# Technical Summary: DocuQuery System

DocuQuery is constructed upon a sophisticated, scalable architecture designed to deliver advanced document analysis and question-answering functionalities. The core technological stack comprises the following elements:

- **Backend Framework**: FastAPI was selected for its asynchronous capabilities, type safety, and robust performance metrics. This choice was pivotal for the efficient management of concurrent document processing and WebSocket connections.
- **Language Model Integration**: The system integrates OpenAI's GPT models via LangChain, thereby establishing a versatile abstraction layer. This integration facilitates the seamless interchange of models while maintaining a uniform API for document-based query resolution.
- **Vector Database**: FAISS (Facebook AI Similarity Search) was implemented to enable rapid and efficient similarity searches of document segments. This implementation substantially enhances retrieval efficacy when compared to conventional search methodologies.
- **Embedding Models**: Sentence transformers from HuggingFace were designated for document embedding, with "all-MiniLM-L6-v2" offering an optimal equilibrium between performance and precision for semantic search functionalities.
- **Document Processing**: PyMuPDF (Fitz) is utilized for PDF extraction, chosen for its speed and comprehensive feature set for document parsing.
- **Real-time Communication**: WebSockets are employed to facilitate streaming responses and real-time processing updates, thereby improving the user experience with immediate feedback.
- **Hybrid Retrieval**: The system employs both BM25 (keyword-based) and vector search, merging the strengths of both approaches to augment retrieval quality.
- **Document Chunking**: A bespoke FastChunker class was developed to execute semantically aware document segmentation, preserving context across segments.

Challenges and Resolutions

- **Document Structure Variability**: The structural diversity of PDF documents poses challenges to consistent extraction. This issue was addressed through the implementation of adaptive segmentation strategies that discern document structure (narrative versus structured) and consequently adjust parameters.
- **Memory Management**: Processing extensive documents with intricate NLP operations engendered memory strain. A session management system with TTL (time-to-live) and cleanup procedures was implemented to mitigate memory leaks and assure resource availability for subsequent requests.
- **Retrieval Quality**: Preliminary evaluations indicated that simple vector search was inadequate for precise retrieval. A multi-stage retrieval procedure was devised,

encompassing:
1. BM25 for keyword relevance.
2. Vector search for semantic similarity.
3. Cross-encoder re-ranking for definitive precision.
- **Context Window Limitations**: The context window restrictions of LLMs constrained the volume of document content that could be integrated into queries. A history-aware retriever was implemented that considers both document content and conversational history when selecting pertinent passages.
- **Metadata Extraction**: Automated identification of document metadata (author, title, date) proved arduous. A multi-stage method involving regular expressions and pattern matching was established to extract metadata from both PDF properties and content.
- **Query Understanding**: Ambiguous queries posed by users were common. Query expansion techniques and specialized handling for metadata questions were introduced to enhance the quality of responses.
- **Performance Optimization**: Initial document processing was protracted for large documents. Concurrent processing via ThreadPoolExecutor was applied to parallelize operations and significantly curtail processing durations.

Future Enhancements

Given additional time, enhancements to DocuQuery would encompass:

1. **Advanced Table Processing**: Improved extraction and querying capabilities for tabular data, which is presently suboptimally managed.
2. **Multi-Modal Support**: Extension of support to extract and analyze images and charts within documents.
3. **Fine-tuned Models**: Training of domain-specific models to refine performance in specialized domains such as legal, medical, or technical documentation.
4. **User Authentication**: Deployment of a robust user authentication system to ensure secure document management.
5. **Improved Caching**: Execution of a distributed caching layer utilizing Redis to boost performance and facilitate horizontal scaling.
6. **Document Version Control**: Incorporation of features to track document version changes and compare differentials.
7. **Integration with Document Management Systems**: Development of connectors for SharePoint, Google Drive, and other prevalent document repositories.
8. **Multilingual Support**: Extension of functionalities to accommodate non-English documents via language-specific models and processing.

Development Experience

The developmental process was both challenging and enriching. Salient aspects included:

1. Balancing precision and performance, especially within retrieval systems.
2. Formulating the document processing pipeline to robustly manage diverse document formats and structures.
3. Implementing real-time feedback through WebSockets for a responsive experience.

The principal challenge lay in constructing a system capable of judiciously identifying the correct context for query resolution, particularly when dealing with extensive documents where pertinent information may be disseminated across multiple sections.

In totality, this project demonstrated the transformative capabilities of modern NLP techniques on document interaction, advancing beyond mere search to encompass intelligent, context-aware question answering. The modular architectural design allows for ongoing enhancement in conjunction with the advancements in LLM and embedding technologies.

Here is the Github Link For the Repository
[Github](Github)