

Problem statement- assignment 2

Q1. Where would you rate yourself on (LLM, Deep Learning, AI, ML). A, B, C [A = can code independently; B = can code under supervision; C = have little or no understanding]

Answer:

B = Can code under supervision

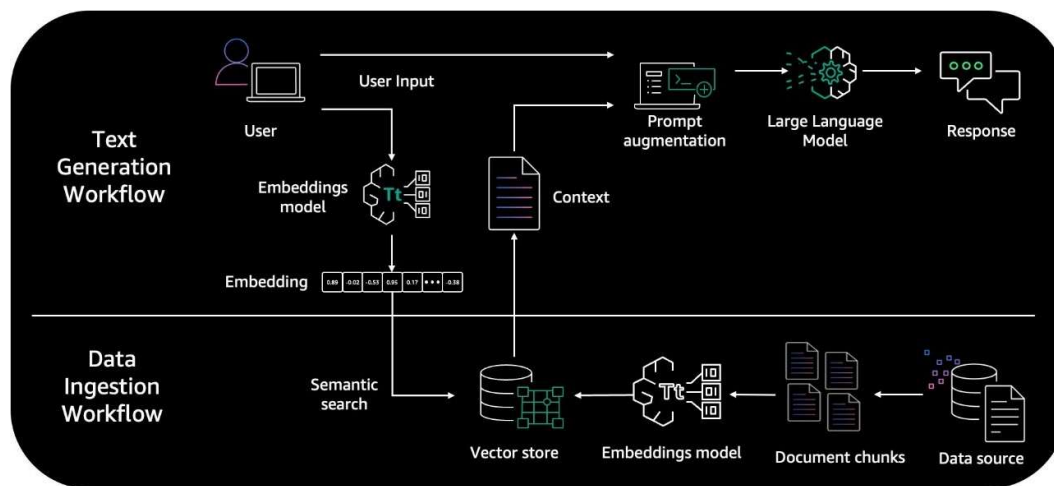
Q2. What are the key architectural components to create a chatbot based on LLM? Please explain the approach on a high-level.

Answer:

A chatbot based on Large Language Models (LLMs) relies on a modular architecture that processing user inputs, retrieves relevant context, and generates coherent responses while maintaining conversation state. The high-level approaches involve orchestrating components for ingestion, retrieval, generation, and safety.

Input layer → Embedding layer → Transformer Architecture → output layer.

LLM Based Chatbot Architecture:



Above diagram represents contextual chatbot application using knowledge bases. The data ingestion workflow uses LLMs to create embedding vectors that represent semantic meaning of texts. Embeddings are created for documents and user questions. The document embeddings are split into chunks and stored as indexes in a vector database. The text generation workflow then takes a question's embedding vector and uses it to retrieve the most similar document chunks based on vector similarity. It augments prompts with these relevant chunks to generate an answer using the LLM.

Key Architectural Components of an LLM-Based Chatbot

1. User Interface (UI)

This is the front-end through which users interact with the chatbot (e.g., web app, mobile app, or messaging platform). It captures the user query and displays the response.

2. Backend / Application Layer

This layer receives user input from the UI and handles:

- API routing
- Authentication (if required)
- Formatting queries and responses
- Managing conversation sessions

Technologies commonly used include Flask, FastAPI, Node.js, etc.

3. LLM Core Engine

This is the central intelligence of the chatbot and can be:

- Hosted LLM APIs (OpenAI GPT-4/Claude/Gemini)
- Self-hosted open-source LLMs (Llama, Mistral, etc.)

The LLM performs tasks such as:

- Natural language understanding
- Reasoning and generation
- Conversation continuity

4. Memory & Context Management

To maintain coherent multi-turn conversations, the chatbot stores and retrieves context. Two types of memory may be used:

- **Short-term memory:** conversation history within the session
- **Long-term memory:** persisted information across sessions

Databases such as Redis, PostgreSQL, or vector databases may be used.

5. Retrieval & Knowledge Base Layer (Optional but common)

For factual accuracy or domain-specific tasks, a Retrieval-Augmented Generation (RAG) pipeline is used. It includes:

- **Embeddings model** to convert documents into vector representations
- **Vector database** (e.g., Pinecone, Weaviate, Chroma)
- **Retriever** to fetch relevant context

The retrieved context is then injected into the LLM prompt.

6. External Tools and Integrations (Optional)

For transactional chatbots, APIs may be integrated for:

- Database queries
 - CRM systems
 - Search tools
 - Payment systems
- This allows the chatbot to perform actions beyond text generation.

7. Observability, Logging & Governance

For production systems, logging, analytics, and monitoring are used to:

- Track performance
- Monitor costs
- Ensure compliance
- Handle errors and auditing

High-Level Approach

The typical workflow of an LLM-based chatbot is:

1. **User sends a message** through the UI.
2. **Backend receives the message**, validates it, and manages session context.
3. **Optional retrieval step** where relevant knowledge is fetched (RAG).
4. **Prompt construction** using system instructions + context + user input.
5. **LLM generates a response** using language understanding and reasoning.
6. **Backend post-processes the output** (formatting or tool execution).
7. **Response is sent back to UI** and displayed to the user.
8. **Conversation history** may be stored for future context.

Summary

In summary, an LLM chatbot architecture consists of:

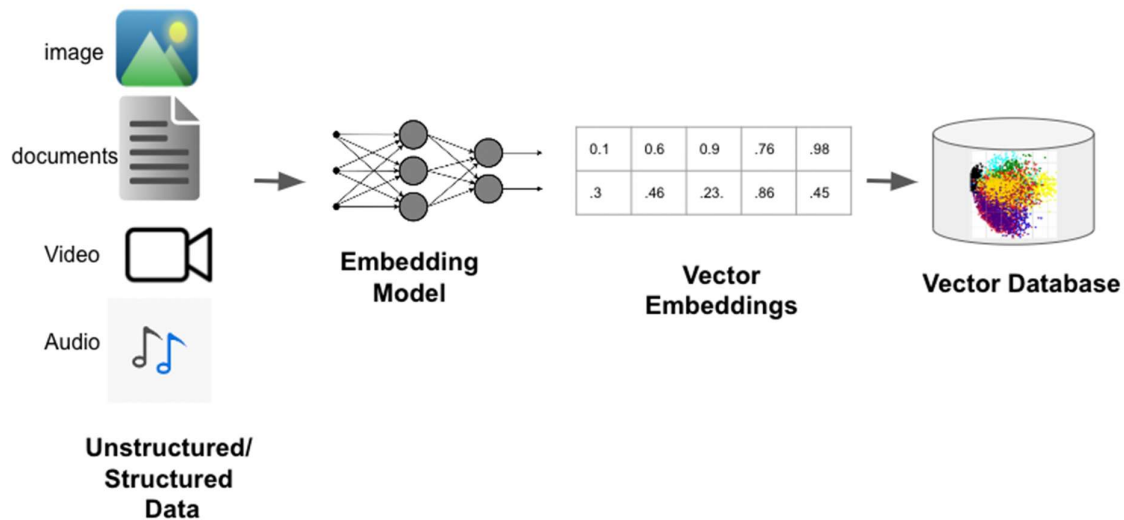
User Interface, Backend Application Layer, LLM Core Engine, Memory / Context Manager, (Optional) Retrieval Knowledge Base, (Optional) External API Integrations, Observability & Governance

These components collectively enable the chatbot to understand input, maintain context, retrieve relevant information, and generate helpful responses in natural language.

Q3. Please explain vector databases. If you were to select a vector database for a hypothetical problem (you may define the problem) which one, will you choose, and why?

Answer:

A vector database is a specialized database system designed to store, index, and retrieve high-dimensional vectors – numerical representations of data such as text, images, audio etc., These vectors also called embeddings. These vectors encode semantic or contextual information about unstructured data. Unlike traditional databases that rely on exact matches using keys and structured schemas, vector databases find data points that are close in numerical space based on similarity metrics like cosine similarity or Euclidean distance. Embeddings convert raw data into dense vectors via models BERT or CLIP, capturing semantic meaning (e.g., “king” - “man” + “woman” ~ “queen”).



The above diagram represents a process of transforming unstructured data, such as images, documents, videos, and audio, into vector embeddings. These vector representations are then stored in a vector database, which can be used for various applications like search, retrieval, and data analysis. The key steps shown are data preprocessing, creating the embedding model, generating vector embeddings, and storing them in the vector database.

Why Vector Databases needed?

Machine learning and AI models convert raw data into dense numerical vectors. These vectors represent meaning:

- Text gets converted into embeddings that capture semantic similarity.
- Images get mapped into vector space based on visual features.
- Audio clips become vectors based on frequency patterns.

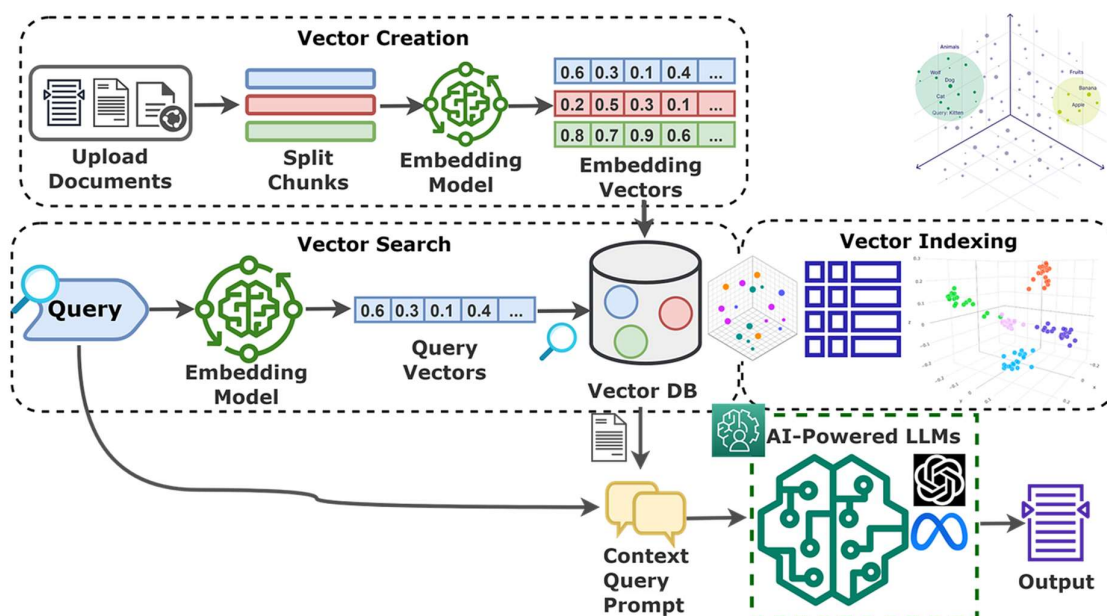
These high-dimensional vectors are large (often hundreds or thousands of dimensions) and do not fit well into relational to NoSQL databases optimized for scalar values. Vector databases are built to efficiently store, index, and query such vectors so that similarity searches are fast and scalable.

2. How Vector Databases Work (High-Level)

Vector databases work by storing data in the form of vectors, which are lists of numbers that capture the meaning of text, images, or other content. When a search is made, the database compares the new vector with the stored vectors to find which ones are most similar. This process is called similarity search.

To make the search faster, vector databases use special indexing methods. They also support filtering based on extra information, such as categories or tags. Overall, vector databases help us quickly find similar information based on meaning, not just exact words. So please refer the below diagram pipeline of how vector database work.

Vector Database work pipeline:



Key Features of Vector Databases

1. **High-Dimensional Data Support:** They can handle data represented in high-dimensional spaces, which is common in AI and ML applications.
2. **Efficient Similarity Search:** Vector databases implement algorithms like Approximate Nearest Neighbor (ANN) search to quickly find similar vectors.
3. **Scalability:** Designed to scale with large datasets, making them suitable for applications with growing data needs.
4. **Integration with ML Models:** They often integrate seamlessly with machine learning frameworks, allowing for easy storage and retrieval of model outputs.

Common Use Cases

- **Recommendation Systems:** Finding similar products or content based on user preferences.
- **Image and Video Search:** Retrieving images or videos that are visually similar.
- **Natural Language Processing:** Searching for documents or text that are semantically similar.

Hypothetical Problem: Customer Support Chatbot

Problem Definition

Imagine a customer support system for an e-commerce platform where users can ask questions about products, track orders, or resolve issues. The goal is to develop a chatbot that can understand user queries and provide relevant answers by retrieving information from a large database of FAQs, product details, and past customer interactions.

I chose Vector Database: Qdrant

Qdrant is an open-source vector database and similarity search engine, built for AI applications, that stores high-dimensional vectors (numerical representations of data like text, images) to enable fast and accurate searches based on meaning (semantics) rather than just keywords.

Why I Choose Qdrant VectorDB?

For above specific problem, I would choose **Qdrant** as the vector database. Here are the reasons:

1. **Efficient Similarity Search:** Qdrant is optimized for high-dimensional vector searches, which is crucial for quickly retrieving the most relevant responses to user queries.
2. **Scalability:** As the e-commerce platform grows, so does the volume of data. Qdrant can scale efficiently to handle increasing amounts of vector data without compromising performance.
3. **Integration with ML Models:** Qdrant allows easy integration with machine learning models that can generate embeddings for user queries and responses, enhancing the chatbot's understanding of natural language.
4. **Real-Time Performance:** The database is designed for low-latency queries, which is essential for providing quick responses in a customer support setting.
5. **Flexibility:** Qdrant supports various data types and can be used for different tasks beyond just text, such as image or video retrieval, if needed in the future.

Conclusion

In summary, vector databases like Qdrant are powerful tools for handling high-dimensional data and enabling efficient similarity searches. For a customer support chatbot in an e-commerce context, Qdrant offers the scalability, performance, and integration capabilities necessary to deliver a seamless user experience.