

## 1. explain about llms and explain the end to end process of how llms work ?

“LLMs, or Large Language Models, are AI models trained on massive amounts of text data to understand and generate human-like language.

They work using the **Transformer architecture**, which helps them understand relationships between words and generate context-aware responses.”

LLMs are built on the **Transformer architecture**, which uses two main components:

- **Encoder** – understands the meaning of the input text
- **Decoder** – generates new text based on that understanding

However, models like **GPT** (from OpenAI) use **decoder-only transformers**, which are focused on predicting the next word efficiently.

### 1. Data Collection

The process begins with collecting **massive datasets** — books, articles, Wikipedia, code, and websites.

### 2. Tokenization

Text data cannot be processed directly by a model.

So, the text is broken into small units called **tokens** (words or subwords).

For example:

“Artificial Intelligence” → [“Artificial”, “Intelligence”]

Each token is assigned a **unique numerical ID**, because the model only understands numbers.

### 3. Embeddings

Each token is converted into a **dense numerical vector** called an **embedding**.

Embeddings capture **semantic meaning** — so words with similar meanings have similar embeddings.

### 7. Vector Databases (used during RAG or Knowledge Storage)

Once embeddings are generated for texts, they can be stored in a **vector database** (like Pinecone, FAISS, or Chroma).

## 8. Cosine Similarity Matching

The system compares this **query embedding** with all the **document embeddings** stored in the **vector database**. It retrieves the most relevant data to the query and then it will generate the relevant output.

## 2. Explain about transformers?

A **Transformer** is the **neural network architecture** that powers all modern LLMs like GPT, Gemini, Claude, and LLaMA.

Let's go step by step, like how the backend pipeline works 🙌

### 1. Input Stage (Tokenization + Embedding Layer)

When a user enters text — say:

“AI is transforming the world.”

**Internally:**

- Each word is **tokenized** → [“AI”, “is”, “transforming”, “the”, “world”]
- Tokens are converted into **numeric IDs** (example: [154, 23, 5621, 9, 341])
- These IDs are passed through an **Embedding Layer**, which converts each token into a **dense vector** (e.g., 768-dimensional vector for GPT-like models).

### 2. Self-Attention Mechanism (The Real Engine 🔥 )

This is the **heart of the Transformer** — what makes it intelligent.

**Developer View:**

For each token, the model creates three internal vectors:

- **Query (Q)**
- **Key (K)**
- **Value (V)**

### 3. Feed-Forward Network (FFN)

After attention, each token vector is passed through a **fully connected neural network**. This step helps the model apply non-linear transformations and build **deeper understanding** of language.

### 4. Output Layer (Decoder)

Finally, the **decoder** predicts the next token.

### 3. explain about pre training and post training

#### 1.Pre-Training Phase (The Foundation Training)

##### ⚙️ Goal:

To make the model **understand the structure and meaning of language** — grammar, context, logic, facts, reasoning, and relationships — purely from text data.

1. **Tokenization** –  
All the text data is split into tokens (subwords or words).
2. **Input → Output Prediction Task (Next Token Prediction):**  
The model learns to **predict the next token** in a sequence.  
Example:  
Input: “Artificial intelligence is changing the”  
Target Output: “world”

This task is repeated **billions of times** until the model learns language patterns, grammar, and semantics.

3. **Architecture Used:**  
Transformer (mostly **decoder-only** for GPT-like models).  
Each layer learns to capture deeper relationships between tokens using **self-attention**.
4. **Output of Pre-Training:**  
After this phase, the model can:
  - Complete sentences
  - Understand context
  - Generate text fluently
  - Recognize relationships between entities

BUT — it’s **not yet aligned with human intentions**.

It may produce random, biased, or unsafe responses — because it just learned *patterns*, not *values*.

---

##### 🌱 Analogy:

It’s like teaching a child **the dictionary and grammar** —  
They can form sentences but don’t yet know **how to behave or answer politely**.

#### 2. Post-Training Phase (Fine-Tuning / Alignment)

##### ⚙️ Goal:

To make the model **follow instructions, act safely, and behave conversationally like a human**.

If pre-training teaches the model the **language**,  
then post-training teaches it **how to communicate like a human** — to listen, follow, and answer responsibly.

### Summary — Pre-training vs Post-training

Aspect	Pre-Training	Post-Training (Fine-Tuning)
Goal	Learn general language understanding	Align with human intent and tasks
Data	Massive unlabeled text data	Curated instruction or human-labeled data
Task	Predict next token	Follow instructions / rank responses
Model Type	Generic language model	Chat or task-oriented model
Output	Fluent language generation	Helpful, safe, human-aligned replies

---

### ✅ How to Conclude in Interview:

“So, pre-training builds the foundation — the model learns language structure by predicting the next token on large-scale data.

Post-training or fine-tuning refines this by aligning the model with human feedback and task-specific goals, using supervised datasets and reinforcement learning.

Together, these stages make an LLM both *intelligent* and *interactive*.”

## 4. LLMs we are using do they contain encoder part or decoder only part?

They are **Decoder-Only Models**.

GPT, LLaMA, Falcon, Mistral

### Why Decoder-Only?

“Most modern Large Language Models like GPT, LLaMA, Mistral, and Claude use a **Decoder-Only Transformer architecture**.

This is because they are **generative**, designed to predict the next token in a sequence using **masked self-attention**.

Think of the decoder as a *storyteller* —

it reads the context so far (previous words) and then decides what word fits next in the story.

## 5. explain about encoder part and decoder part?

Transformers consist of an Encoder and a Decoder.

The Encoder's job is to understand the input using bidirectional self-attention, generating contextual embeddings.

The Decoder's job is to generate outputs using masked self-attention and sometimes cross-attention to the encoder output.

Both are built using multi-head attention + feed-forward layers, but they serve different purposes.

### 1. Encoder — “Understanding Part”

 Goal:

To understand and encode the input sequence into contextual representations (embeddings).

Input Tokenization:

Example: “Artificial intelligence is powerful.”

→ Tokens → ['Artificial', 'intelligence', 'is', 'powerful']

Embeddings + Positional Encoding:

Each token gets converted into a numeric vector (embedding) and positional information is added.

Self-Attention Mechanism:

Every token attends to every other token — meaning it looks at all other words in the sentence to understand context.

Example:

In “bank of the river”, the encoder learns that *bank* relates to *river*, not *money*.

### 2. Decoder — “Generating Part”

 Goal:

To generate or predict the next token in a sequence using:

- The encoder's output (in encoder-decoder models)
- Or its own previous outputs (in decoder-only models like GPT)

## 6. Explain about rag architecture end to end

RAG or Retrieval-Augmented Generation is an architecture that enhances LLMs by connecting them with external knowledge sources.

It works in two stages — first, documents are chunked, embedded, and stored in a vector database. Then, during a query, the user input is embedded and compared using cosine similarity to retrieve the most relevant chunks. These chunks are appended to the user prompt and passed to the LLM, which generates a context-aware and factual response.

This approach removes hallucinations, reduces fine-tuning cost, and allows dynamic updates to the knowledge base

RAG = Retrieval + Generation.

It's a hybrid AI architecture that combines:

- Information retrieval (from external knowledge sources like vector databases), and
- Text generation (from LLMs).

### Why We Need RAG

- LLMs like GPT or LLaMA are trained only up to a certain date and cannot know recent or private data.
- They sometimes hallucinate (make up answers).
- Fine-tuning for every new dataset is expensive and slow.

### PHASE 1: Preprocessing / Data Indexing (Offline Phase)

This happens before queries come in.

#### Step 1. — Data Collection

Collect your knowledge base:

- PDFs, documents, websites, CSVs, research papers, product manuals, etc.

#### Step 2. — Chunking

- Documents are split into smaller chunks (e.g., 300–500 words each).
- Because LLMs have context length limits (like 8K or 32K tokens).
- Each chunk keeps logical meaning (like paragraph or section).

#### Step 3. — Tokenization

- Each chunk is tokenized → converted into tokens that models understand.

#### Step 4. — Embeddings Generation

- Each chunk is passed through an embedding model (like OpenAI's text-embedding-3-small or sentence-transformers).
- Embeddings = numerical vector representations of meaning.

Example:

"Python is a programming language" → [0.12, 0.87, -0.34, ...]

Now, text data is converted into vectors in a high-dimensional space.

Step 5. — Store in Vector Database

- All embeddings are stored in a Vector Database such as:
  - Pinecone
  - FAISS
  - Chroma
  - Weaviate

PHASE 2: Query / Inference Flow (Online Phase)

This is what happens when a user asks a question.

Step 1. — User Query

User asks:

“What is the difference between supervised and unsupervised learning?”

Step 2. — Query Embedding

- The query is also converted into a vector embedding using the same embedding model.  
Example: [0.11, 0.85, -0.30, ...]

Step 3. — Vector Similarity Search (Retrieval)

- The system compares the query vector with all document vectors in the vector DB.
- It finds the most similar chunks using cosine similarity.

👉 Cosine Similarity Formula:

$\text{similarity} = \frac{A \cdot B}{|A| |B|}$  Where A and B are embedding vectors.

It measures the angle between two vectors — smaller angle = higher similarity.

Step 4. — Fetch Top-k Relevant Chunks

- Top k (say 3–5) most relevant chunks are retrieved.
- Example:
  - Chunk 1: “Supervised learning uses labeled data...”
  - Chunk 2: “Unsupervised learning finds hidden patterns...”

Step 5. — Context Construction (Prompt Augmentation)

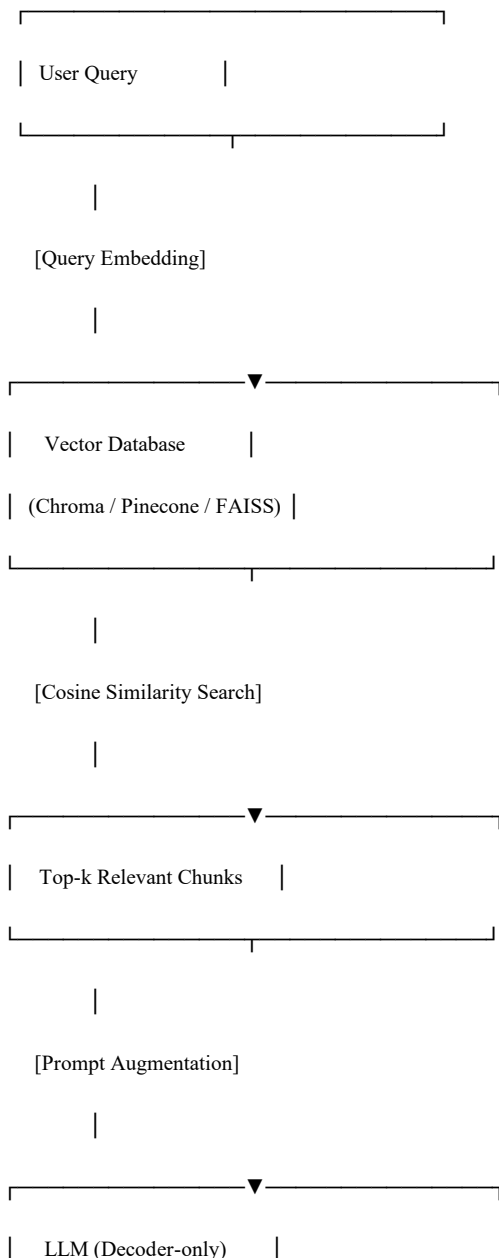
- These chunks are attached to the user’s question in the prompt.

#### Step 6. — LLM Generation

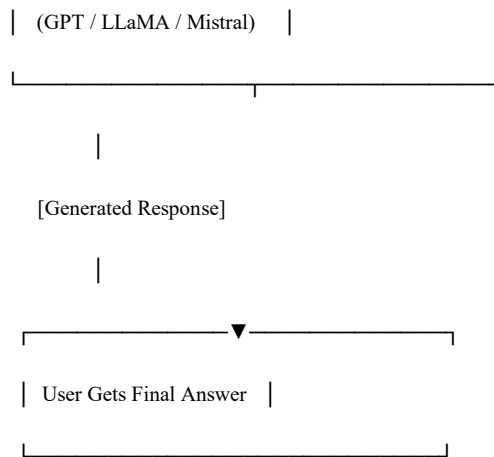
- This prompt is sent to the LLM (like GPT, LLaMA, or Falcon).
- The LLM uses:
  - Its pre-trained knowledge (general language)
  - Plus the retrieved context
- It then generates a factual, context-aware answer.

#### Step 7. — Response Returned

“Supervised learning uses labeled datasets to train models, while unsupervised learning finds structure in unlabeled data.”







## 7.explain how retriever works?

### Interview-Ready Answer (Concise Form)

“In RAG, the retriever is the component responsible for finding the most relevant information from a vector database.

When a user asks a question, the retriever first converts the query into an embedding vector using the same embedding model used during indexing.

It then performs a **vector similarity search** (commonly cosine similarity) to compare the query vector with stored document vectors, ranks them by similarity, and retrieves the top-k relevant chunks.

These chunks are then provided to the LLM as additional context for generating accurate, factual answers.”

## 8.what is temperature?

Definition:

Temperature is a value between 0 and 1 (sometimes higher) that determines how deterministic or creative the model should be while predicting the next word.

Temperature controls how “confident” or “creative” a model behaves while generating text. Lower = predictable & factual, Higher = imaginative & varied.

## 9.what is an agent?

An **Agent** is an intelligent system built around an LLM that can autonomously decide, plan, and perform tasks using external tools.

## 10.how multi agents will interact each other?

Multi-Agent Systems involve multiple autonomous agents that collaborate, communicate, and coordinate through message passing to achieve a shared goal.

In **multi-agent systems**, a **handoff** means the **transfer of control and context** from one agent to another.

Each agent specializes in a subtask, processes information using an LLM, and passes results to other agents.

Communication and orchestration are handled via frameworks like LangGraph, LangChain, or AutoGen.

### Task Assignment

A **Supervisor Agent (or Orchestrator)** receives the user query.

Example: “Generate a report on AI trends in 2025.”

It divides the task into sub-tasks:

- Agent 1 → Search latest AI news
- Agent 2 → Summarize the findings
- Agent 3 → Write a final report

## 11.which library you will use to build multi agents?

For building multi-agent systems, I use LangGraph, which is built on top of LangChain.

It allows me to create a graph-based workflow where each agent acts as a node, and control is handed off through defined edges.

## 12.explain about chunking, tokenization, embeddings, cosine similiarity?

### 1. Chunking

#### ♦ What It Is:

**Chunking** means **splitting large text data into smaller, manageable pieces (chunks)** before sending it to the LLM or embedding model.

LLMs have a **token limit** (e.g., 8k, 32k, 128k tokens), so we **divide documents** into smaller sections.

---

#### ♦ Why We Do It:

- LLMs can't process large files at once.

- Smaller chunks = faster embedding, better retrieval accuracy.
- Each chunk gets its **own embedding vector**, which helps in **precise search** later.

## 2. Tokenization

### ♦ What It Is:

**Tokenization** means **converting raw text into tokens** — the smallest units the model understands.

Each token may represent:

- A word (cat)
  - A subword (tion)
  - Or even a punctuation mark (.)
- 

### ♦ Why:

LLMs don't read words directly — they process **numbers (token IDs)**.  
So tokenization is how text becomes a sequence of **numeric tokens** for the model to process.

## 3. Embeddings

### ♦ What It Is:

An **embedding** is a **numerical vector representation** of text, capturing its **semantic meaning**.

Similar meanings → similar embeddings.

Each **chunk** (from step 1) is converted into a **high-dimensional vector** (e.g., 768 dimensions).

## 4. Cosine Similarity

### ♦ What It Is:

**Cosine similarity** measures **how close two embedding vectors are** — based on the **angle between them** in vector space.

Chunking splits large text into smaller parts.

Each chunk is tokenized into numerical tokens.

Those tokens are converted into **embeddings**, which are high-dimensional vectors representing meaning.

These embeddings are stored in a **vector database**.

During retrieval, we compare the query embedding with stored embeddings using **cosine**

**similarity** to find the most semantically relevant chunks.

Those chunks are then passed to the LLM for context-aware answer generation