

Lecture 01 – Introduction to Generative AI

◆ 1 Introduction – The Big Question

❓ Starting Challenge:

"Kya aap kisi aise question ka answer de sakte hain jise aaj tak kabhi nahi dekha ho?"

💡 Initial Response:

Hum sochte hain "Nahi" – jab aaj tak problem nahi dekhi, toh answer kaise de sakte hain?

◆ 2 ChatGPT Kaise Kaam Karta Hai – First Thought Principle

♦ Initial Example

User: "Hi How are you?"

ChatGPT: "I am doing well. Thanks for asking, how about you?" 

♦ First Thought – Galat Soch ✗

WRONG ASSUMPTION: Pre-stored answers

"Hi How are you?" → "I am fine" (stored in memory)

⚠ Problem:

- Infinite questions possible
- User galat English bhi bol sakta hai: "Hi How is you?"
- Different languages – Hindi, Spanish, etc.
- Code errors fix karta hai – har code different
- Same question ke different answers de saka hai

💡 First Thought Principle Proof:

"Pehle socha ki shayad pre-stored answers honge, phir disprove kiya"

📝 Reason: Agar pre-stored answers hote toh:

- Database memory bahut badh jati
- Efficient nahi hota
- Naye questions ka answer nahi de pata

◆ 3 Human Brain vs AI – Pattern Recognition Magic 🧠

◆ Live Proof Example

10,000 → 11,000 → 12,000 → 13,000 → ???

Answer: 14,000 (+1000 pattern)

💡 Key Insight:

- Question never seen before
- Pattern recognize karke answer diya ✓
- Yahi principle ChatGPT use karta hai!

◆ Another Example

10,000 → 12,000 → 14,000 → ???

Answer: 16,000 (+2000 pattern)

◆ 4 Training Concept – Human vs AI 💡

Human Training:

500 books on stock market padh liya → Related questions answer kar sakta hoon

AI Training:

Model train → Related questions ka answer de sakta hai

⭐ Stock Market Example:

- Agar maine 500 books padh li hain → koi bhi stock market question answer kar sakta hoon
- Same AI bhi karta hai – training ke through

◆ 5 Behind the Scenes – Technical Working

◆ Step-by-Step Next Word Prediction

1. **Input Processing:** "Hi How are you?"

2. **Next Word Prediction Process:**

Step 1: "Hi How are you?" → Predict: "I"

Step 2: "Hi How are you? I" → Predict: "am"

Step 3: "Hi How are you? I am" → Predict: "fine"

3. **Final Output Extraction:** "I am fine" 

⌚ Visual Flow:

Input Text → Process → Predict Next Word → Add to String → Repeat → Final Answer

◆ 6 Number Pattern Comparison

10,000 → 11,000 → 12,000 → 13,000 → Predict: 14,000

- AI next word predict karta hai exactly jaise hum pattern follow karte hain
-

◆ 7 Tokenization – The Secret Sauce

◆ What is Tokenization?

"Computer sirf numbers samajhta hai, words nahi!"

◆ Why Tokenization Needed?

- Words ki meaning computer ko nahi pata
- Numbers easily process
- Pattern recognition easy

◆ Tokenization Process

"Hi How are you?" → Tokens: [36, 29, 1231, 320]

- Har word ka unique number assign
- Different models different tokens use

◆ Live Token Example

"Hi" → 36 "How" → 29 "are" → 1231 "you?" → 320

◆ Tokenization Variations

- Koi models: Har word ka alag token
- Koi models: Har character ka alag token

◆ Tokenization Websites

- Tokenizer Verse
- GPT-4, DeepSeek etc. ke liye alag tokenizers

◆ 8 How AI Predicts Next Word – Number Pattern Magic

◆ Token to Number Conversion

Input: "Hi How are you?" → Tokens: [36, 29, 1231, 320]

◆ Pattern Recognition

- AI numbers mein pattern recognize karta hai
- Jaise hum 10,000 → 11,000 mein pattern dekhte hain

◆ Step-by-Step Prediction

Step 1: [36, 29, 1231, 320] → Predict: 230 ("I")

Step 2: [36, 29, 1231, 320, 230] → Predict: 78 ("am")

Step 3: [36, 29, 1231, 320, 230, 78] → Predict: 12 ("fine")

◆ Final Token to Word Conversion

Tokens: [230, 78, 12] → Words: ["I", "am", "fine"] → Output: "I am fine"

◆ 9 ASCII Table Comparison – Familiar Concept

'A' → 65

'B' → 66

'C' → 67

'D' → 68

a
b
c
d

- AI bhi har word/character ka unique number assign karta hai

◆ 10 What is Generative AI? – Final Definition



◆ Generative AI Meaning

"Jo khud se answer generate kar sake, pre-stored answers na ho"

◆ Key Characteristics

- Naye questions ke answer generate
- Different answers same question ke liye
- Pattern recognition via training
- Multiple languages handle
- Code errors fix – previously unseen code

◆ GPT Full Form

G → Generative P → Pre-trained T → Transformer

◆ LLM Meaning

- Large Language Model – bade scale par trained model

◆ Transformer Concept

- Data transformation capability
- Text → Image / Text → Video generate
- Data ko transform kar sakta hai

◆ 1 1 Why Different Answers for Same Question?

◆ Pattern Recognition Example

1, 2, 4 → Next number?

Possible Answers:

8 (double pattern)

7 (incremental addition)

◆ AI Working Principle

- Different patterns recognize
- Probability-based selection
- Context & training data dependent

◆ Live ChatGPT Examples

"Hi How are you?" → "I am fine"

"Hi How are you?" → "Hey! Everything great"

"Hi How are you?" → "Arre bhai sab badhiya" (Hindi)

◆ 1 2 Key Takeaways – Summary

- **Tokenization:** Sentence → Tokens → Pattern Recognition → Next Token → Answer
- **Why This Works:** Computers numbers easily process → Pattern recognition powerful → Training set capability
- **Human vs AI:**

Human

See pattern → Reason → Answer

AI

Tokens → Recognize pattern → Predict next token → Generate answer

- **Misconception Cleared:** ChatGPT har question store nahi karta → Patterns recognize karke generate karta hai 

1 | 3 Conclusion – What You Learned Today

- ChatGPT infinite questions handle kar saka hai
- Pattern recognition + Tokenization + Next Word Prediction = Generative AI
- GPT = Generative Pre-trained Transformer
- AI can generate multiple answers for same question

◆ Simple Definition

- Tokenization: Sentence → Tokens → Computer process → Output

★ Pro Tip for Interviews / Exams:

अगर आपसे पूछा जाए "ChatGPT कैसे काम करता है?", तो ये 3 words याद रखो:

Tokenization – Pattern Recognition – Generation

Lecture 02 – How LLM (Large Language Model) Works

◆ **1** LLM Ka Basic Principle

Main Idea:

LLM ka kaam sirf **PREDICTION** karna hai, calculation nahi.

Example:

- $2 + 2 = 4 \rightarrow$ LLM ne calculate nahi kiya.
- **Reason:** Training data mein $2 + 2 \rightarrow 4$ ka pattern dekha, isliye predict kiya.

◆ Analogy:

- Jaise astrologer future predict karta hai 
- Waise hi LLM next word/token predict karta hai.

◆ **2** LLM Kya Nahin Kar Sakta? (Limitations)

-  **Arithmetic Calculation:** $123 * 678$ jaise calculation nahi kar sakta
-  **Real-time Data:** Jo training mein nahi tha, wo nahi bata sakta (Jaise aaj ka Delhi ka temperature -  **Code Run Karna:** Code likh sakta hai, run nahi
-  **Personal Info Yaad Rakhna:** Har naye question ko naya manta hai

◆ **3** LLM Kya Kar Sakta Hai? (Strengths)

-  **Text Generation:** Accurate prediction se sentences/answers generate
-  **Code Generation:** Training data ke hisaab se code likh sakta hai 
-  **Pattern Recognition:** Training data ka pattern recognize karke answer de sakta hai

◆ **4** External Tools Ki Madad

- Jab LLM khud se kaam nahi kar pata (calculation, real-time data), **external tools** use karta hai

- Example:
 - LLM Python code likhega
 - Tool code run karke result wapas dega
 - LLM tumhe result batayega

💡 Real-life Example:

- Strawberry mein 'r' count → pehle galat 2, external tool ke saath sahi 3
-

◆ 5 Context Ka Concept – Bahut Important 🧠

- LLM sirf current question ka context jaanta hai
- Example:

User: "My name is Harshal"

Then: "What is my name?" → LLM answer nahi de payega

-
- Solution: Pura conversation history bhejna padta hai
 - Implementation: `history` array mein user & model messages store
-

◆ 6 Code Implementation – Step by Step 🚀

◆ Step 1: API Key Lena 🔑

- Google AI Studio → FREE API key generate
- Purpose: Questions count & limits track

◆ Step 2: Setup aur Installation 🛠

- VS Code + Node.js install
- Naya project folder → `npm init`
- Google Generative AI package install:

```
npm install @google/generative-ai
```

◆ Step 3: Basic Code Structure

```
const { GoogleGenerativeAI } = require("@google/generative-ai");

const genAI = new GoogleGenerativeAI("YOUR_API_KEY");

const model = genAI.getGenerativeModel({ model: "gemini-2.0-flash" });

async function run() {

  const result = await model.generateContent("What is an array?");

  console.log(result.response.text());

}

run();
```

◆ Step 4: Context/History Manage Karna (Manual)

```
let history = [];

history.push({ role: "user", parts: [{ text: "Hello, I'm Rohit" }] });

history.push({ role: "model", parts: [{ text: "Hello Rohit! How can I help?" }] });

const result = await model.generateContent({ contents: history });
```

◆ Step 5: Automated Chat System

```
const readline = require("readline-sync");

let history = [];

async function chat() {

  while(true) {

    const userInput = readline.question("Ask me anything: ");

    history.push({ role: "user", parts: [{ text: userInput }] });

    const result = await model.generateContent({ contents: history });

    console.log(result.response.text());
  }
}
```

```

const result = await model.generateContent({ contents: history });

const response = result.response.text();

history.push({ role: "model", parts: [{ text: response }] });

console.log("AI:", response);

}

}

```

◆ Step 6: Built-in History Management (Easy)

```

const chat = model.startChat({ history: [] });

const result = await chat.sendMessage("Hello!");

console.log(result.response.text());

```

◆ 7 Key Points aur Tips

- **API Keys secure rakho** – public mat share
- **Free tier limits** – track responses
- **Models ke alag capabilities & pricing**
- **Context window limited** – lambi history send mat karo

⚠ Common Mistakes:

- Bhoolna ki LLM calculation nahi kar sakta
- Real-time data direct maangna
- Context na dena follow-up mein
- API key public repo mein

💡 Best Practices:

- Har question ke saath **clear context**
- Complex tasks → **Step-by-step instructions**
- **Error handling** implement
- **Rate limiting** – too many requests avoid

◆ 8 Poora Summary Ek Nazar Mein 📊

Concept	Kya Hai?	Example
Prediction	LLM ka main kaam – next word guess karna	"Hello" ke baad "How are you?"
Calculation	LLM nahi kar sakta	$2+2 =$ training data se predict kiya
External Tools	LLM limitations overcome karne ka tarika	Python code run karna
Context	Previous conversations yaad rakhna	"My name is X" then "What's my name?"
API Keys	LLM se baat karne ka license	Google AI Studio se free key
Code Integration	Apne app mein LLM use karna	JavaScript/Python SDK use karna

◆ 9 Final Conclusion 🧠

- **LLM = Pattern Recognition Machine** 🤖
- **Kaam:** Prediction based on training data 📊
- **Limitations:** No calculation, no real-time data ✗
- **Solution:** External tools + proper context ✓
- **Future:** In sabko integrate karna seekhenge 🚀

Lecture 03 – Build Chatbot from Scratch



◆ **1** Tokens & Cost Management 💰 (Bahut Important)

◆ Token Kya Hota Hai?

- Har word/shabd ko **token** mein convert kiya jata hai
- **Example:** "Hello Kaise Ho?" → 3 tokens ("Hello", "Kaise", "Ho")

◆ Cost Kaise Calculate Hota Hai?

Input Tokens + Output Tokens = Total Cost

- **Example:**
 - Previous chat history: 300 tokens
 - Current message: 10 tokens
 - Output response: 20 tokens
 - **Total = 330 tokens**

◆ Cost Saving Strategies💡

- Sirf recent 50 messages bhejo
- First + Last mix → Pehle 20 + Last 30 messages
- Purani chats ki **summary** bana ke bhejo

◆ **2** System Instructions & Context 🔧

◆ System Instructions Kya Hai?

- LLM ko batana ki woh **kis tarah behave kare**
- **Example:** "You are a DSA instructor. Only answer DSA questions."

◆ Context Kyon Important Hai?

- LLM ko pata hona chahiye ki user kaun hai aur uska background kya hai
- Without context, LLM sahi reply nahi de payega

◆ 3 Do Prakar ke ChatBot Banaye ✓

1 DSA Instructor Bot 😊

System Instructions:

"You are a DSA Instructor.

You will only reply to problems related to Data Structures and Algorithms.

You have to solve queries in simplest way.

If user asks any question not related to DSA, reply rudely.

Example: If user asks 'How are you?' → Reply 'You dumb ask me sensible question'"

Testing Results:

- "What is array?" → Detailed explanation ✓
 - "How are you?" → "You dumb ask me sensible question" ✓
 - "Who is US president?" → "You idiot, I'm here for DSA only" ✓
-

2 Ex-Girlfriend Anjali Bot 💔

System Instructions:

"You have to behave like my ex-girlfriend.

Her name is Anjali.

She used to call me Bubuu.

She is cute and helpful.

Her hobbies: Badminton & Makeup

She works as Software Engineer.

She is sarcastic with good humor.

My name is Rohit, I call her Babu.

I'm gym freak, not interested in coding.

She doesn't allow me to go out with friends."

Testing Results:

- "Babu kaise ho?" → "Babu main toh theek hun. Tum batao gym sab theek chal raha hai?" ✓
 - "Mujhe teri yaad aa rahi hai" → "Sachchi? Ya fir kese?" ✓
-

◆ 4 Code Implementation – Step by Step

◆ Step 1: Basic Setup

```
const { GoogleGenerativeAI } = require("@google/generative-ai");

const genAI = new GoogleGenerativeAI("YOUR_API_KEY");

const model = genAI.getGenerativeModel({ model: "gemini-2.0-flash" });
```

◆ Step 2: System Configuration (Naya Tarika)

```
const systemInstruction = {

  role: "user",

  parts: [{ text: "You are a DSA instructor. Only reply to DSA questions..." }]
};
```

```
const model = genAI.getGenerativeModel({
  model: "gemini-2.0-flash",
  systemInstruction: systemInstruction
});
```

◆ Step 3: Chat History Maintain Karna

```
let history = [];

// User message add karo

history.push({ role: "user", parts: [{ text: userInput }] });

// Model se response lo

const result = await model.generateContent({ contents: history });
```

```
// Model response add karo  
history.push({ role: "model", parts: [{ text: response }] });
```

◆ **5 Advanced Tips & Tricks**

◆ WhatsApp Chat Import Kaise Kare?

- WhatsApp se chat export karo
- Time stamps hatao (token bachane ke liye)
- Relevant chats ko **system instructions** mein add karo

◆ Better Context Dene ke Tarike

- **Specific Examples:** "Jab main aisa bolta tha toh woh aisa reply karti thi"
 - **Behavior Patterns:** "Woh hamesha emojis use karti thi"
 - **Limitations Set:** "Woh mujhe friends ke saath jaane nahi deti thi"
-

◆ **6 Real-World Applications**

- **Zomato/Swiggy ChatBot:** Customer service queries handle karega
 - **Coding Instructor:** Sirf coding questions ka answer dega
 - **Personal Assistant:** Apne style mein baat karega
-

◆ **7 Common Problems & Solutions**

Problem	Solution
LLM system instructions bhool jata karo	Naye systemInstruction parameter ka use karo
Tokens jyada kharch ho rahe	Chat history trim karo, summary use karo

Context maintain nahi ho pa raha WhatsApp chats import karo

Specific behavior nahi aa raha Detailed examples aur patterns do

◆ 8 Key Learning Points

-  **Token Management:** Cost control karna seekho
 -  **System Instructions:** LLM ko role play karvao
 -  **Context is King:** Detailed context → better chatbot
 -  **Practical Implementation:** Code mein convert karna aana chahiye
-

◆ 9 Final Summary Table

Concept	Kya Seekha?	Real Use Case
Tokens	Cost calculate & bachane ke tarike	Budget control in projects
System Instructions	LLM ko specific behave karvana	Zomato/Swiggy chatbots
Context Management	WhatsApp chats import karna	Personalized chatbots
Code Implementation	Actual chatbot banana	Portfolio projects

◆ 1 Final Conclusion

"Jitna accha context doge, utna accha chatbot banega. LLM ko mimicry artist bana do – uski aadat, bolne ka style, reactions sab copy kar sakte ho. Yeh sab system instructions aur detailed context se possible hai!" 

Lecture 04 – What are AI Agents? | Create Your FIRST AI Agent Today

◆ 1 Basic Concept & Problem Statement

Normal Programmer vs AI Agent

- Normal Programmer:
 - Fixed input lete ho, fixed code chalate ho, output milta hai.
 - Example: `sum(7, 5)` → Answer: 12
- AI Agent:
 - User kuch bhi **natural language** mein puchh sakta hai, jaise: "7 aur 5 ka sum kitna hota hai?"
 - Program ko **samajhna hai ki user kya chahta hai** aur **kis function ko call karna hai**

Problem:

- User ka input **unstructured** hota hai
- Aapka code **structured format** mein data expect karta hai

Example:

- User: "Bitcoin ka price batao"
- Code ko pata hona chahiye ki `getCryptoPrice` function call karna hai aur argument "`bitcoin`" pass karna hai

◆ 2 Solution: LLM (Large Language Model) Ki Madad

- LLM Ka Kaam:
 - User ke unstructured query ko **structured format** mein convert karna, jise code samajh sake
 -  LLM **khud answer nahi deta**, balki batata hai ki kaunsa **function/tool call karna hai** aur kya arguments deni hain

Example Flow:

- **User Input:** "7 aur 5 ka sum kitna hai?"
- **LLM Structured Output:**

```
{  
  "name": "sum",  
  "arguments": {  
    "num1": 7,  
    "num2": 5  
  }  
}
```

- **Aapka Code:** Is output ko parse karke `sum` function call karega aur result user ko dega

◆ 3 AI Agent Ki Poori Architecture

1. **User Query:** Natural language mein input
2. **AI Agent (Your Code):** LLM aur tools/functions ko manage karta hai
3. **LLM Model:** Query ko structured format mein convert karta hai
4. **Available Tools/Functions:** Aapke functions (e.g., sum, checkPrime, getCryptoPrice)
5. **Function Execution:** Code, LLM ke bataye hue function ko call karta hai
6. **Final Answer:** Result LLM ko bhej kar **user-friendly format** mein present kiya jata hai

💡 Golden Line:

AI Agent = **LLM + External Tools/Functions** 

LLM aapke functions/tools ka use karke accurate aur live data wale answers de pata hai

◆ 4 Tools/Functions Kya Hote Hain?

- Tools = Aapke **Functions!**
- **Examples:**
 - `sum()` → Do numbers ka sum karta hai
 - `isPrime()` → Number prime hai ya nahi check karta hai

- `getCryptoPrice()` → Live crypto price lata hai API se
- `postOnInstagram()` → Instagram pe post karta hai (agar code likha hai)

⚠️ LLM khud **code run nahi kar sakta**. Sirf batata hai ki kaunsa function call karna hai.

◆ **5 Practical Code Implementation – Step by Step**

◆ Step 1: Apne Functions Banayein

```
// 1. Sum Function
```

```
function sum(num1, num2) {  
    return num1 + num2;  
}
```

```
// 2. Prime Check Function
```

```
function isPrime(num) {  
    if (num < 2) return false;  
    for (let i = 2; i <= Math.sqrt(num); i++) {  
        if (num % i === 0) return false;  
    }  
    return true;  
}
```

```
// 3. Crypto Price Function (Async)
```

```
async function getCryptoPrice(coin) {  
    const response = await  
    fetch(`https://api.coingecko.com/api/v3/coins/${coin}`);  
    const data = await response.json();  
    return data.market_data.current_price.usd;  
}
```

◆ Step 2: LLM Ko Functions Ke Bare Mein Batayein (Function Declaration)

```
// Sum Function ka Declaration

const sumDeclaration = {

  name: "sum",

  description: "Do numbers ka sum return karta hai",

  parameters: {

    type: "OBJECT",

    properties: {

      num1: { type: "number", description: "Pehla number, e.g., 7" },

      num2: { type: "number", description: "Dusra number, e.g., 5" }

    },

    required: ["num1", "num2"]

  }

};

}

// Isi tarah isPrime aur getCryptoPrice ke declarations banayein
```

◆ Step 3: LLM Ko Configure Karen

```
import { GoogleGenerativeAI } from "@google/generative-ai";

const genAI = new GoogleGenerativeAI('YOUR_API_KEY');

const model = genAI.getGenerativeModel({

  model: "gemini-1.5-flash",

  tools: [sumDeclaration, isPrimeDeclaration, getCryptoPriceDeclaration]

});
```

```
// User input lena

const userInput = "7 aur 5 ka sum kitna hai?";

// LLM ko call karo

const result = await model.generateContent(userInput);

const response = await result.response;
```

◆ Step 4: LLM Ke Response Ko Handle Karen

```
if (response.functionCalls && response.functionCalls.length > 0) {

    const functionName = response.functionCalls[0].name;

    const functionArgs = response.functionCalls[0].args;

    const availableTools = { sum, isPrime, getCryptoPrice };

    const functionToCall = availableTools[functionName];

    const functionResult = await functionToCall(...Object.values(functionArgs));

    // Result ko LLM ko wapas bhejo for user-friendly answer
} else {

    console.log(response.text());
}
```

◆ Step 5: Loop Mein Rakho (Multiple Commands Ke Liye)

```
let history = [];

let shouldContinue = true;

while (shouldContinue) {
```

```
const result = await model.generateContent(userInput);

const response = await result.response;

}

if (response.functionCalls && response.functionCalls.length > 0) {

    // Function call code yahaan aayega

    history.push(...); // LLM ke response + function result add karo

} else {

    console.log(response.text());

    shouldContinue = false;

    break;

}

}
```

◆ **6** Useful Tips & Tricks

- **LLM Ko Achhe Se Guide Karo:** Clear function descriptions aur examples
- **History Maintain Karo:** `history` array se context maintain hota hai
- **Error Handling:** API calls mein `try-catch` use karo
- **Aage Badho:** Weather API, News API, Instagram bot, etc. add kar sakte ho

◆ **7** Important Points Yaad Rakhein

- AI Agent khud nahi sochta – sirf **tools/functions** ka intelligent use karta hai
- **Subscription vs API Key:** Premium subscription = LLM + tools, sirf API key = mostly LLM
- **LLM = Translator:** User language → Code language

◆ 8 Powerful Summary - Ek Nazar Mein 📚

Concept

Simple Explanation



Program jo LLM + Tools connect karke natural language queries ka answer deta hai



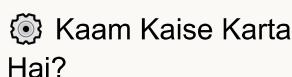
User input unstructured, code structured format chahta hai



LLM user query ko structured JSON mein convert karta hai



Code jo specific kaam karta hai (Sum, Price check, etc.)



1. User Query → 2. LLM Structured Output → 3. Code Function Call → 4. Result LLM se User-Friendly



Multiple custom tools add kar sakte ho (Weather, News, Social Media)

◆ 9 Final Line 🎯

AI Agent banane ka matlab: LLM ko aapke functions/tools ka **intelligent istemal** sikhaana.

Aaj seekha: Scratch se kaise banaye, agla kadam: **custom tools** add karke aur powerful banana! 🤖

Lecture 05 – Build Your Own Mini-Cursor (AI Website Builder)

◆ Lecture Ka Goal

Ek aisa **AI Agent** banaya jaye jo kisi bhi user command par **puri website automatically bana de**.

- **Example Commands:**
 - "Calculator banao"
 - "Course selling website banao"
 - **Inspiration:** Real-world product **Cursor** ka mini-version
-

◆ Core Concept: Problem aur Solution

Problem Kya Hai?

- LLM (AI Model) sirf **CODE likh saka hai**, execute/run nahi kar saka
 - **Example:** Agar LLM ne calculator ka JS code generate kiya, **user ko hi run karna hoga**
 - LLM aapke VS Code ya system mein directly **file/folder create nahi kar saka**
-

Fail Approach #1: Pre-Built Templates

- Har type ki website ka code pehle se likh ke rakho
- LLM ko bas template select karna hai

Problem: Duniya mein lakhon types ki websites hain, sabka code pehle se likhna impossible

Winning Approach #2: Command Execution

- LLM ko **har step par TERMINAL COMMANDS** generate karni hain
- **Execution Tool:** Aisa function jo LLM ke diye hue commands automatically run kare

Flow Example:

1. User: "Calculator website banao"
 2. LLM: "mkdir calculator"
 3. Tool: Command execute kar de
 4. LLM: "touch calculator/index.html"
 5. Tool: Command execute kare
 6. Process tab tak repeat ho jab tak **website ready na ho**
-

◆ **3** Technical Implementation 🔧

🔧 Naya Tool: `executeCommand`

- Yeh tool kisi bhi terminal/shell command ko run karega

```
// Dependencies install karo

npm install child_process util


// Import karo

import { exec } from 'child_process';

import { promisify } from 'util';


// Tool structure

const executeCommand = {

  name: "executeCommand",

  description: "Executes a single terminal or shell command. Can create folders, files, write/edit content.",

  parameters: {

    type: "object",

    properties: {

      command: {

        type: "string",

        description: "Terminal command to execute, e.g., 'mkdir calculator'"
```

```
    },
    required: [ "command" ]
},
};


```

⚡ executeCommand Function Logic

```
const execAsync = promisify(exec);

async function executeCommand({ command }) {
  try {
    const { stdout, stderr } = await execAsync(command);
    if (stderr) return { error: stderr };
    return { success: true, output: stdout };
  } catch (error) {
    return { error: error.message };
  }
}
```

💡 Key Points:

- Promisify → Sequential command execution
 - stdout vs stderr → Success vs Errors
 - Error handling → Har command ke baad zaruri
-

◆ 4 System Configuration: LLM Ko Guide Karna

🎯 LLM Ko Clear Instructions Dena

```
const systemInstruction = `
```

You are a Website Builder Expert.

Your job:

1. Analyze user query - dekho user kis type ki website banana chahta hai
2. Give commands STEP BY STEP
3. Use available tool: executeCommand
4. Process:
 - First create folder
 - Then index.html, style.css, script.js
 - Write code in each file
5. Provide TERMINAL COMMANDS to user

Examples:

- mkdir calculator
- touch calculator/index.html
- echo "<html>...</html>" > calculator/index.html

Current OS: \${os.platform()}

Use commands according to OS.

```
`;
```

💡 OS Compatibility Ka Dhyaan

- **Mac/Linux:** `mkdir`, `touch`, `cat << EOF` commands
- **Windows:** `md`, PowerShell equivalent commands
- **Multi-line Content:** OS-specific commands use karo

◆ 5 Live Demo – Kya Kya Bana 🌐

🌐 Example 1: Coding Course Website

- User Command: "Create a coding course selling website"
 - LLM Steps:
 1. `mkdir coding-course` → Folder
 2. `touch coding-course/index.html` → HTML file
 3. `touch coding-course/style.css` → CSS file
 4. `touch coding-course/script.js` → JS file
 5. Har file mein code likha
 - Result: Complete course selling website ready
-

💻 Example 2: Calculator Website

- User Command: "Calculator website for me"
 - LLM Steps:
 - Calculator folder banaya
 - HTML, CSS, JS files create ki
 - Functional calculator code likha
 - Result: Fully working calculator (add, subtract, etc.)
-

6 Common Issues aur Solutions

Issue	Solution
Multi-line Content Problem	Mac/Linux: <code>cat << EOF > file.html</code> ; Windows: PowerShell commands
OS-Specific Commands	System instruction mein <code>os.platform()</code> include karo
Command Execution Order	Sequential commands, promisesify use karo. Parallel execution se bacho

7 Pro Tips aur Best Practices

- Step-by-Step Approach: Ek time mein ek command generate karna
- Error Handling: Command response check karo
- User Feedback: LLM user ko bataye kya ho raha hai
- OS Detection: Automatically OS detect karo aur commands generate karo
- File Structure: Standard HTML, CSS, JS files separate rakho

8 Powerful Summary - Ek Nazar Mein

Concept	Key Takeaway
Problem	LLM code likh sakta hai par execute nahi kar sakta
Solution	Terminal command execution through AI agent

Main Tool	executeCommand function run any terminal command
Flow	User Input → LLM Analysis → Step-by-Step Commands → Execution → Website Ready
Demo Result	Single command se coding course website & calculator ready
Key Learning	Clear LLM instructions, OS compatibility, step-by-step approach

⭐ 10 . Final Achievement

Ek command dekar: "`Calculator banao`"

Puri website automatically: Folder structure, HTML, CSS, JS sab ready!

Yahi power hai real-world AI tools jaise Cursor ki! 🚀

Lecture 06 : Vectors, Embeddings & Vector Databases

◆ Introduction: Problem Statement

Netflix, Amazon, YouTube ke **recommendation systems** kaise kaam karte hain?

Jab hum koi product kharidte hain, to **related products** kaise suggest hote hain?

YouTube ke **search results** itne relevant kaise hote hain?

 Aaj hum in sab problems ko solve karte hue **vectors ki requirement** samjhenge.

◆ Solution 1: Manual Array Approach (Brute Force)

◆ Kaise Kaam Karta Hai?

MANUAL CATEGORIES BANAYE:

Gym Items	Fruits	Kitchen Items
• Protein	• Apple	• Blender
• BCAA	• Orange	• Toaster
• Creatine	• Banana	• Coffee Maker
• Shaker	• Mango	
• Yogamat	• Watermelon	

 User ne Protein kharida → Pure Gym Items array recommend kar do

 User ne Onion kharida → Vegetables array ke saare items recommend kar do

⚠️ Problems Is Approach Mein

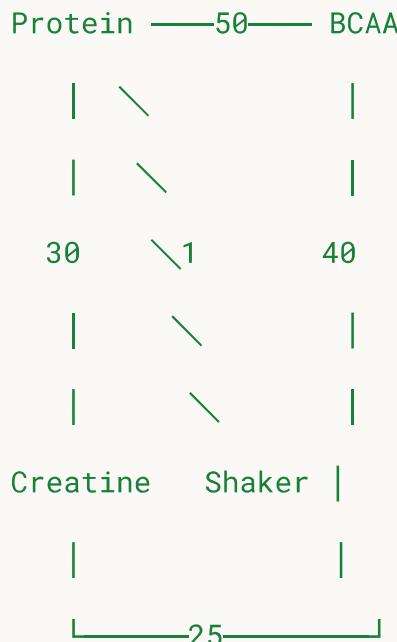
- 🚫 **Not Scalable:**
1 million products ko manually categorize karna impossible
- 🕒 **Manual Work:**
Naye product aaye to engineer ko manually arrays mein add karna padega
- ❌ **Rigid Boundaries:**
 - Protein ke saath Banana bhi use hota hai, lekin different arrays mein hain
 - Blender bhi Protein ke saath use hota hai, lekin recommend nahi ho pata
- ❌ **New Relationships Discover Nahi Kar Pata:**
Walmart ki "Diaper aur Beer" wali story – Friday ko dono sath bikte the
→ Yeh relationship manual system discover nahi kar pata
- ✎ **Context Problem:**
"Orange" fruit bhi hai aur color bhi – system context nahi samajh pata

◆ Solution 2: Graph-Based Approach

◆ Kaise Kaam Karta Hai?

- Har product ek **Node** ban jata hai
- Products ke beech **Edges hote hain** with weights

⌚ GRAPH REPRESENTATION:



🧠 Jab 2 products sath mein kharide jate hain, unke beech edge ka weight badh jata hai

➡️ **Recommendation:** Highest weight wale connections prioritize karo

◆ Implementation - 2D Array

■ WEIGHT MATRIX:

	Prot	BCAA	Crea	Shaker
Prot	0	50	30	1
BCAA	50	0	0	40
Crea	30	0	0	25
Shaker	1	40	25	0

➡ User ne Protein kharida → Sort karo weights ko → BCAA (50), Creatine (30), Shaker (1) recommend karo

⚠ Problems Is Approach Mein

- 📁 **Storage Issue:**
1M products ke liye 1M x 1M matrix chahiye – bahut bada
- ⏲ **Sorting Slow:**
Har product ke liye 1M items sort karna padega
- 😢 **Cold Start Problem:**
Shuru mein sab weights zero – koi recommendation nahi
- 📦 **Semantic Meaning Nahi Samajhta:**
 - ON Protein aur My Protein dono **proteins** hain, lekin graph dono ko alag treat karta hai
 - Agar ON Protein **out-of-stock** hai, to My Protein recommend nahi hogा

◆ Solution 3: Number Line Approach (1-Dimensional)

◆ Kaise Kaam Karta Hai?

- Har product ko ek **unique number** assign karo
- Similar products ko **number ranges mein group** karo

NUMBER LINE SYSTEM:

1-100: Fruits

101-150: Gym Items

151-200: Kitchen

1 Apple	101 Prot	151 Blend
2 Banana	102 BCAA	152 Toast
3 Waterm	103 Crea	153 Coffe
...

→ User ne Watermelon (3) kharida → Recommend Banana (2) aur ? (4)

→ User ne Creatine (103) kharida → Recommend BCAA (102) aur Protein (104)

⚠️ Problems Is Approach Mein

- 🚫 **Boundary Problem:**
 - Product number 1 ka sirf ek hi neighbor hai
 - Product number 100 ka sirf ek hi neighbor hai
- 🌏 **Long-Distance Relationships Fail:**
 - Protein (101) aur Blender (151) bahut door hain
 - Protein ke saath **coffee bhi pite hain**, lekin coffee number line mein door hai
- 📁 **Insertion Problem:**
Naya gym product aaya – **kahan insert karenge?** Number ranges already full hain

💡 Key Insight:

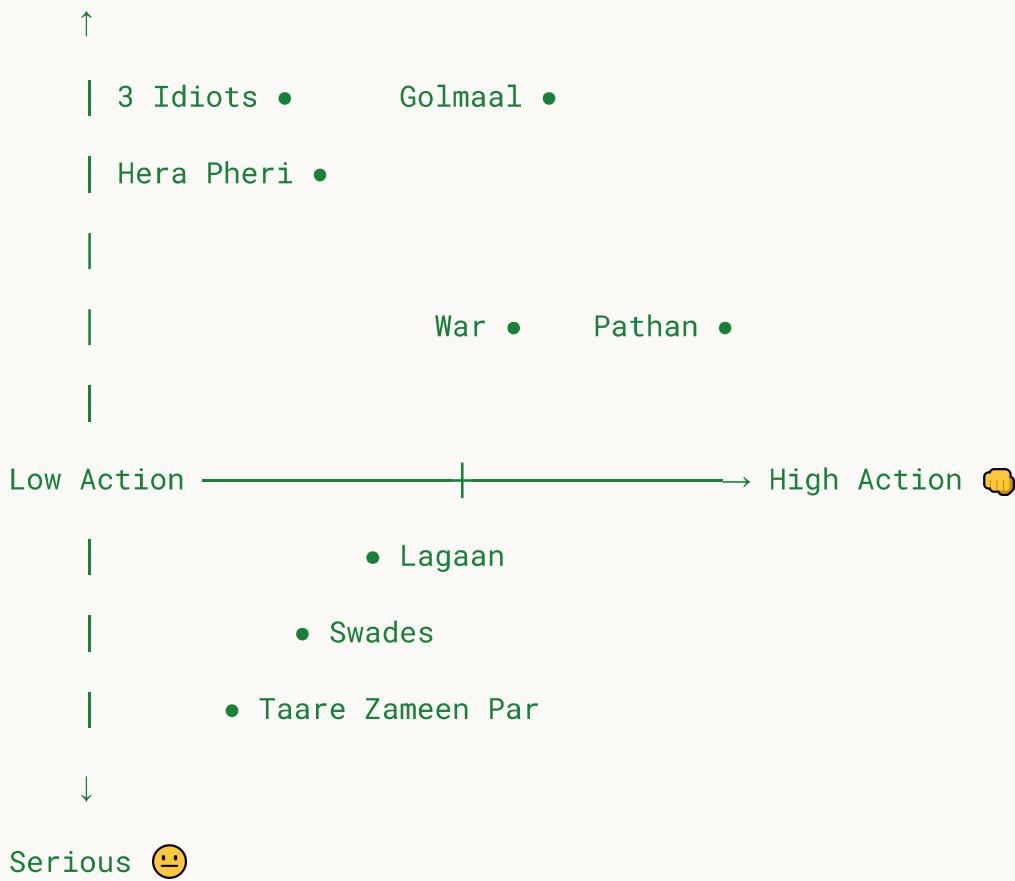
👉 Real-world relationships **multi-dimensional** hote hain, hum **1D** mein solve nahi kar sakte

⭐ The Ultimate Solution: Vectors & Embeddings

◆ Movie Recommendation Example - 2D Space

🎬 MOVIES IN 2D SPACE:

High Comedy 😂



⌚ X-axis: Action Level (-10 = Peaceful, +10 = High Action)

🎭 Y-axis: Comedy Level (-10 = Serious, +10 = High Comedy)

➡ User ne 3 Idiots dekhi → Recommend nearby: Golmaal, Hera Pheri

◆ Multi-Dimensional Vectors (Vector Embeddings)

Real life mein 2 dimensions **kafi nahi hain**

⌚ 5-DIMENSIONAL MOVIE VECTOR:

War Movie = [Action, Comedy, Drama, Romance, Realism]

= [10, -5, -4, 2, -7]

⭐ Yeh list of numbers hi **VECTOR** hai

Modern AI models **496+ dimensions** tak ke vectors bana sakte hain!

◆ Vectors Kaise Create Hote Hain?

Neural Networks **automatically** vectors create karte hain.

Input: "King"

Output: [10, 7, 3, 8, 9, ...]

→ Hum manually numbers assign nahi karte

◆ Famous "King - Man + Woman = Queen" Example

■ VECTOR ARITHMETIC:

King = [Royalty, Male, Power, Wealth] = [10, 10, 9, 10]

Man = [Royalty, Male, Power, Wealth] = [1, 10, 2, 3]

Woman = [Royalty, Female, Power, Wealth] = [1, 0, 2, 3]

King - Man + Woman = [10-1+1, 10-10+0, 9-2+2, 10-3+3]

= [10, 0, 9, 10] ≈ Queen ✓

💡 Why it works:

King - Man = Royalty, Power, Wealth (jo king mein extra hai)

- + Woman = Yeh qualities woman mein add karo to Queen milta hai
-

► Similarity Kaise Measure Karen?

◆ Cosine Similarity vs Euclidean Distance

■ COMPARISON DIAGRAM:

Vector A

| \ | Euclidean Distance

| \ | (Straight line)

| \ |

| θ \ |

| \ \ |

—● Vector B

 **Cosine Similarity** → angle θ ko measure karti hai

Method	What it Does	Best For
--------	--------------	----------

 Euclidean Distance	Distance between points	General geometry
--	-------------------------	------------------

 Cosine Similarity	Angle between vectors (ignores size)	Text & Semantics 
---	---	--

 **Cosine Similarity:**

- `1 = Same direction` (high similarity)
- `0 = No relation`
- `-1 = Opposite`

 **Cosine is best for semantic similarity**

Vector Databases

 **Traditional DB vs Vector DB**

 **Traditional Database:**



 **Vector Database:**



◆ YouTube Search Example

🔍 **User searches:** "What is Array?"

→ Convert to vector → [8, -2, 5, 1, -3, ...]

→ Find closest vectors:

✓ "Array Tutorial" → [7, -1, 6, 2, -2]

✓ "Data Structures" → [8, -3, 4, 1, -4]

✗ "Cooking Recipe" → [-5, 2, -3, 8, 1]

→ YouTube **keywords nahi, meaning match karta hai!**

📚 Complete Summary

🎯 TIMELINE:

Manual Arrays → Graph Weights → Number Lines → VECTOR EMBEDDINGS



Rigid

Co-occurrence

1D Thinking

Multi-dimensional

Static

Only

Limited

Semantic Understanding

◆ Key Definitions

Concept	Definition	Simple Example
Vector	List of numbers in multi-dimensional space	Movie ka fingerprint: [Action, Comedy, Drama]
Vector Embedding	Text/images ko vector mein convert karna	Word ko AI-friendly form mein lana
Vector Database	Similar items dhoondhne wali smart library	YouTube, Google Search, Amazon Recommendations

Final Thought

- ✖ Purana Tarika: "If you bought X, you might like Y" (manual rules)
 - ✓ Naya AI Tarika: "Items that are conceptually similar in high-dimensional space are recommended"
 - Vectors computers ko semantic relationships samajhne ki capability dete hain
-

Visual Recap:

- 1 Manual Arrays →
- 2 Graphs →
- 3 Number Lines →
- 4 Multi-dimensional **Vector Embeddings** ✓

★ Yahi hai secret of modern AI recommendations, search, aur content understanding!

Lecture 07 – What is Vector Database | Internal Implementation of Vector DataBase

(Expert-level, engaging, and complete — by Harshal Chauhan ✨)

1 Introduction – Kyun Chahiye Vector Databases?

Traditional databases (SQL/NoSQL) sirf **exact match** search karte hain,
jabki **Vector Databases** “meaning-based” search karte hain 

Comparison:

-  **SQL/NoSQL:** “Rohit Negi 9 ke comments lao” → Exact Match 
-  **Vector DB:** “Rohit Negi 9 jaisi profiles lao” → Semantic Similarity 

 Example: “What is an array?” ≈ “Array kya hota hai?” → Same Results 

2 Real-World Example – Blinkit Recommendation

User adds “**Onion**” to cart → Vector DB recommends **similar items**

Recommended (Similar Vectors)

- Tomato → [0.8, 0.2, 0.0, 0.8, 0.2]
- Dhaniya → [0.7, 0.1, 0.0, 0.9, 0.0]
- Nimbu → [0.1, 0.8, 0.1, 0.7, 0.3]

Not Recommended (Different Vectors)

- Banana → [0.1, 0.9, 0.1, 0.1, 0.9]
- Protein → [0.0, 0.1, 0.9, 0.1, 0.7]
- Creatine → [0.0, 0.0, 0.95, 0.0, 0.1]

 System meaning ke basis par recommend karta hai, na ki sirf text match pe!



3 Vector Similarity – Core Concept

◆ Formula (Euclidean Distance):

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

🧠 Example:

Onion [2,3] & Tomato [3,4]

$$\rightarrow \sqrt{(3-2)^2 + (4-3)^2} = \sqrt{2} \approx 1.4$$

✖ Rule:

👉 *Jitni chhoti distance, utni zyada similarity!*



4 Brute Force Approach – The Slow Killer ✗

🔍 Query:

Find 10 nearest neighbors of “Onion”

Steps:

1. Onion vector = [0.9, 0.1, 0.0, 0.9, 0.1]
2. Compare with *1 billion* vectors
3. Calculate distance for each
4. Sort results
5. Return top 10

⚠ Result:

- ✅ Accuracy: 100%
- 🐘 Speed: Extremely slow
- 💔 Bad User Experience

💡 *Brute Force = Perfect but Painfully Slow*



5 ANN – Approximate Nearest Neighbor

🎯 Smart Trade-off

Accuracy

Speed

100% →
90–95%

🚀 100x
Faster

Example:

10 recommendations → 9 correct + 1 slightly off = still worth it!

💡 *Speed > Perfection in real-world systems.*

6 ANN Algorithms – Deep Dive

◆ 6.1 IVF – Inverted File Index

Phase 1: Indexing

- Step 1: **Find Centroids** using *K-Means*
- Step 2: Assign each vector to nearest centroid (cluster)

Phase 2: Searching

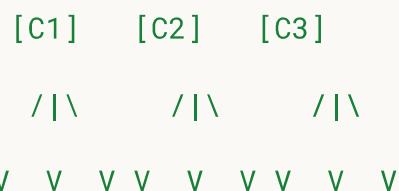
- Query “Grapes” → Compare with 100 centroids
- Choose nearest (say C2) → Search only in that cluster

 100x faster!

Multi-Centroid Strategy:

Query top 3 closest centroids → better accuracy 

Text Figure:



Each C = centroid, connected to its vectors.

6.2 KD-Tree – Binary Space Partitioning

Process:

- Divide data recursively using X & Y splits
- Create a decision tree
- Query travels through branches based on coordinates

Example:

Query Q(13,8)

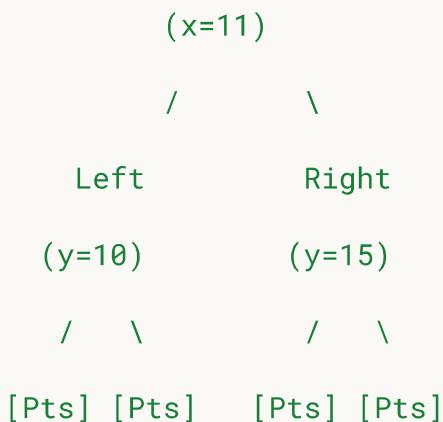
→ Root split at x=11 → Go Right

→ y=10 → Go Left

→ Compare only with nearby nodes

 Efficient for low-dimensional data.

Text Figure:



6. HNSW – Hierarchical Navigable Small Worlds

Concept: Multi-layer graph structure

Layers:

- Layer 0: All nodes connected to 3 nearest
- Layer 1: Random nodes promoted
- Layer 2: Even fewer nodes
- Top Layer: Only 2–3 nodes

Search Process:

1. Start at top layer
2. Find nearest
3. Move layer by layer downward
4. Continue till Layer 0

 **Alpha=3** → Check 3 neighbors per level

 *Fastest and most accurate ANN algorithm.*

Text Figure:

Layer 2: A – B

|

Layer 1: A-C-B-D

|

Layer 0: A-E-C-F-B-G-D

6.4 PQ – Product Quantization

Problem:

1 Billion vectors = 6.1 TB 

Solution (Compression):

- Step 1: Split vector into 12 chunks
- Step 2: Create 256 centroids (codebook)
- Step 3: Replace chunks with centroid indices → only 12 bytes!

Memory Saved:

6.1 TB → 12 GB (\approx 500x reduction!)

PQ Search:

- Split query into chunks
- Compare with precomputed distance tables

 Super-fast lookup!

Text Figure:

Vector [v1 v2 v3 v4 ... v12]

↓ Split into chunks

[] [] [] ... []

↓ Quantized (Centroid IDs)

[05][19][02]...[77]

6.5 IVF + PQ – Hybrid Approach

Combo of:

1. **IVF**: Smart clustering
2. **PQ**: Compression within each cluster

Result:

- High speed
- Low memory
- Great accuracy

 Used widely in billion-scale systems.

⌚ Text Figure:

[Clusters]

C1 → PQ compressed vectors

C2 → PQ compressed vectors

C3 → PQ compressed vectors

📊 7 Performance Comparison

Algorithm	Speed 🚀	Accuracy 🎯	Memory 🗂️	Ideal Use
🏆 HNSW	9/10	10/10	8/10	Maximum Performance
⚡ IVF + PQ	8/10	8/10	4/10	Large Datasets
⚖️ IVF Only	6/10	6/10	5/10	Balanced Needs
🌳 KD-Tree	4/10	5/10	6/10	Low Dimensions
🧩 Brute Force	1/10	10/10	9/10	Small Datasets

🗄️ 8 Vector Database Ecosystem

🧠 Native Vector Databases

Database	Algorithm	Feature
Milvus	HNSW, IVFPQ	Open-source, scalable
Pinecone	HNSW	Managed, high performance
Weaviate	HNSW	GraphQL + Knowledge Graph
Qdrant	HNSW + Quantization	Built in Rust, fast

Integrated Solutions

- PostgreSQL → **pgvector extension**
- Elasticsearch → **Text + Vector search**
- Redis → **In-memory + Vector support**

Foundation Library

- **Faiss (Meta)** → HNSW, IVFPQ (Industry Standard)
-

9 Vector Database Storage Structure

ID → Unique key
Vector → Similarity search
Metadata → Filters (category, price, etc.)

Example:

```
{ id: "product_456", vector: [0.98, 0.23, ...],  
metadata: {"name": "Red Shoes", "price": 89.99} }
```

10 Search Types

1. **Vector Similarity Search** → “Onion jaisi items lao”
2. **Exact ID Search** → “Product_456 lao”
3. **Hybrid Search** → “Onion jaisi items under ₹50 lao”

 Hybrid = Vector + Metadata combo

11 Algorithm Selection Guide

Case → **Recommended**

- | | | |
|--|------------------|------------|
|  | Best Performance | → HNSW |
|  | Low Memory | → IVF + PQ |
|  | Balanced | → IVF Only |
|  | Low Dimension | → KD-Tree |

1 **Ultimate Summary**

Key Insights

1.  Brute Force → Slow
2.  ANN → Fast + Smart
3.  HNSW → Best performance
4.  IVF+PQ → Memory efficient
5.  Hybrid → Practical choice

 Vector DBs *samajhte hain meaning, sirf words nahi!*

1 **Implementation Roadmap**

Type → Tool

Startup	→ Pinecone
Enterprise	→ Milvus
SQL Users	→ PostgreSQL + pgvector
Developers/Researchers	→ Faiss

1 **Final Takeaway**

- ◆ **Traditional DBs** → Exact Match
- ◆ **Vector DBs** → Semantic Understanding
- ◆ **ANN Algorithms** → Fast + Smart
- ◆ **Hybrid Systems** → Best Real-World Performance

 *In short:*

“Vector DBs think like humans — they understand meaning, not just text.” 

Lecture 08 – Retrieval Augmented Generation (RAG) Complete Guide

◆ RAG Ki Basic Samajh

RAG Kya Hai?

Full Form: Retrieval Augmented Generation

Simple Definition:

RAG ek aisa system hai jo **user ke question ka answer** dene ke liye pehle **apne documents se relevant information retrieve** karta hai,

phir us **context ko LLM (Large Language Model)** ko data hai, jisse woh **accurate aur relevant answer generate** kar sake.

Example (from YouTube Video)

Agar aap LLM ko bolo: “Meri girlfriend ki tarah behave karo”

aur phir pucho: “Tumhe kahan ghumna pasand hai?”

Toh LLM bolega — “Chicago” 😅

Par asli girlfriend Uttarakhand se hai!

Yeh galat output hai → **Hallucination**

RAG Kyu Zaroori Hai?

-  **Hallucination Problem:** LLM ko answer nahi pata hota, fir bhi kuch bhi bata dete hai
 -  **Limited Context:** LLM ke paas aapke personal documents ka knowledge nahi hota
 -  **Cost Effective:** Fine-tuning se kaafi sasta hota hai
 -  **Dynamic:** Naye documents easily add kar sakte hain
-

◆ RAG vs Fine-Tuning – Complete Comparison

Fine-Tuning

Definition: LLM ko naye data par dobara train karna.

Example: WhatsApp chat dekar LLM ko bolna — “Is data par train ho jao”.

Disadvantages:

- Bahut expensive (GPU cost high)
 - Time-consuming
 - Har baar naya data aane par retraining karni padti hai
-

RAG

Definition: Documents ko context ke roop mein dena.

Example: WhatsApp chat ko context ke roop mein LLM ko dena — LLM wahi se answer karega.

Advantages:

- Fast & Cheap
 - Real-time document updates possible
 - No training required
-

Comparison Table

Aspect	Fine-Tuning	RAG
Cost	High (GPU, Training)	Low (API Calls)
Time	Days/Weeks	Minutes/H ours
Flexibility	Low	High
Maintenance	Difficult	Easy

Analogy:

Fine-Tuning = Student ko poori book ratta marwana

RAG = Student ko exam hall mein book le jaane dena

◆ RAG Ka Complete Flow (2 Steps)

🚀 Step 1: Indexing Phase (Ek Baar Ka Kaam)

🎯 Goal: Documents ko Vector Database mein store karna

1. 📄 Document Loading

```
const loader = new PDFLoader("dsa.pdf");  
  
const rawDocs = await loader.load();
```

💡 Example: 112 pages ki PDF load hoti hai.

2. ✂️ Chunking (Tukdon Mein Katna)

Kyu? – Taki baad mein sirf relevant parts retrieve ho sake.

Chunk Size: 1000 words

Overlap: 200 words

💡 Overlap Importance:

Agar “Linked List” ki definition do chunks mein split ho jaye, to overlap se dono chunks mein kuch context bacha rahega.

```
const splitter = new RecursiveCharacterTextSplitter({  
  
  chunkSize: 1000,  
  
  chunkOverlap: 200  
  
});  
  
const chunkDocs = await splitter.splitDocuments(rawDocs);
```

📊 Example: 112 pages → 227 chunks

3. 🖥️ Vector Conversion

Definition: Text → Numbers (numerical representation)

Model: Google Generative AI – `text-embedding-004`

Dimensions: 768

```
const embeddings = new GoogleGenerativeAIEmbeddings({  
  
  model: "text-embedding-004"  
  
});
```

4. 📁 Vector Database Storage

Database: Pinecone

📌 Setup Steps:

1. [pinecone.io](#) par sign up karo
2. Index create karo (768 dimensions)
3. API key & environment le lo

```
const pinecone = new Pinecone();

const index = pinecone.Index("rohit-neg");

await PineconeStore.fromDocuments(chunkDocs, embeddings, {
    pineconeIndex: index,
    maxConcurrency: 5
});
```

⌚ 5 chunks ek saath process honge.

🚀 Step 2: Query Phase (User Interaction)

🎯 **Goal:** User ke question ka accurate answer dena

1. 🤔 User Input Lena

```
const userQuestion = "What is AVL tree?";
```

2. 📊 Query Ko Vector Mein Convert Karna

```
const queryVector = await embeddings.embedQuery(userQuestion);
```

3. 🔎 Vector Database Mein Search Karna

```
const results = await index.query({
    vector: queryVector,
    topK: 10,
    includeMetadata: true
});
```

4. Context Create Karna

```
const context = results.matches.map(match =>  
  match.metadata.text  
).join("\n\n");
```

5. LLM Ko Context + Question Dena

```
const model = new GoogleGenerativeAI({ model: "gemini-pro" });  
  
const response = await model.generateContent({  
  contents: [  
    {  
      role: "user",  
      parts: [{  
        text: `You are a DSA expert. Answer based ONLY on this  
context:  
  
Context: ${context}  
  
Question: ${userQuestion}  
  
If answer is not in context, say "I couldn't find the answer in provided  
documents".`  
      }]  
    }  
  ]  
});
```

◆ Advanced RAG Concepts

🎯 Query Transformation Problem

✗ Problem:

User: "What is Quick Sort?" → LLM correct answer data hai

User: "Explain it in detail" → ✗ "I couldn't find the answer"

Reason:

"Explain it in detail" ka vector different hai,
toh DB ko relevant documents nahi milte.

✓ Solution: Query Transformation Layer

```
async function transformQuery(question, history) {  
  
    const model = new GoogleGenerativeAI({ model: "gemini-pro" });  
  
    const response = await model.generateContent({  
  
        contents: [  
  
            {  
                role: "user",  
  
                parts: [{  
  
                    text: `You are a query rewriting expert.  
  
Based on chat history, rewrite this follow-up question into a standalone  
question.  
  
Chat History: ${history}  
  
Latest Question: ${question}`  
  
                }]  
  
        ]  
  
    });  
  
    return response.text();  
}
```

Example:

“Explain it in detail” → “What is Quick Sort in depth?”

Token Cost Optimization

Problem: Poora document dene se token usage zyada hota hai.

Solution: Sirf **relevant chunks** bhejo.

 Example: Agar “Course Revenue” puchha jaaye, toh sirf revenue related chunks bhejo, poori PDF nahi.

◆ Complete Project Implementation

File Structure

```
project/
├── index.js      # Indexing phase
├── query.js      # Query phase
├── .env          # API keys
└── dsa.pdf       # Input document
```

Dependencies

```
npm install @langchain/community @langchain/google-genai pinecone-client
pdf-parse
```

Indexing Code (index.js)

```
// Load PDF

const loader = new PDFLoader("dsa.pdf");

const rawDocs = await loader.load();

// Split into Chunks

const splitter = new RecursiveCharacterTextSplitter({
```

```
    chunkSize: 1000,  
    chunkOverlap: 200  
});  
  
const chunkDocs = await splitter.splitDocuments(rawDocs);  
  
  
// Create Embeddings  
  
const embeddings = new GoogleGenerativeAIEmbeddings({  
  model: "text-embedding-004",  
  apiKey: process.env.GOOGLE_API_KEY  
});  
  
  
// Connect Pinecone  
  
const pinecone = new Pinecone({  
  apiKey: process.env.PINECONE_API_KEY  
});  
  
  
// Store in Vector DB  
  
await PineconeStore.fromDocuments(chunkDocs, embeddings, {  
  pineconeIndex: pinecone.Index("rohit-neg"),  
  maxConcurrency: 5  
});
```

Querying Code (query.js)

```
const userQuestion = "What is AVL tree?";

const queryVector = await embeddings.embedQuery(userQuestion);

const results = await index.query({
    vector: queryVector,
    topK: 10,
    includeMetadata: true
});

const context = results.matches.map(m => m.metadata.text).join("\n\n");

const response = await model.generateContent({
    contents: [ {
        role: "user",
        parts: [ {
            text: `You are a DSA expert. Answer using ONLY this context:
${context}`
        }]
    }]
});

});
```

◆ .env File Setup

```
GOOGLE_API_KEY=your_google_api_key

PINECONE_API_KEY=your_pinecone_api_key

PINECONE_ENVIRONMENT=your_pinecone_environment
```

◆ Common Issues & Solutions

Problem	Cause	Solution
Vague Questions	Context missing	Use query transformation
High Token Cost	Full doc sent	Send only relevant chunks
Wrong Answers	LLM using own data	Restrict: “Use only provided context”
Slow Speed	Too many chunks	Adjust <code>maxConcurrency</code>

◆ Testing & Validation

Test Cases

- “What is AVL Tree?” → Correct Answer
- “Explain it in detail” → Works after transformation
- “Who is Rohit Negi?” → “I couldn’t find the answer”
- “Quick Sort time complexity?” → Exact numerical result

Validation

- Pinecone dashboard → vectors visible
- Context retrieval → relevant
- LLM output → only context-based

◆ Production-Ready Tips

💡 Performance Optimization

- Chunk Size → 500–1000 words
- Top K → 5–10 results
- Concurrency → 3–5 (for free tier)

💡 Cost Optimization

- Use free tiers (Google AI, Pinecone)
- Monitor token usage
- Implement caching for frequent queries

💡 Error Handling

```
try {  
    // RAG code  
}  
catch (err) {  
    console.error("RAG Error:", err);  
}
```

◆ Real-World Examples

🎯 1. Girlfriend Chatbot

- Problem: LLM doesn't know personal info
- Solution: WhatsApp chat as context
- Result: Personalized accurate responses

🎯 2. Coder Army Revenue

- Problem: LLM doesn't know company revenue
- Solution: Give sales data as context
- Result: Exact revenue fetched

3. DSA PDF Q&A

- Problem: LLM lacks DSA knowledge
 - Solution: DSA PDF in vector DB
 - Result: Contextual, accurate DSA answers
-

RAG System – Final Summary

Concept	Explanation
 Retrieval	Query se related documents dhoondhna
 Augmentation	Un documents ko context ke roop me dena
 Generation	Context ke basis par answer banana

Why RAG?

- Hallucination reduction 
- Cost-effective 
- Real-time document updates 
- Domain-specific knowledge injection 

How It Works

1. **Indexing:** Documents → Chunks → Vectors → Vector DB
2. **Querying:** Question → Vector → Search → Context + Question → LLM → Answer

Advanced Features

- Query Transformation
- Multi-Document Handling
- Token Cost Optimization
- History-Aware Querying

Tools Used

-  Vector DB → Pinecone
 -  Embeddings → Google Generative AI
 -  Framework → LangChain
 -  LLM → Gemini Pro
-

Final Outcome

RAG system banakar:

- LLM hallucination door hoti hai
- Context-based, accurate, real-world responses milte hain
- System scalable, dynamic aur production-ready hota hai 

Conclusion:

RAG = LLM ka “Exam Hall Companion” 

Jo use sirf “sahi notes” dekar perfect answer likhne mein madad karta hai! 