

# L2 AI Engineer Interview Preparation Guide

Preparing for a Level 2 Technical Interview in Generative AI

May 15, 2025

Prepared by: Candidate Name

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Overview of L2 Discussion for AI Engineer Role</b>	<b>2</b>
2.1	What is an L2 Discussion?	2
2.2	Key Areas of Focus	2
<b>3</b>	<b>Expected Questions and Sample Responses</b>	<b>3</b>
3.1	GenAI Project Deep Dive	3
3.1.1	Walk us through the end-to-end architecture of your GenAI chatbot.	3
3.1.2	What chunking strategy and embedding model did you use, and why?	3
3.1.3	What is RAG, and how does it help in your project?	4
3.1.4	How is access control implemented in your solution?	4
3.1.5	Why did you choose embeddings over fine-tuning?	4
3.2	Technical Concepts and Implementation	5
3.2.1	What is the role of vector databases, and how does Faiss compare to others?	5
3.2.2	How do embeddings differ from fine-tuning, and when would you use each?	5
3.2.3	How did you use MLflow in your LLMOps pipeline?	6
3.2.4	How do you handle hallucinations in LLMs?	6
3.3	Prompt Engineering	6
3.3.1	Share a specific prompt you used and explain its effectiveness.	6
3.3.2	How do system prompts help in structuring responses?	6
3.4	System Design and Scalability	7
3.4.1	Design a scalable RAG-based chatbot system.	7
3.4.2	How would you optimize latency for real-time responses?	7
3.5	Comparative Questions	8
3.5.1	What distinguishes GenAI from traditional AI?	8
3.5.2	How does GPT-3.5 compare to GPT-4?	8
3.5.3	How does fine-tuning differ from prompt engineering?	8
3.6	Behavioral Questions	8
3.6.1	Describe a time you faced a technical challenge.	8
3.6.2	How do you stay updated on AI advancements?	9
<b>4</b>	<b>Preparation Strategies</b>	<b>9</b>
4.1	Review Your Project	9
4.2	Study Key Concepts	9
4.3	Practice System Design	9
4.4	Code Practice	9
4.5	Mock Interviews	9
4.6	Ask Questions	9
<b>5</b>	<b>Sample Practice Questions</b>	<b>9</b>
<b>6</b>	<b>Conclusion</b>	<b>10</b>

# 1 Introduction

This document serves as a comprehensive guide to prepare for a Level 2 (L2) technical interview for an AI Engineer role, with a focus on Generative AI (GenAI). The L2 discussion is typically the second round of technical interviews, designed to assess your technical depth, problem-solving skills, and practical experience in building AI systems. Drawing from your experience with a GenAI chatbot project using Retrieval-Augmented Generation (RAG), Azure services, and OpenAI APIs, this guide outlines what to expect, key topics, sample questions, and preparation strategies. It includes detailed sample responses to anticipated questions, ensuring you are well-equipped to demonstrate your expertise.

## 2 Overview of L2 Discussion for AI Engineer Role

### 2.1 What is an L2 Discussion?

The L2 interview is a deeper evaluation compared to the initial screening (L1). It focuses on:

- **Technical Depth:** Understanding of AI/ML concepts, particularly GenAI, embeddings, and LLMOps.
- **Practical Experience:** Ability to apply knowledge to real-world projects, such as your RAG-based chatbot.
- **Problem-Solving:** Designing systems, optimizing performance, and addressing challenges.
- **Communication:** Explaining complex concepts clearly to technical and non-technical audiences.

The panel typically includes senior engineers, data scientists, or engineering managers and lasts 45–60 minutes.

### 2.2 Key Areas of Focus

Based on your GenAI project and the topics provided, expect the following areas:

- **Project Deep Dive:** Detailed questions about your chatbots RAG architecture, Azure integration, and OpenAI API usage.
- **Technical Concepts:** Vector databases, embeddings vs. fine-tuning, LLMOps, and prompt engineering.
- **System Design:** Designing scalable AI systems and optimizing performance.
- **Coding/Problem-Solving:** Pseudocode or debugging tasks related to retrieval or text processing.
- **Comparative Questions:** Comparing GenAI vs. traditional AI, GPT-3.5 vs. GPT-4, etc.
- **Behavioral Questions:** Teamwork, handling challenges, and staying updated with AI trends.

## 3 Expected Questions and Sample Responses

### 3.1 GenAI Project Deep Dive

#### 3.1.1 Walk us through the end-to-end architecture of your GenAI chatbot.

**Response:** My GenAI chatbot is built using a Retrieval-Augmented Generation (RAG) architecture to provide accurate, contextually relevant responses based on a companys internal documentation. The pipeline includes:

- **Document Ingestion:** Documents (e.g., PDFs, Word files) are uploaded to Azure Blob Storage. Text is extracted using Azures Document Intelligence service.
- **Chunking and Embedding:** Documents are split into 512-token chunks with 20% overlap to preserve context. Each chunk is embedded using OpenAIs `text-embedding-ada-002` model, producing 1536-dimensional vectors.
- **Vector Storage:** Embeddings are indexed in Faiss using the HNSW (Hierarchical Navigable Small World) index for fast similarity search.
- **Query Processing:** User queries are embedded and compared against the Faiss index to retrieve the top-5 relevant chunks via cosine similarity.
- **Response Generation:** Retrieved chunks and the query are passed to GPT-4 via a system prompt that enforces concise, professional responses. The prompt includes instructions to avoid speculation.
- **Frontend:** The chatbot is deployed using Streamlit, integrated with Azure Active Directory (AAD) for authentication.

This architecture balances accuracy, scalability, and cost, leveraging Azures secure storage and OpenAIs robust models.

#### 3.1.2 What chunking strategy and embedding model did you use, and why?

**Response:** I used a chunking strategy that splits documents into 512-token segments with a 20% overlap. This size was chosen to fit within the context window of the `text-embedding-ada-002` model and to ensure each chunk contained meaningful context. The overlap prevented loss of context across chunk boundaries, especially for complex documents like technical manuals. I experimented with chunk sizes (256, 512, 1024 tokens) and found 512 tokens optimal based on retrieval accuracy.

For embeddings, I used OpenAIs `text-embedding-ada-002` due to its high performance in semantic similarity tasks, cost-effectiveness, and 1536-dimensional output, which captures nuanced meanings. Compared to alternatives like Sentence-BERT, it offered better trade-offs between quality and API costs. The choice was validated through A/B testing, where `ada-002` outperformed BERT-based models by 15% in retrieval precision.

### 3.1.3 What is RAG, and how does it help in your project?

**Response:** Retrieval-Augmented Generation (RAG) combines information retrieval with generative AI to produce contextually accurate responses. It has two components:

- **Retriever:** Uses a vector database (Faiss in my project) to fetch relevant document chunks based on query embeddings.
- **Generator:** A large language model (GPT-4) generates responses using the retrieved chunks and the query.

RAG helped my project by:

- **Reducing Hallucinations:** Grounding responses in retrieved documents ensured factual accuracy.
- **Scaling Knowledge Bases:** It handled large document sets without exceeding GPT-4s context window.
- **Cost Efficiency:** Avoided fine-tuning by relying on embeddings and retrieval.
- **Flexibility:** Allowed easy updates to the knowledge base by adding new documents.

For example, when users asked about company policies, RAG retrieved exact policy excerpts, enabling precise answers.

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

**Response:** Access control was implemented at multiple layers:

- **User Authentication:** Streamlit integrated with Azure Active Directory (AAD) to ensure only authorized users could access the chatbot.
- **Data Security:** Documents in Azure Blob Storage were secured with role-based access control (RBAC), granting the chatbots service account read-only access to specific containers.
- **Query Restrictions:** The system checked user roles to limit access to sensitive documents (e.g., HR policies for HR staff only).
- **Audit Logging:** Azure Monitor logged all queries and responses for compliance and debugging.

These measures protected sensitive data while ensuring a seamless user experience.

### 3.1.5 Why did you choose embeddings over fine-tuning?

**Response:** I opted for embeddings over fine-tuning because:

- **Scalability:** Embeddings allowed adding new documents without retraining, unlike fine-tuning, which requires dataset updates.
- **Cost and Time:** Generating embeddings with `text-embedding-ada-002` was faster and cheaper than fine-tuning, which needs labeled data and GPU resources.
- **Generalization:** Pre-trained embeddings generalized well across domains, while fine-

tuning risked overfitting.

- **RAG Compatibility:** Embeddings enabled efficient retrieval in the RAG pipeline.

Fine-tuning would be preferred for highly specialized tasks (e.g., legal document analysis), but embeddings suited the projects need for flexibility and rapid deployment.

## 3.2 Technical Concepts and Implementation

### 3.2.1 What is the role of vector databases, and how does Faiss compare to others?

**Response:** Vector databases store high-dimensional embeddings for fast similarity searches, critical for RAG systems. Their roles include:

- **Fast Retrieval:** Enabling sub-second query responses via approximate nearest neighbor (ANN) search.
- **Scalability:** Handling millions of vectors efficiently.
- **Semantic Search:** Supporting cosine similarity for context-aware retrieval.

Faiss, compared to Pinecone, Weaviate, or Milvus:

- **Pros:** Open-source, CPU/GPU-optimized, supports flexible index types (e.g., HNSW, IVF). It achieved 0.1-second retrieval times in my project.
- **Cons:** Lacks built-in real-time updates or distributed scaling, unlike Milvus.
- **Why Chosen:** Its performance and zero cost suited the projects needs, and integration with LangChain simplified development.

### 3.2.2 How do embeddings differ from fine-tuning, and when would you use each?

**Response:**

- **Embeddings:** Convert text into fixed-length vectors capturing semantic meaning, used for retrieval or similarity tasks. Theyre generated by pre-trained models and require no training.
- **Fine-Tuning:** Retrains a models weights on a task-specific dataset to improve performance, altering its behavior.

**Differences:**

- Embeddings are faster and cheaper but less precise for niche tasks.
- Fine-tuning requires labeled data and is resource-intensive but excels in specialized domains.

**When to Use:**

- Embeddings: For RAG, semantic search, or dynamic knowledge bases (as in my project).
- Fine-Tuning: For tasks like sentiment analysis in a specific industry or generating domain-specific jargon.

### 3.2.3 How did you use MLflow in your LLMOps pipeline?

**Response:** MLflow was integral to my LLMOps pipeline:

- **Experiment Tracking:** Logged chunk sizes, embedding models, and prompt variations with metrics like retrieval accuracy.
- **Model Management:** Stored and versioned embedding models and prompt templates.
- **Performance Monitoring:** Tracked API latency and token usage to optimize costs.
- **Reproducibility:** Ensured experiments were reproducible by saving configurations.

For example, MLflow helped identify that a 512-token chunk size outperformed 1024 tokens, guiding the final configuration.

### 3.2.4 How do you handle hallucinations in LLMs?

**Response:** Hallucinations occur when LLMs generate incorrect or fabricated information. I mitigated them by:

- **Enhanced Retrieval:** Increased retrieved chunks from 3 to 5 for richer context.
- **Prompt Engineering:** Added instructions like Only use provided context and State if information is missing.
- **Validation:** Post-processed responses to cross-check against retrieved chunks.
- **User Feedback:** Incorporated user-flagged errors to refine prompts, reducing hallucinations by 70%.

## 3.3 Prompt Engineering

### 3.3.1 Share a specific prompt you used and explain its effectiveness.

**Response:** Heres a system prompt used:

You are a professional assistant for a companys knowledge base. Provide concise, accurate answers based on the provided document chunks. Avoid speculation, and if context is insufficient, say so politely. Use clear paragraphs.

Context: {retrieved<sub>chunks</sub>}

Question: {user<sub>query</sub>} **Effectiveness:**

- **Clarity:** Defined the assistants role and tone, ensuring professional responses.
- **Accuracy:** Emphasized reliance on context, reducing hallucinations.
- **Structure:** Enforced paragraph formatting for readability.
- **Fallback:** Handled missing information gracefully, improving user trust.

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

**Response:** System prompts guide LLM behavior by:

- **Defining Role:** Setting tone and purpose (e.g., professional assistant).
- **Reducing Errors:** Limiting responses to provided context to avoid fabrication.
- **Ensuring Consistency:** Standardizing response formats (e.g., paragraphs, lists).
- **Task Focus:** Tailoring output to specific tasks like summarizing or answering FAQs.

In my project, system prompts ensured uniform, accurate responses aligned with company standards.

### 3.4 System Design and Scalability

#### 3.4.1 Design a scalable RAG-based chatbot system.

**Response:** A scalable RAG-based chatbot system includes:

- **Data Ingestion:** Azure Blob Storage for document storage, with Azure Document Intelligence for text extraction.
- **Embedding Pipeline:** Batch processing to generate embeddings, stored in Faiss on Azure VM clusters.
- **Retrieval:** Faiss with HNSW index, scaled via Azure Kubernetes Service (AKS) for high traffic.
- **Generation:** OpenAI GPT-4, with caching (Redis) for frequent queries.
- **Frontend:** Streamlit on Azure App Service, with auto-scaling based on traffic.
- **Monitoring:** Azure Monitor for latency and token usage, with alerts for anomalies.

**Scalability Features:**

- Horizontal scaling of Faiss and Streamlit instances.
- Load balancing via Azure Application Gateway.
- Asynchronous embedding updates to avoid downtime.

#### 3.4.2 How would you optimize latency for real-time responses?

**Response:** To optimize latency:

- **Caching:** Store frequent query embeddings in Redis to skip retrieval.
- **Index Optimization:** Use Faiss HNSW for faster searches (0.1s vs. 0.5s for Flat index).
- **Parallel Processing:** Run retrieval and generation concurrently using asyncio in Python.
- **Model Selection:** Use GPT-3.5 for low-complexity queries to reduce API latency.
- **Edge Deployment:** Deploy Streamlit on Azure CDN for faster UI loading.

These reduced average response time from 2 seconds to 0.8 seconds.



## 3.5 Comparative Questions

### 3.5.1 What distinguishes GenAI from traditional AI?

**Response:**

- **GenAI:** Generates content (text, images) using models like LLMs or diffusion models. It excels in creative tasks but may hallucinate.
- **Traditional AI:** Includes rule-based systems and ML (e.g., regression, SVM) for predictive tasks. Its deterministic and data-driven.

**Examples:**

- GenAI: Chatbots, image generators.
- Traditional AI: Fraud detection, recommendation systems.

### 3.5.2 How does GPT-3.5 compare to GPT-4?

**Response:**

- **Output Quality:** GPT-4 is more coherent and accurate, especially for complex queries, with fewer hallucinations.
- **Latency:** GPT-3.5 is faster (0.5–1s vs. 1–2s for GPT-4) due to its smaller size.
- **Cost:** GPT-3.5 is cheaper per token.

In my project, GPT-4 was used for its superior accuracy, as response quality was prioritized.

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

**Response:**

- **Fine-Tuning:** Retrains model weights on specific data, improving task-specific performance but requiring time and resources.
- **Prompt Engineering:** Crafts input prompts to guide behavior without retraining, offering flexibility but less precision.

**Use Cases:**

- Fine-Tuning: Domain-specific tasks (e.g., medical diagnosis).
- Prompt Engineering: General tasks or rapid prototyping.

## 3.6 Behavioral Questions

### 3.6.1 Describe a time you faced a technical challenge.

**Response:** In my chatbot project, hallucinations were a challenge when retrieved chunks lacked context. I addressed this by increasing retrieved chunks to 5, refining prompts to avoid speculation, and validating responses against context. This reduced hallucinations by 70%, improving user trust.

### 3.6.2 How do you stay updated on AI advancements?

**Response:** I follow ArXiv for research papers, read blogs like Towards Data Science, and monitor X for GenAI trends (e.g., multimodal RAG). I also attend conferences like NeurIPS and experiment with new tools via side projects.

## 4 Preparation Strategies

### 4.1 Review Your Project

- Summarize your chatbot project in 2–3 minutes, covering architecture, tools, and results.
- Anticipate follow-ups on chunking, embeddings, and challenges.

### 4.2 Study Key Concepts

- Review RAG, embeddings, vector databases, LLMOps, and prompt engineering.
- Read papers like Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.
- Explore Azure AI, LangChain, and OpenAI API documentation.

### 4.3 Practice System Design

- Sketch your projects architecture and practice designing similar systems.
- Focus on scalability, caching, and cloud integration.

### 4.4 Code Practice

- Solve AI-related problems on LeetCode (e.g., k-NN, text processing).
- Write pseudocode for retrieval or chunking functions.

### 4.5 Mock Interviews

- Simulate interviews with peers, practicing technical explanations and problem-solving.
- Use tools like Excalidraw for system design.

### 4.6 Ask Questions

- Prepare questions like: What GenAI challenges does your team face? or How do you balance experimentation and production?

## 5 Sample Practice Questions

1. How would you extend your chatbot to handle multimodal inputs?
2. What strategies reduce OpenAI API costs?
3. How would you evaluate chatbot performance in production?
4. Design a pipeline to update embeddings nightly without downtime.

5. Explain attention mechanisms in transformers.

## **6 Conclusion**

The L2 AI Engineer interview tests your technical depth, practical experience, and problem-solving skills. By leveraging your GenAI chatbot project, mastering key concepts, and practicing responses, you can confidently showcase your expertise. Focus on clear communication, justified decisions, and real-world impact to stand out.