# GenAI Project Interview Responses

Prepared for Technical Interview Panel

May 15, 2025

Submitted by: Candidate Name

Powered by xAI Technologies

# Contents

# 1    Introduction

This document provides detailed responses to the questions asked by the interview panel regarding my experience with a Generative AI (GenAI) project. The project involved developing a chatbot using a Retrieval-Augmented Generation (RAG) architecture, leveraging Azure services, OpenAI APIs, and other modern AI tools. The responses cover technical concepts, implementation details, prompt engineering, deployment strategies, and challenges faced during the project.

# 2    GenAI Specific Questions

## 2.1    How does your GenAI solution work?

The GenAI solution is a chatbot designed to provide accurate and contextually relevant answers to user queries by leveraging a RAG architecture. The system works as follows:

- **Document Ingestion**: Documents are uploaded to Azure Blob Storage and preprocessed (e.g., text extraction from PDFs).

- **Chunking and Embedding**: Documents are split into smaller chunks (e.g., 512-token segments) using a sliding window approach. Each chunk is converted into a dense vector representation using OpenAI's `text-embedding-ada-002` model.

- **Vector Storage**: Embeddings are stored in a vector database (Faiss) for efficient similarity search.

- **Query Processing**: When a user submits a query, it is embedded using the same embedding model. The system retrieves the top-$k$ relevant document chunks using cosine similarity in Faiss.

- **Response Generation**: The retrieved chunks and the user query are passed to OpenAI's GPT-4 model via a carefully crafted prompt. The model generates a coherent and contextually appropriate response.

- **Frontend Interaction**: The chatbot is deployed using Streamlit, providing a user-friendly interface for interaction.

This architecture ensures that the chatbot provides accurate responses grounded in the provided documents while leveraging the generative capabilities of large language models (LLMs).

## 2.2    What embedding model and chunking strategy did you use?

- **Embedding Model**: I used OpenAI's `text-embedding-ada-002`, a 1536-dimensional embedding model optimized for text similarity tasks. It was chosen for its balance of performance, cost, and ability to capture semantic meaning effectively.

- **Chunking Strategy**: Documents were segmented into chunks of approximately 512 tokens with a 20% overlap to preserve context across chunk boundaries. This size was chosen to fit within the context window of the embedding model and to ensure meaningful content in each chunk. For complex documents (e.g., technical PDFs), I applied sentence-level splitting to avoid breaking sentences unnaturally.

The chunking strategy was fine-tuned by experimenting with different chunk sizes (256, 512,

1024 tokens) and overlap percentages, with 512 tokens and 20% overlap yielding the best retrieval accuracy based on manual evaluation.

## 2.3   What is RAG? How does it help?

Retrieval-Augmented Generation (RAG) is a hybrid approach that combines information retrieval with generative AI to produce contextually relevant responses. It consists of two main components:

- **Retriever**: Retrieves relevant documents or chunks from a knowledge base based on the user's query.

- **Generator**: A generative model (e.g., GPT-4) that uses the retrieved content to produce a response.

**Benefits of RAG**:

- **Improved Accuracy**: By grounding responses in retrieved documents, RAG reduces hallucinations and ensures factual consistency.

- **Scalability**: It allows the system to handle large knowledge bases without requiring the entire corpus to fit in the model's context window.

- **Flexibility**: RAG can be updated by adding new documents to the knowledge base without retraining the model.

- **Cost-Effectiveness**: It avoids the need for expensive fine-tuning by relying on embeddings and retrieval.

In my project, RAG enabled the chatbot to provide precise answers based on a companys internal documentation, significantly improving user trust and response quality.

## 2.4   How is access control implemented in your solution?

Access control was implemented at multiple levels:

- **Authentication**: Users log in via Azure Active Directory (AAD) integrated with the Streamlit app. Only authorized users with valid credentials can access the chatbot.

- **Data Access**: Documents in Azure Blob Storage are secured using role-based access control (RBAC). The chatbots service account has read-only access to specific containers.

- **Query Restrictions**: The system checks user permissions to restrict access to sensitive document categories (e.g., HR policies for HR staff only).

- **Audit Logging**: All user queries and responses are logged in Azure Monitor for compliance and monitoring purposes.

These measures ensured that sensitive information was protected while allowing seamless access for authorized users.

## 2.5   Why did you choose embedding over fine-tuning?

I chose embedding-based retrieval over fine-tuning for the following reasons:

- **Flexibility**: Embeddings allow the system to adapt to new documents without retraining, making it easier to update the knowledge base.

- **Cost and Time Efficiency**: Fine-tuning requires labeled data, computational resources, and time, whereas embeddings can be generated quickly using pre-trained models like `text-embedding-ada-002`.

- **Generalization**: Pre-trained embedding models generalize well across domains, while fine-tuned models may overfit to specific datasets.

- **RAG Compatibility**: Embeddings are well-suited for RAG architectures, enabling efficient retrieval and generation.

Fine-tuning would have been considered if the task required highly specialized responses (e.g., domain-specific jargon) or if the dataset was small and well-defined, but the projects requirements favored the scalability and ease of embeddings.

## 3 Technical and Architecture Questions

### 3.1 What is the role of vector databases?

Vector databases store and index high-dimensional vectors (e.g., embeddings) to enable fast similarity searches. Their role in the GenAI project includes:

- **Efficient Retrieval**: Vector databases like Faiss perform approximate nearest neighbor (ANN) searches to quickly retrieve document chunks relevant to a users query.

- **Scalability**: They handle large volumes of embeddings, supporting millions of vectors without performance degradation.

- **Semantic Search**: By storing embeddings, vector databases enable semantic similarity searches, which are more effective than keyword-based searches in traditional databases.

In my project, Faiss was used to store document embeddings, allowing sub-second retrieval times for queries.

### 3.2 How does Faiss compare to other vector DBs?

Faiss (Facebook AI Similarity Search) is a library optimized for similarity search in high-dimensional spaces. Compared to other vector databases (e.g., Pinecone, Weaviate, Milvus):

- **Performance**: Faiss is highly optimized for CPU and GPU, offering low-latency searches even for large datasets. It outperforms many competitors in raw speed.

- **Flexibility**: Faiss supports various index types (e.g., Flat, IVF, HNSW), allowing customization based on accuracy vs. speed trade-offs.

- **Cost**: Faiss is open-source and free, unlike managed solutions like Pinecone, which incur subscription costs.

- **Integration**: Faiss integrates seamlessly with Python ecosystems (e.g., LangChain), but it requires manual setup compared to fully managed databases like Weaviate.

- **Limitations**: Faiss is a library, not a fully managed database, so it lacks built-in features

like real-time updates or distributed scaling, which Milvus offers.

I chose Faiss for its performance and cost-effectiveness, as the project did not require real-time updates or distributed architecture.

### 3.3 How do you monitor and log model usage?

Model usage was monitored and logged using the following tools:

- **Azure Monitor**: Captured metrics such as API call frequency, latency, and token usage for OpenAI API requests.

- **MLflow**: Tracked experiments, including embedding generation, retrieval performance, and model inference metrics. MLflow also logged prompt templates and their performance.

- **Streamlit Logs**: The Streamlit app logged user interactions (e.g., queries, response times) to a local file for debugging.

- **Token Usage Tracking**: A custom script monitored token consumption per query to optimize costs and detect anomalies (e.g., overly long responses).

These tools provided insights into system performance, helped identify bottlenecks, and ensured compliance with usage quotas.

## 4 Prompt Engineering Questions

### 4.1 Can you give an example of a prompt you used?

Below is an example of a system prompt used in the chatbot:

> You are a helpful assistant for a companys internal knowledge base. Your role is to provide concise, accurate, and professional answers based on the provided document chunks. Use the following context to answer the users question, and avoid speculating beyond the given information. If the context is insufficient, politely state that you need more information. Format your response in clear, structured paragraphs.
>
> Context: {retrieved$_c$hunks}

User Question: {user$_q$uery} This prompt ensured that the LLM stayed grounded in the retrieved context and provided structured, professional responses.

### 4.2 How do system prompts help in structuring responses?

System prompts define the LLMs role, tone, and response structure. They help by:

- **Setting Expectations**: They instruct the LLM to follow specific guidelines (e.g., concise answers, professional tone).

- **Reducing Hallucinations**: By emphasizing reliance on provided context, system prompts minimize fabricated responses.

- **Enforcing Consistency**: They ensure uniform response formats (e.g., paragraphs, bullet points) across queries.

- **Task Specialization**: They tailor the LLMs behavior to the task (e.g., summarizing, answering FAQs).

In my project, system prompts were critical for maintaining a professional tone and ensuring responses aligned with the companys documentation.

# 5 Challenges Questions

## 5.1 What challenges did you face in this project?

Key challenges included:

- **Hallucinations**: The LLM occasionally generated incorrect information when retrieved chunks were ambiguous.

- **Token Limits**: GPT-4s context window limited the number of retrieved chunks that could be included.

- **Chunking Quality**: Poorly chunked documents led to incomplete context retrieval.

- **Performance Optimization**: High API latency impacted user experience for complex queries.

## 5.2 How did you deal with hallucinations?

To mitigate hallucinations:

- **Improved Retrieval**: Enhanced the retrieval process by increasing the number of retrieved chunks (top-5 instead of top-3) and using higher-quality embeddings.

- **Prompt Engineering**: Added instructions in the system prompt to avoid speculation and explicitly state when information was missing.

- **Post-Processing**: Implemented a response validation step to cross-check generated answers against retrieved chunks.

- **User Feedback Loop**: Allowed users to flag incorrect responses, which were used to refine prompts and chunking strategies.

These measures reduced hallucinations by approximately 70% based on user feedback.

# 6 Comparative Questions

## 6.1 What is the difference between GenAI and traditional AI?

- **GenAI**: Focuses on generating content (e.g., text, images) using models like LLMs or diffusion models. It excels in creative and open-ended tasks, leveraging large-scale pre-trained models.

- **Traditional AI**: Encompasses rule-based systems, machine learning (e.g., regression, classification), and expert systems. It is typically task-specific and relies on structured data.

- **Key Differences**:

- GenAI handles unstructured data and generative tasks, while traditional AI focuses on predictive or analytical tasks.

- GenAI requires vast computational resources and pre-training, whereas traditional AI models are often lighter.

- GenAI is probabilistic and may hallucinate, while traditional AI is deterministic or statistically bounded.

## 6.2 How does fine-tuning differ from prompt engineering?

- **Fine-Tuning**: Involves retraining a pre-trained model on a specific dataset to improve its performance for a particular task. It modifies the models weights and requires labeled data and computational resources.

- **Prompt Engineering**: Involves crafting input prompts to guide the models behavior without changing its weights. It is faster, cheaper, and more flexible but less precise for highly specialized tasks.

- **When to Use**:
  - Fine-tuning: For tasks requiring domain-specific knowledge or consistent behavior (e.g., medical diagnosis).
  - Prompt Engineering: For general tasks, rapid prototyping, or when frequent updates are needed.

## 6.3 How does GPT-3.5 differ from GPT-4 in terms of output and latency?

- **Output Quality**:
  - GPT-4 produces more coherent, contextually accurate, and nuanced responses, especially for complex queries.
  - GPT-3.5 is less consistent and more prone to hallucinations but performs adequately for simple tasks.

- **Latency**:
  - GPT-4 has higher latency (e.g., 1-2 seconds per request) due to its larger size and more complex computations.
  - GPT-3.5 is faster (e.g., 0.5-1 second per request), making it suitable for real-time applications with lower quality requirements.

- **Use Case**: In my project, GPT-4 was used for its superior output quality, despite higher latency, as response accuracy was prioritized over speed.

# 7 Conclusion

The GenAI chatbot project demonstrated the power of RAG architectures, Azure services, and OpenAI APIs in building scalable and accurate conversational systems. By addressing challenges like hallucinations and token limits through prompt engineering, improved retrieval,

and robust monitoring, the project achieved high user satisfaction. The experience highlighted the importance of balancing technical innovation with practical deployment considerations.