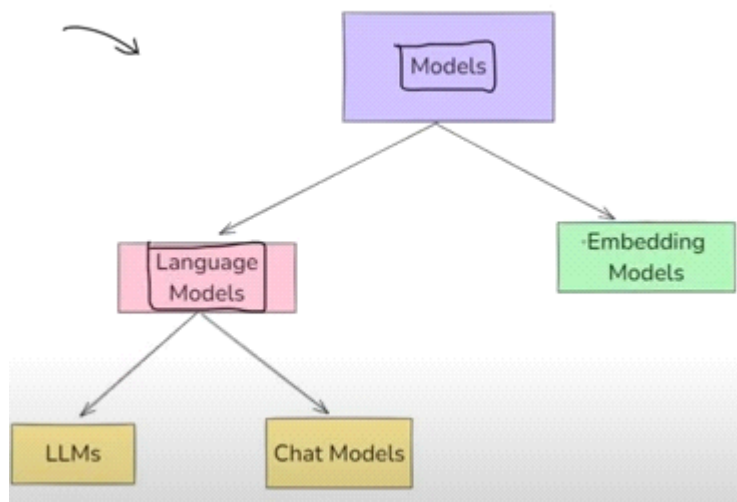# LangChain Models

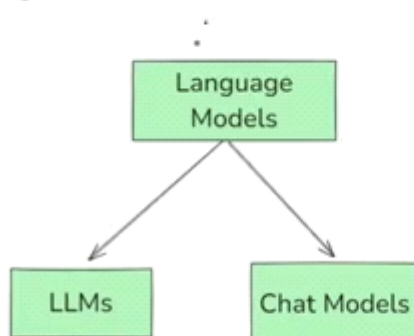22 May 2025        22:24



## What are Models

07 January 2025        23:15

The Model Component in LangChain is a crucial part of the framework, designed to facilitate interactions with various language models and embedding models.

It abstracts the complexity of working directly with different LLMs, chat models, and embedding models, providing a uniform interface to communicate with them. This makes it easier to build applications that rely on AI-generated text, text embeddings for similarity search, and retrieval-augmented generation (RAG).

## Language Models

Language Models are AI systems designed to process, generate, and understand natural language text.



LLMs - General-purpose models that is used for raw text generation. They take a string(or plain text) as input and returns a string( plain text). These are traditionally older models and are not used much now.

Chat Models - Language models that are specialized for conversational tasks. They take a sequence of messages as inputs and return chat messages as outputs (as opposed to using plain text). These are traditionally newer models and used more in comparison to the LLMs.

| Feature | LLMs (Base Models) | Chat Models (Instruction-Tuned) |
|---|---|---|
| Purpose | Free-form text generation | Optimized for multi-turn conversations |
| Training Data | General text corpora (books, articles) | Fine-tuned on chat datasets (dialogues, user-assistant conversations) |
| Memory & Context | No built-in memory | Supports structured conversation history |
| Role Awareness | No understanding of "user" and "assistant" roles | Understands "system", "user", and "assistant" roles |
| Example Models | GPT-3, Llama-2-7B, Mistral-7B, OPT-1.3B | GPT-4, GPT-3.5-turbo, Llama-2-Chat, Mistral-Instruct, Claude |
| Use Cases | Text generation, summarization, translation, creative writing, code generation | Conversational AI, chatbots, virtual assistants, customer support, AI tutors |

`temperature` is a parameter that **controls the randomness of a language model's output**. It affects how **creative or deterministic** the responses are.

- **Lower values** ($0.0 - 0.3$) → More **deterministic** and predictable.
- **Higher values** ($0.7 - 1.5$) → More **random**, creative, and diverse.

| Use Case | Recommended Temperature |
|---|---|
| Factual answers (math, code, facts) | $0.0 - 0.3$ |
| Balanced response (general QA, explanations) | $0.5 - 0.7$ |
| Creative writing, storytelling, jokes | $0.9 - 1.2$ |
| Maximum randomness (wild ideas, brainstorming) | $1.5+$ |

## Open Source Models

11 February 2025    08:58

Open-source language models are freely available AI models that can be downloaded, modified, fine-tuned, and deployed without restrictions from a central provider. Unlike closed-source models such as OpenAI's GPT-4, Anthropic's Claude, or Google's Gemini, open-source models allow full control and customization.

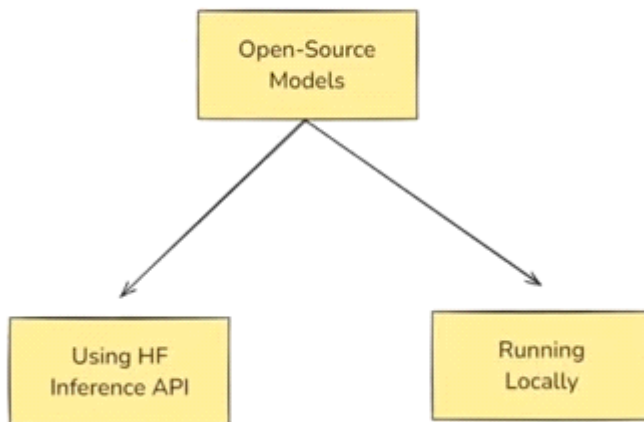| Feature | Open-Source Models | Closed-Source Models |
|---|---|---|
| Cost | Free to use (no API costs) | Paid API usage (e.g., OpenAI charges per token) |
| Control | Can modify, fine-tune, and deploy anywhere | Locked to provider's infrastructure |
| Data Privacy | Runs locally (no data sent to external servers) | Sends queries to provider's servers |
| Customization | Can fine-tune on specific datasets | No access to fine-tuning in most cases |
| Deployment | Can be deployed on **on-premise** servers or cloud | Must use vendor's API |

## Some Famous Open Source Models

| Model | Developer | Parameters | Best Use Case |
|---|---|---|---|
| LLaMA-2-7B/13B/70B | Meta AI | 7B - 70B | General-purpose text generation |
| Mixtral-8x7B | Mistral AI | 8x7B (MoE) | Efficient & fast responses |
| Mistral-7B | Mistral AI | 7B | Best small-scale model (outperforms LLaMA-2-13B) |
| Falcon-7B/40B | TII UAE | 7B - 40B | High-speed inference |
| BLOOM-176B | BigScience | 176B | Multilingual text generation |
| GPT-J-6B | EleutherAI | 6B | Lightweight and efficient |
| GPT-NeoX-20B | EleutherAI | 20B | Large-scale applications |
| StableLM | Stability AI | 3B - 7B | Compact models for chatbots |

Where to find them?

**HuggingFace** - The largest repository of open-source LLMs

Ways to use Open-source Models

```
        ┌──────────────┐
        │ Open-Source  │
        │   Models     │
        └──────┬───────┘
         ╱            ╲
 ┌──────────────┐  ┌──────────────┐
 │  Using HF    │  │  Running     │
 │ Inference API│  │  Locally     │
 └──────────────┘  └──────────────┘
```

## Disadvantages

| Disadvantage | Details |
|---|---|
| High Hardware Requirements | Running large models (e.g., LLaMA-2-70B) requires expensive GPUs. |
| Setup Complexity | Requires installation of dependencies like **PyTorch, CUDA, transformers.** |
| Lack of RLHF | Most open-source models don't have **fine-tuning with human feedback,** making them weaker in instruction-following. |
| Limited Multimodal Abilities | Open models don't support **images, audio, or video** like GPT-4V. |