# YouTube Chatbot using LangChain: A RAG System

**1. Core Problem Statement: YouTube Chatbot**

The primary problem addressed is the difficulty of quickly extracting information from lengthy YouTube videos (e.g., podcasts, lectures). The proposed solution is a RAG-based system that allows users to "chat" with any YouTube video in real-time.

- **Problem:** "अगर आपको सारा का सारा कंटेंट समझना है उस वीडियो का तो आपको पूरा वीडियो देखना पड़ेगा" (If you want to understand all the content of that video, you have to watch the entire video.)

- **Solution:** Users can ask questions like "क्या इस पॉडकास्ट में एआई के बारे में बात हो रही है?" (Is AI discussed in this podcast?) or "क्या आप इस पूरे वीडियो को समराइज कर सकते हो इन फाइव बुलेट पॉइंट्स?" (Can you summarize this entire video in five bullet points?). The RAG system will provide instant, relevant answers or summaries without requiring the user to watch the entire video.

**2. System Architecture: The RAG Pipeline**

The YouTube chatbot is built following a standard RAG architecture, which involves four main steps: Indexing, Retrieval, Augmentation, and Generation.

**2.1. Indexing (Data Preparation)**

This phase involves preparing the YouTube video content for efficient retrieval.

- **Step 1: Load YouTube Video Transcript.** The goal is to fetch the full transcript of a given YouTube video.

- While LangChain's YTLoader (YouTube Loader) exists, the tutorial opts to use YouTube's own API due to YTLoader being "थोड़ा सा बगी है" (a bit buggy) for some videos. The API returns transcripts with timestamps, which are then concatenated into a single, large string.

- Example: For Hindi videos, the language parameter in the API call needs to be set to 'hindi' instead of 'english'.

- **Step 2: Split Transcript into Chunks.** Since video transcripts can be very long, they are divided into smaller, manageable chunks.

- The RecursiveCharacterTextSplitter is used for this purpose.

- **chunk_size**: Set to 1000.

- **chunk_overlap**: Set to 200.

- Example: A 2-hour podcast transcript resulted in 168 chunks.

- **Step 3: Generate Embeddings and Store in a Vector Store.** Each chunk is converted into a numerical vector (embedding) using an embedding model.

- **Embedding Model:** OpenAIEmbeddings.

- **Vector Store:** FAISS (Facebook AI Similarity Search) is used for local storage. Each chunk gets a unique ID in the vector store.

**2.2. Retrieval**

This phase involves finding relevant information from the indexed data based on the user's query.

1. **Process:**A Retriever is formed (initially a simple similarity_search based retriever using the FAISS vector store).

2. The user's query is sent to the retriever.

3. The retriever embeds the query, performs a semantic search in the vector store to find the "closest" vectors, and returns the corresponding relevant document chunks.

- **"रिट्रीवर क्या करेगा कि इस क्वेरी को एंबेड करेगा वेक्टर के फॉर्म में लेके आएगा और फिर इस वेक्टर स्टोर में जाकर के सर्च करेगा कि इस गिवन वेक्टर के क्लोजेस्ट कौन से वेक्टर्स हैं इस वेक्टर स्टोर में और उसके कॉरेस्पोंडिंग जो भी चंक्स हैं जो भी डॉक्यूमेंट्स हैं उनको लाके आप तक देगा"** (The retriever will embed this query, bring it in vector form, and then search in this vector store for the closest vectors to this given vector, and bring you the corresponding chunks or documents.)

- The retriever is configured to return the top 4 most similar documents.

**2.3. Augmentation**

This phase combines the user's query with the retrieved context to form a comprehensive prompt for the Language Model (LLM).

1. **Process:**The retrieved documents (list of strings/page content) are concatenated into a single large context string.

2. A PromptTemplate is created with placeholders for context and question.

3. **Prompt Template:** "यू आर अ हेल्पफुल असिस्टेंट आंसर ओनली फ्रॉम द प्रोवाइडेड ट्रांसक्रिप्ट कॉन्टेक्स्ट इफ द कॉन्टेक्स्ट इज़ इनसफिशिएंट जस्ट से यू डोंट नो" (You are a helpful assistant. Answer only from the provided transcript context. If the context is insufficient, just say you don't know.)The concatenated context and the original question are inserted into the prompt template to create the "final prompt".

**2.4. Generation**

The final phase where the LLM processes the augmented prompt and generates the answer.

- **LLM:** OpenAI's models are used (OpenAI class is instantiated).

- **Process:** The final_prompt is sent to the LLM. The LLM generates a response based on the provided context and question.

- **"हमारा एलएलएम क्वेरी को समझेगा इस कॉन्टेक्स्ट को समझेगा और पलट के मुझे एक रिस्पांस जनरेट करके देगा"** (Our LLM will understand the query, understand this context, and generate a response for me.)

- The generated answer can then be extracted (e.g., response.content).

**3. Building a LangChain Chain for Simplified Execution**

Initially, each RAG step (retrieval, augmentation, generation) is called manually. To streamline the process, a LangChain Chain is constructed.

- **Problem:** "ये सारे स्टेप्स अलग-अलग तरीके से काम कर रहे हैं मतलब हमें हर स्टेप को मैनुअली कॉल करना पड़ रहा है" (All these steps are working separately, meaning we have to manually call each step.)

- **Solution:** A composite chain is built using RunnableParallel, RunnablePassthrough, and RunnableLambda.

- **Parallel Chain:** Handles the retrieval and context formatting.

- It takes the user's question as input.

- It uses the retriever to get documents, and then RunnableLambda with a custom format_docs function to concatenate these into a context string.

- It also passes through the original question.

- Output: A dictionary with context and question keys.

- **Main Chain:** Combines the prompt, LLM, and parser.

- **prompt**: The PromptTemplate defined earlier.

- **llm**: The instantiated LLM.

- **parser**: StringOutputParser to get the output as a clean string.

- The parallel_chain is then connected to the main_chain using the | operator, ensuring seamless data flow where the output of the parallel chain becomes the input for the main chain's prompt.

- **Result:** A single main_chain.invoke() call with the user's question now executes the entire RAG pipeline automatically, making the code "much cleaner" and easier to manage.

## 4. Potential Improvements and Advanced RAG Concepts

Implemented system is a "basic level RAG system" and discusses numerous ways to improve it to an "industry-grade level." This introduces the concept of **Advanced RAG**.

- **"इंडस्ट्री में जो रैक सिस्टम्स इंप्लीमेंट किए जाते हैं वो काफी कॉम्प्लेक्स होते हैं और उनमें बहुत तरह की टेक्निक्स यूज़ की जाती हैं"** (The RAG systems implemented in the industry are quite complex and use many types of techniques.)

### 4.1. UI-Based Enhancements

- **Streamlit Website:** Allow users to paste a YouTube URL and chat within a web interface.

- **Chrome Plugin:** The ideal solution, providing a chat interface directly on the YouTube video page. Requires HTML, CSS, JavaScript knowledge.

### 4.2. Evaluation

Crucial for assessing system performance and guiding improvements.

- **Libraries:** Ragas (popular library for RAG evaluation) and Langsmith (for tracing and debugging RAG pipelines).

- **Metrics (from Ragas):Faithfulness:** Is the generated answer truly supported by the context?

- **Answer Relevancy:** Is the answer relevant to the question?

- **Context Precision:** How useful was the retrieved context in answering the question?

- **Context Recall:** Was all useful information from the vector store retrieved?

### 4.3. Indexing Improvements

- **Document Ingestion:Error Correction:** Fixing errors in auto-generated transcripts.

- **Translation:** Translating transcripts into a desired language (e.g., English) if the original is in a different language (like Hindi).

- **Text Splitting:** Using Semantic Chunker instead of RecursiveCharacterTextSplitter to avoid splitting paragraphs mid-sentence and maintain semantic coherence.

- **Vector Store:** Using cloud-based solutions like Pinecone for production environments instead of basic FAISS.

## 4.4. Retrieval Improvements

- **Pre-Retrieval:Query Rewriting:** Using an LLM to improve vague or ambiguous user queries.

- **Multi-Query Generation:** Generating multiple queries from a single original query to capture different perspectives and improve retrieval.

- **Domain-Aware Routing:** For complex systems with multiple retrievers, routing queries to the most appropriate retriever based on the query's domain.

- **During Retrieval:MMR (Maximal Marginal Relevance) Search:** Retrieves diverse and relevant results.

- **Hybrid Retrieval:** Combining semantic search with keyword search.

- **Post-Retrieval:Re-ranking:** Using an LLM to re-rank the retrieved documents based on their relevance to the query, improving the order of context provided to the LLM.

- **Contextual Compression:** Trimming unnecessary parts of retrieved documents to only keep the most meaningful information, reducing token usage for the LLM.

## 4.5. Augmentation Improvements

- **Prompt Templating:** More sophisticated prompt engineering, including providing examples to the LLM.

- **Answer Grounding:** Explicitly instructing the LLM to generate answers *only* from the provided context and avoid hallucination or fabricating facts.

- **Context Window Optimization:** Trimming context dynamically to ensure it fits within the LLM's token limit, preventing errors and optimizing cost.

## 4.6. Generation Improvements

- **Answer with Citation:** Asking the LLM to cite the source (specific part of the context) from which it derived the answer.

- **Guard Railing:** Implementing mechanisms to prevent the LLM from generating inappropriate, harmful, or incorrect outputs.

## 4.7. System Design Improvements

- **Multi-Modal RAG:** Extending the RAG system to process and respond to inputs beyond text, such as images and videos.

- **Agentic RAG:** Building an AI agent that can perform additional actions (e.g., web browsing) beyond just answering questions to gather more context.

- **Memory-Based RAG:** Incorporating memory into the RAG system to allow for personalized conversations and recall past interactions with the user.

## 5. Future Outlook

The discussed RAG system is a foundational example. There is a "Advanced RAG" that have more complex techniques and challenges in building industrial-grade RAG systems.

- **"अभी इंडस्ट्री में एक कंप्लीटली नई फील्ड बन के उभर के आई है जिसका नाम है एडवांस्ड रैग"** (Currently, a completely new field has emerged in the industry, named Advanced RAG.)