documentation and testing

Ready to revolutionize document-based question answering! 🚀