

# The Key Terminologies

## What is an LLM (Large Language Model)?

- A **Large Language Model (LLM)** is an AI model trained on vast amounts of text to generate human-like responses.
- Examples: **GPT (OpenAI), DeepSeek, Claude (Anthropic), Mistral**
- LLMs **understand, process, and generate text** by predicting the next word in a sequence based on probability.

## DeepSeek: The Rising Star of LLMs

DeepSeek is a **powerful open-source LLM** designed to rival models like GPT and Mistral.

### ◆ Why is DeepSeek unique?

- ✓ Open-source & community-driven
- ✓ Strong capabilities in **coding, reasoning, and problem-solving**
- ✓ Can run **locally** without internet dependency
- ✓ Supports **fine-tuning & customization** for specific tasks

## Key Terminologies You Need to Know

### 1 Tokens & Tokenization

- **Token:** A piece of text (word, subword, or character) that an LLM processes.
- **Example:** "DeepSeek is powerful" → Might be split as ["Deep", "Seek", "is", "powerful"]
- **Tokenization:** The process of breaking text into tokens for AI models to process.

### 2 Parameters

- The **"brain size"** of an AI model, determining its ability to learn and process information.
- DeepSeek, GPT, and Mistral have different numbers of **parameters** affecting their power.

- **More parameters = More intelligence (but also more computing power required).**

### 3 Training & Fine-tuning 🏆

- **Pre-training:** The initial phase where an LLM learns from large datasets.
- **Fine-tuning:** A process where an AI model is trained on specific tasks (e.g., medical, finance, coding).
- **Example:** Fine-tuning DeepSeek for Python programming assistance.

### 4 Inference vs. Training ⚡

- **Training:** Building the AI model by processing huge amounts of data.
- **Inference:** Using the trained model to generate text responses.
- **Example:** When you type a question in ChatGPT and get a response → That's **inference!**

### 5 Latency & Speed 🚗

- **Latency:** The delay between input (question) and output (AI response).
- **Optimizing for low latency** is crucial for real-time AI applications like chatbots.

## 🖥️ DeepSeek vs. OpenAI GPT vs. Mistral (Quick Comparison)

Feature	DeepSeek ✅	GPT (OpenAI) 🧠	Mistral 🔥
<b>Open Source?</b>	✅ Yes	❌ No (Closed)	✅ Yes
<b>Local Setup?</b>	✅ Yes	❌ No (Cloud Only)	✅ Yes
<b>Fine-Tuning?</b>	✅ Yes	❌ Limited	✅ Yes
<b>Best For</b>	Custom AI Apps, Local Use	General AI Chat, API-based Apps	Custom AI Models
<b>Cost</b>	💰 Free (Self-hosted)	💰 Paid API	💰 Free & Paid

### 💡 Why Learn DeepSeek?

- **Run AI locally 🏠** without needing OpenAI APIs.
- **Develop custom AI-powered apps 🖥️** using Python & Gradio.
- **Build chatbots, coding assistants, and automation tools 🤖.**

- **Fine-tune models for specific tasks** 🎯 without relying on cloud services.