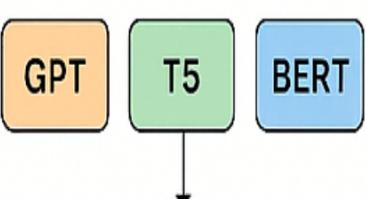


Language Models (LLMs)



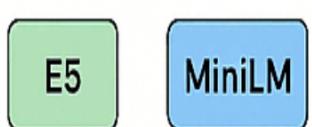
Chat Models

Dialogue Tasks

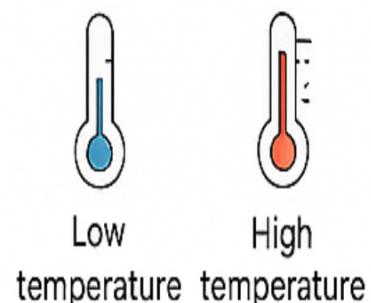


Embedding Models

Sentence Transformers



Temperature
In randomness
and creativity in response



VS

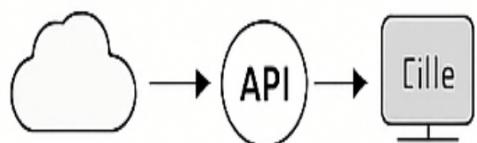
Open vs. Closed Source

Randomness

Random Cost

GPT-3	GPT-Neo
GPT-4	GPT-J
Advantages	Advantages
Access	Customization
Disadvantages	Cost reflective
Hardware requirements	Disadvantages
	Hardware requirement

How to Use Models

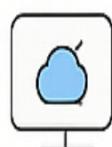


Closed Models

Utilize

Embedding Models Usage

Closed Source



OpenAI Embeddings

Open Source



Hugging Face

Open Source

GPT-3 GPT-J

closed

Via

API

Affordable

Cost-effective

ADVANTAGES

Access-Cost

Disadvantages

DISADVANTAGES

Hardware requirements

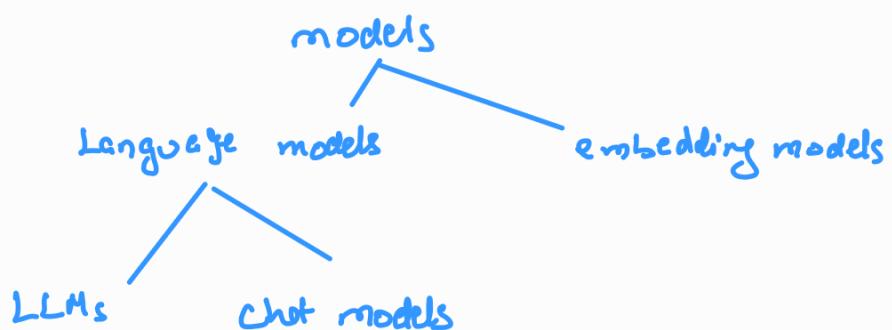
Text Generation Models Usage

• Text similarity

• Other

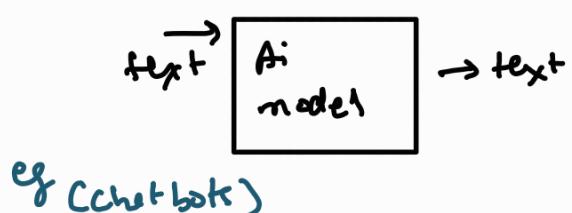
Models:

- 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)



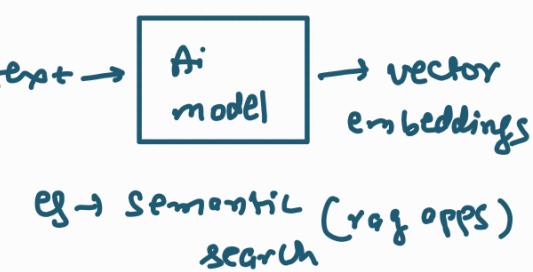
Language models:

input → text
output → text

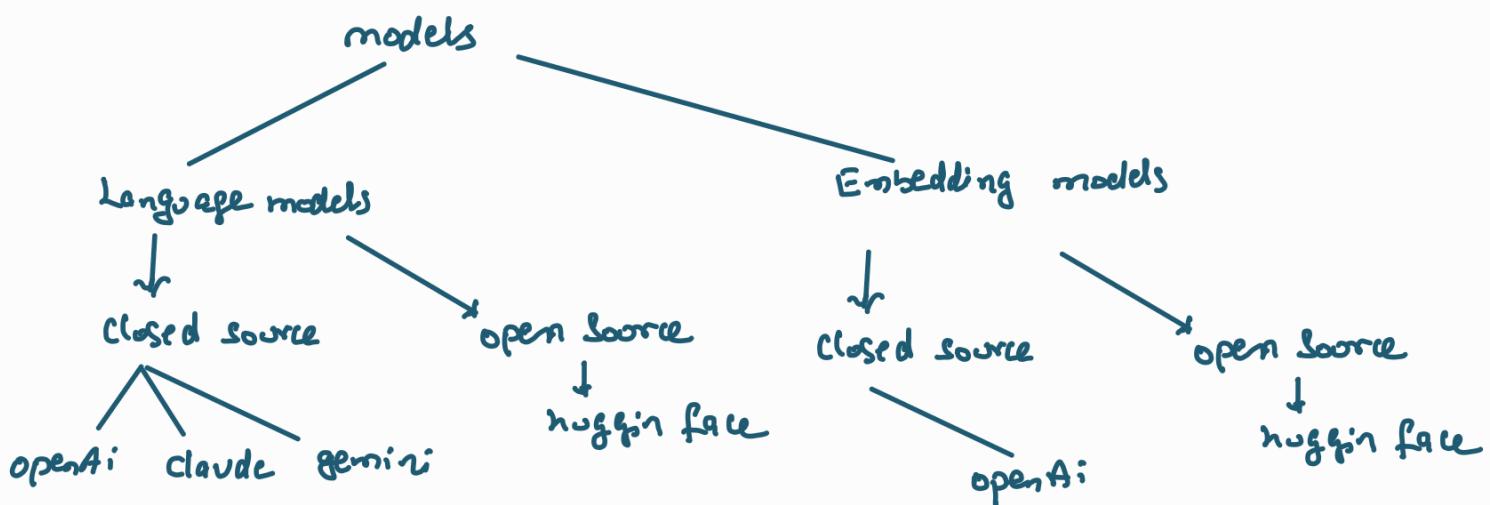


Embedding models:

input → text
output → embeddings / vectors

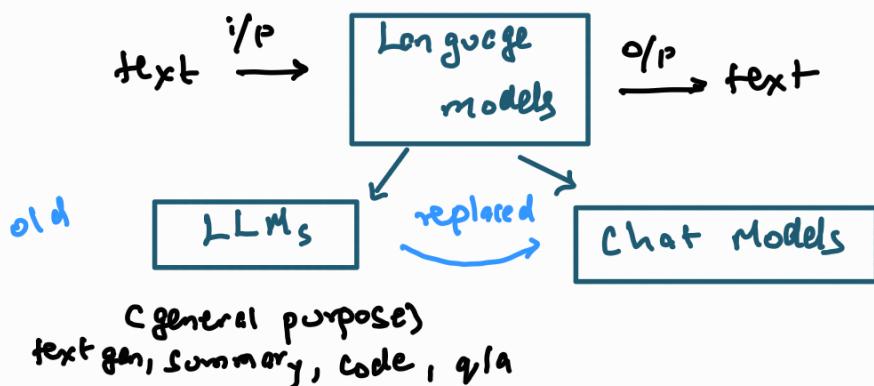


Plan of action



Language models :

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



LLMs :

- General purpose models that are used for raw-text generation.
- takes a string (or plain text as input) and returns a string (or plain text as output)
- 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 message as output (as opposed to using plain text)

feature	LLMs (Base Models)	ChatModel & InstructionModel
Purpose	free form text generation	optimised for multi-turn conversation
Training data	General text corpora (books, articles)	fine tuned on chat databases (dialogues, user-agent: start conversation)
Memory & context	no built-in memory	Supports structured conversation history
Role awareness	no understanding of user and assistant	understands "System", "user" and "assistant" roles
Example models	GPT-3, llama-2-7B, mistral 7B	GPT-4, GPT 3.5, llama-2-chat-mistral, instruct
use cases	Text generation, summarization Creative writing, code generation	Conversational AI, chatbots, virtual assistants, customer support, AI tutors

Set up:

- i) create a file
- ii) open in vs code
- iii) create a new venv
- iv) `python -m venv venv`

iv) activate

venv\scripts\activate

v) create requirements.txt

vi) pip install -r requirements.txt

vii) verify langchain installation.

Requirements.txt

#langchain

langchain

langchain-core

#openai

langchain-openai

openai

#Anthropic

langchain-anthropic

#google

langchain-google-gptai

langchain-generativeai

#huggingface

langchain-huggingface

transformers

huggingface-hub

#ML utils

numpy

scikit-learn

env management

python-dotenv

test

```
import langchain  
print(langchain.__version__)
```

LLMs

chat models

Embedding Models.

1) LLMS

↳ create account

↳ buy \$ & credits

↳ create an API key.

Demo [do not use this ↓] Cuz old , use chat models instead.

```
from langchain.openai import OpenAI  
from dotenv import load_dotenv
```

Load - default

```
lm = OpenAI(model = "gpt-3.5-turbo-instruct")
```

```
result = lm.invoke("what is the capital of India")
```

→ the capital of india is New Delhi

Conclusion

→ input is string

→ output is string

ChatModel

OpenAI → Models

```
from langchain - openai import ChatOpenAI
from dotenv import load_dotenv
load_dotenv()
```

model = ChatOpenAI(model = "gpt-4", temperature=0)

result = model.invoke("what is the capital of India")

print(result) → returns an object to object

print(result.content) → contains the output

→ the capital of India is New Delhi

input - text

output: object {object}

max_completion_tokens = 10
words

used to limit the number of tokens to be received.

#temperature (0 - 2)

defines how random the output will be compared to previous output

→ 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 & predictable
- higher values (0.7 - 1.5) → more random, creative, diverse

use cases:

recommended temp

factual answers (math, code) → 0.0 - 0.3

Balanced (Q/A, explanations) → 0.5 - 0.7

Creative writing, Storytelling → 0.9 - 1.2

Maximum randomness (wild ideas, brainstorming) → 1.5 +

Anthropic - Claude API use.

```
from langchainAnthropic import ChatAnthropic
```

```
from dotenv import load_dotenv
```

```
Load_dotenv()
```

```
model = ChatOpenAI(model="claude-2-sa-2022-04-1022")
```

```
result = model.invoke("what is the capital of India")
```

```
print(result) → Returns an object to object
```

```
print(result.content) → contains the output
```

∴ the capital of India is New Delhi

Google Gemini API

```
from langchainGoogleGemini import ChatGoogleGenerativeAI
```

```
from dotenv import load_dotenv
```

```
Load_dotenv()
```

```
model = ChatGoogleGenerativeAI(model="gemini-2.5-pro", temperature=1)
```

```
result = model.invoke("what is the capital of India")
```

```
print(result) → Returns an object to object
```

```
print(result.content) → contains the output
```

∴ the capital of India is New Delhi

Open Source Models

→ OS models language models are freely available AI models that can be downloaded, modified and deployed without any restrictions from a central provider.

→ Unlike closed source models such as OpenAI GPT-4 or Google Gemini, open source models allow full control and customization

feature	open-source model	closed source model
cost	free	paid API usage
control	modify, fine-tune deploy anywhere	locked to provider's infrastructure
Data privacy	run locally (no data sent to external servers)	sends queries to servers
Customization	can fine tune on specific datasets	very limited to no access to fine tuning up dataset
Deployment	can be deployed on-premise servers or cloud	must use vendor API

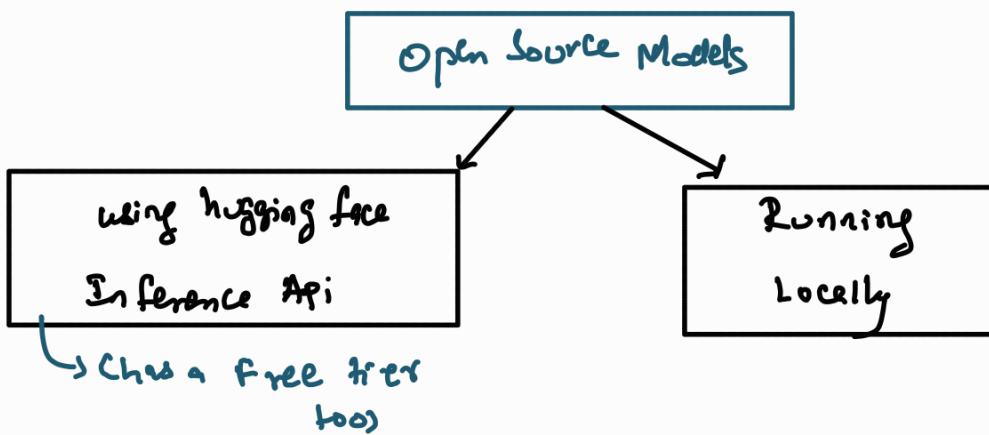
Some famous Open Source Models

- 1) Llama
- 2) bloom
- 3) mistral

Where to find Open Source AI models?

- i) huggingface → largest repository of open-source LMs
- ii) Open roster → consists of both open source and closed source LMs

how to use?



Disadvantages

high hardware requirements : running a large model requires GPUs

setup complexity : requires installation of dependencies like pytorch, cuda, transformers

Lack of RLHF : most open source models don't have fine tuning with
(↳ reinforcement learning from human feedback)

human feedback , making them weaker in instruction-following

Limited Multimodal
Abilities: } open model's do not support images, videos
} audios like GPT-4 V

Using Open Source Models using Hugging Face Inference API

```
from langchain_huggingface import ChatHuggingface, HuggingFaceEndPoint
from dotenv import load_dotenv
load_dotenv()

lm = HuggingfaceEndpoint(
    repository_id="tiiuae/tiiuae-1.1B-chatv1.0",
    task="text-generation"
)
model = ChatHuggingface(lm=lm)
result = model.invoke("What is the capital of India")
print(result.content)
```

Using local LLM

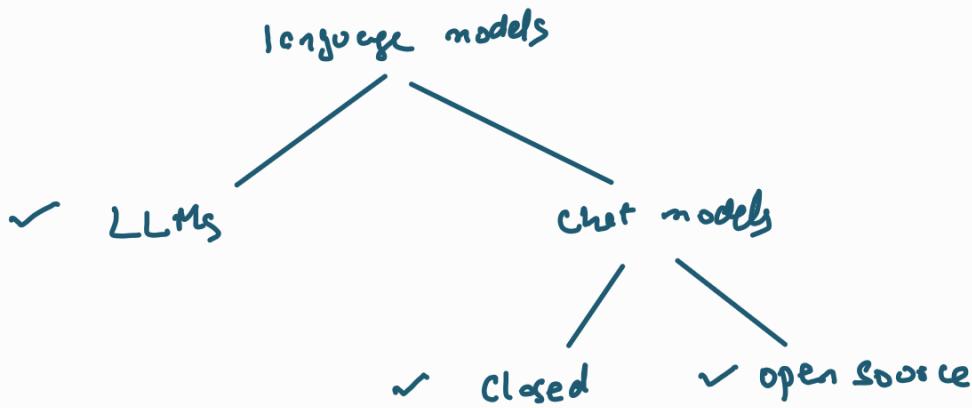
```
from langchain_huggingface import ChatHuggingface, HuggingFacePipeline
from dotenv import load_dotenv
load_dotenv()

lm = HuggingfaceEndpoint(
    repository_id="tiiuae/tiiuae-1.1B-chatv1.0",
    task="text-generation", pipeline_kwargs=(temperature=1,
    max_new_tokens=100))
```

```

model = ChatHuggingface(lm=llm)
result = model.invoke("what is the capital of India")
print(result.content)

```



Embedding Models:

- used to convert words (texts) to vectors, in a way
- these vectors are usually continuous values and capture semantic relationship between the data.
- the goal is to represent the data in a way that preserves its important characteristics while simplifying it into a more manageable form.

common types of embeddings..

- i) word Embeddings (e.g. word2vec, glove)
- ii) Sentence or document Embeddings : e.g (sentence-BERT)
- iii) Image Embeddings : e.g (Resnet, ViT)
- iv) Multimodal Embeddings : e.g (CLIP)

Note: The larger the dimension, the more semantic understanding (contextual meaning)

Sample code

```
from langchain.openai import OpenAIEmbeddings
from dotenv import load_dotenv
load_dotenv()

embeddings = OpenAIEmbeddings(model="text-embedding-gecko",
                               dimension=4)

result = embeddings.embed_query("Delhi is the capital of India")
print(str(result))

→ can pass a single sentence.
```

→ [-0.294366 , 0.501512 , -0.01360584 , 0.213417]

ii) Embedding-docs

```
from langchain.openai import OpenAIEmbeddings
from dotenv import load_dotenv
load_dotenv()

embeddings = OpenAIEmbeddings(model="text-embedding-gecko",
                               dimension=4)

docs = [ "Delhi is the capital of India",
         "Kolkata is the capital of West Bengal",
         "Chennai is the capital of Tamil Nadu" ]

result = embeddings.embed_documents(docs)
```

```
print(str(result))
```

>>

```
[  
    [-0.299366, 0.501512, -0.01360584, 0.2134177],  
    [-0.299322, 0.501514, -0.01360588, 0.2134178],  
    [-0.299367, 0.501516, -0.01360586, 0.2134179]  
]
```

Embedding using Open Source Model.

(Sentence Transformer)

```
from langchain_huggingface import HuggingFaceEmbeddings
```

```
embedding = HuggingFaceEmbeddings(model="sentence-transformers/all-MiniLM-L6-V2")
```

```
text = "Delhi is the capital of India"
```

```
docs = ["delhi is the capital of India",  
        "Kolkata is the capital of West Bengal",  
        "Chennai is the capital of Tamil Nadu"]
```

```
vector = embedding.embed_documents(text)
```

```
vector = Embeddings_Embbed_documents(docs)
```

```
print(vector)
```

```
print([vector])
```

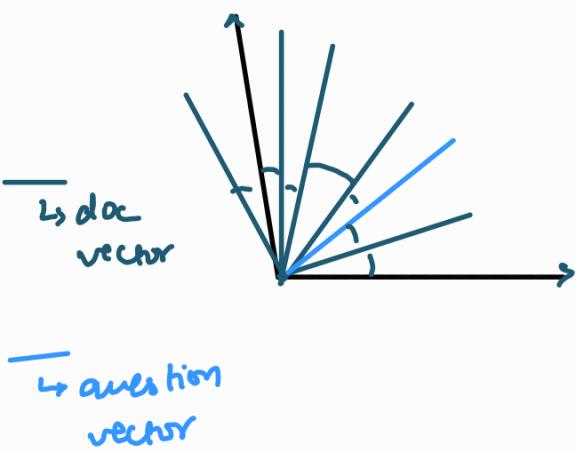
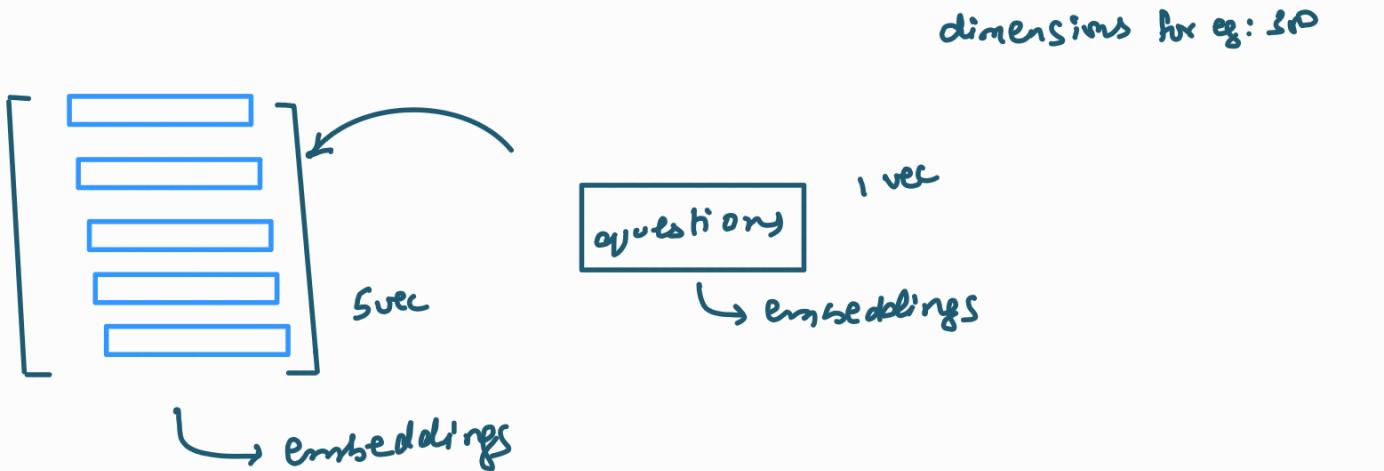
>>

O/P
384
dimension vectors.



Document Similarity Application.

5 documents



Cosine similarity. [a part of scikit learn metric]

(the one with highest similarity score is the answer)

[the same concept is used in RAG applications]

→ let's say we have 4 documents and 1 question related to the document, and we need to find which document it is.

→ how do we do that??

- 1) generate Embeddings for documents and store it
- 2) generate Embeddings for the query and store it
- 3) find the cosine-similarity between the documents and query.

- 4) The vector with highest cosine similarity
- 5) map it with the document to find the document.

Code

```
from sklearn.metrics.pairwise import cosine_similarity
from langchain.openai import OpenAIEmbddings
from dotenv import load_dotenv

load_dotenv()

embedding = OpenAIEmbddings(model="text-embedding-3-large",
                             dimension=300)
```

documents: [

- "Sachin is the God of cricket",
- "Virat is an Aggressive batsmen",
- "Jasprit is an Amazing bowler",
- "Rohit has lots of double centuries"]

query = "tell me about Virat"

```
documentEmbeddings = embedding.embed_documents(documents)

queryEmbedding = embedding.embed_query(query)

scores = cosine_similarity([queryEmbedding], documentEmbeddings)[0]
```

```
index, score = sorted(list(zip(range(len(scores)), scores)), key=lambda x: x[1])[-1]
```

print(f"Document {index+1})

print(f"Similarity Score is {score}")

27 Vint is an Aggressive bkhnan
Similarity score is 0.62809