

S. No	File Name	Retriever Model	Generator Model	Retriever Issues	Generator Issues
1	RAG-LangChain-Base	<b>HuggingFace embeddings (hkunlp/instructor-base)</b> were used to create semantic embeddings from the PDF document, which were subsequently stored in a <b>FAISS vector database</b> to build the retriever model.	The generator model used is <b>fastchat-t5-3b-v1.0</b> , a 3-billion-parameter <b>T5</b> variant from <b>Hugging Face</b> , fine-tuned specifically for conversational tasks and instruction-following.	The text splitting strategy (chunk size and overlap) may split meaningful context, resulting in incomplete retrieval, causing a miss in crucial context needed for accurate responses	FastChat-T5 typically has a relatively small token limit (~512 to 1024 tokens). Large contexts or extensive prompt data might exceed this limit, causing essential context information to get truncated, causing incomplete or irrelevant responses.
2	RAG-LangChain-OpenAI	The retriever model utilizes <b>OpenAI embeddings</b> to generate semantic embeddings from the PDF document, which are stored in a <b>FAISS vector database</b> for efficient semantic retrieval.	The generator model used is <b>OpenAI GPT-4o-mini</b> , integrated with a predefined <b>LangChain prompt template</b> designed to guide the model's responses for accurate and context-aware generation.	OpenAI embeddings, while effective, might generalize concepts excessively, resulting in loss of subtle context distinctions, causing retrieval of superficially similar but contextually irrelevant content.	GPT-4o-mini, being a smaller variant, may have tighter context length constraints compared to standard GPT-4, causing longer contexts or extensive retrieval data to be truncated, affecting response quality.