

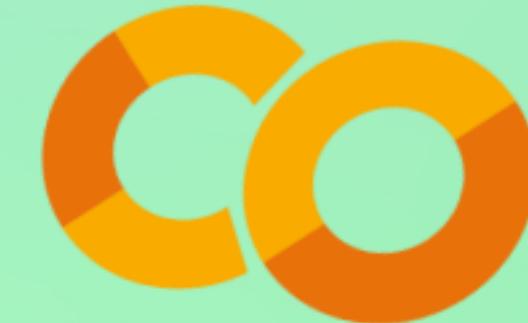
LLM Based Medical Support Chatbot using LangGraph



Hugging Face



LangChain



Tavily

together.ai

LangChain

LLaMa-3-70B

LangGraph

Prompt Engineering

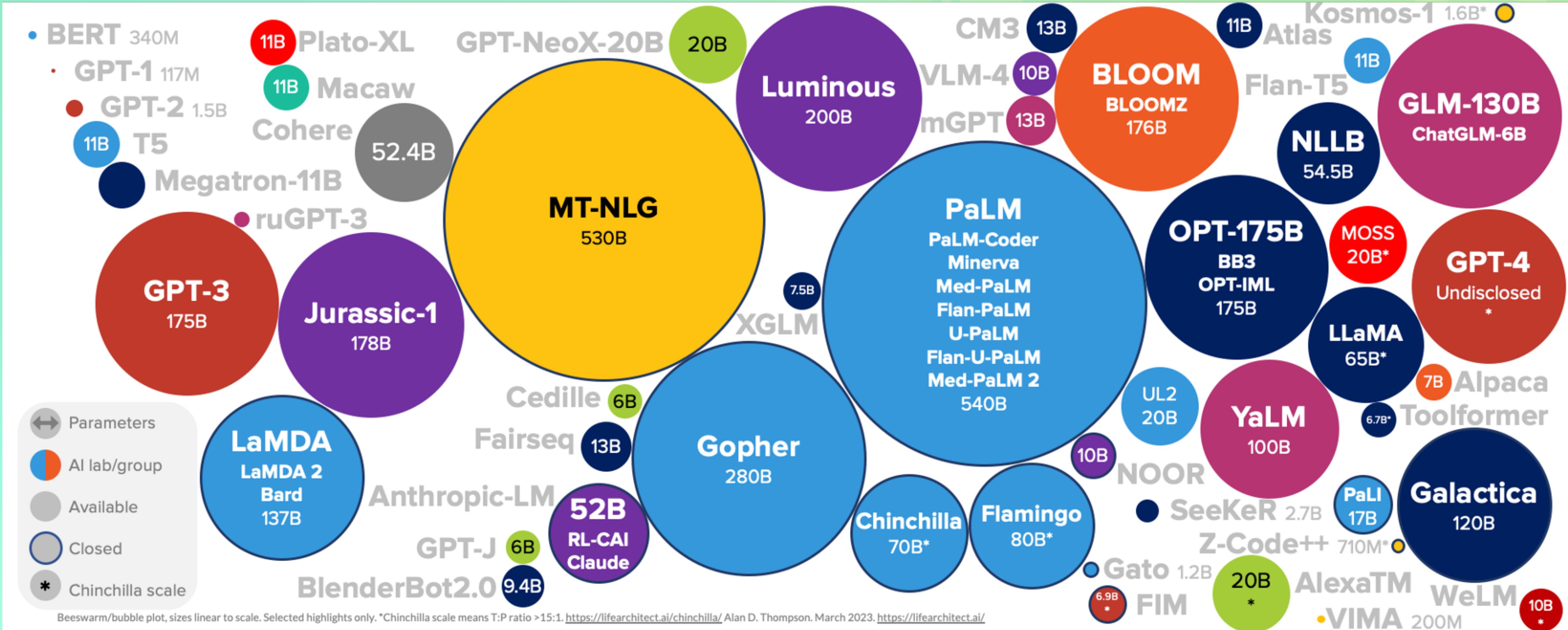
TavilySearch

Pydantic

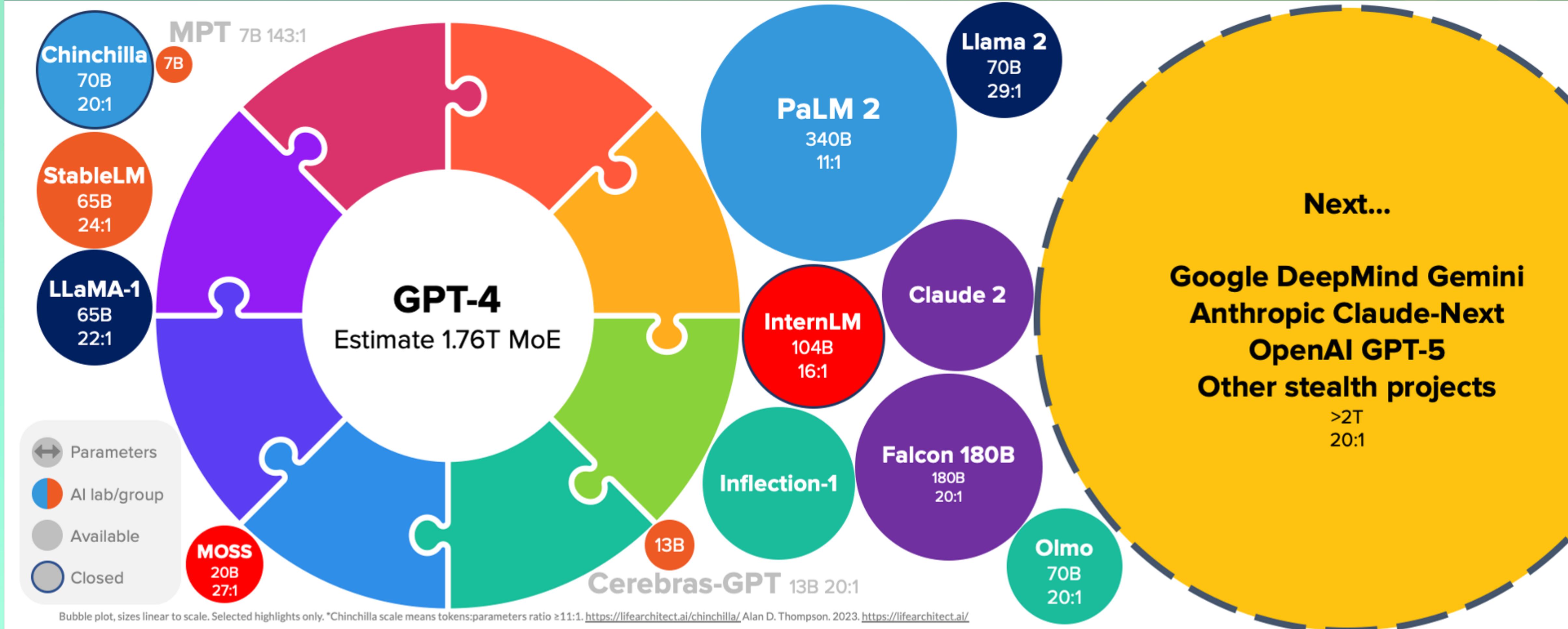
LCEL

FAISS

What and Why of LLMs

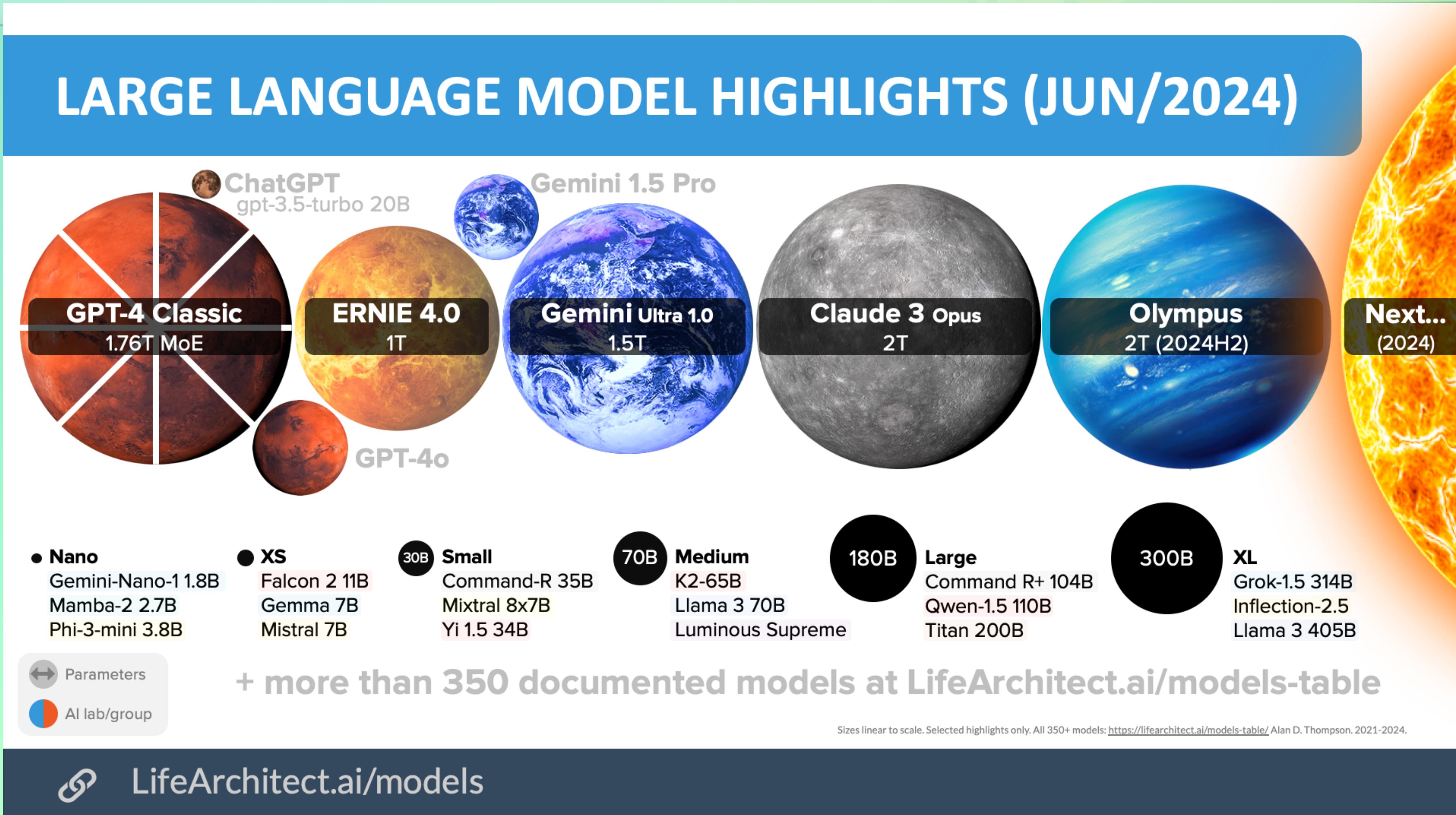


What and Why of LLMs

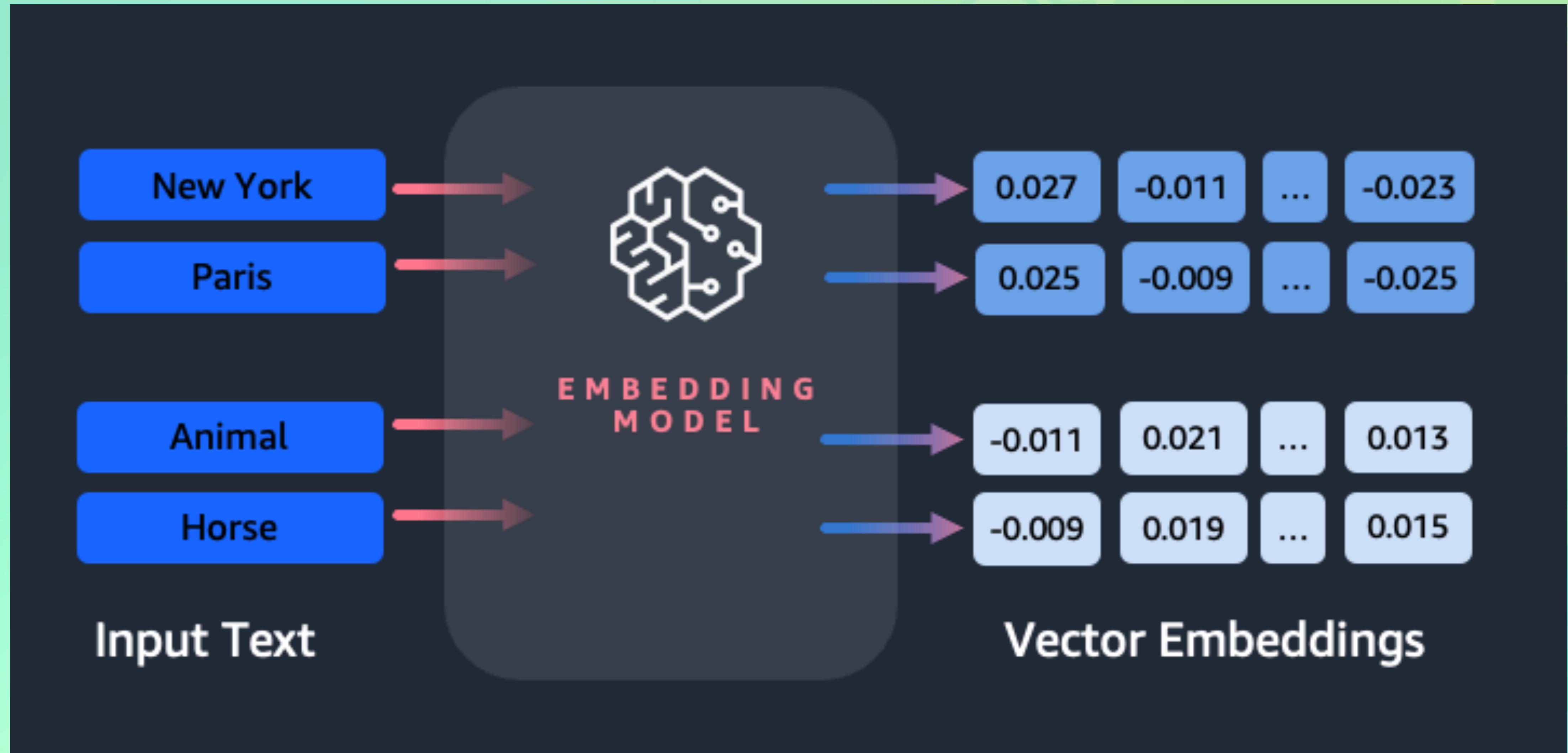


<https://lifearchitect.ai/models/>

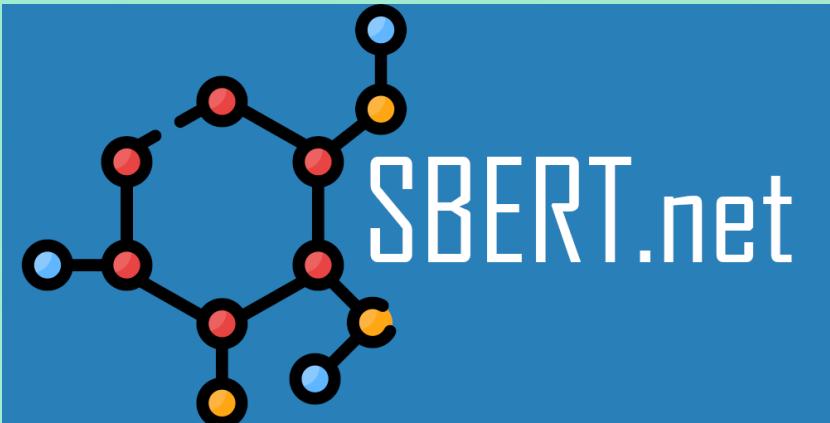
What and Why of LLMs



Embeddings



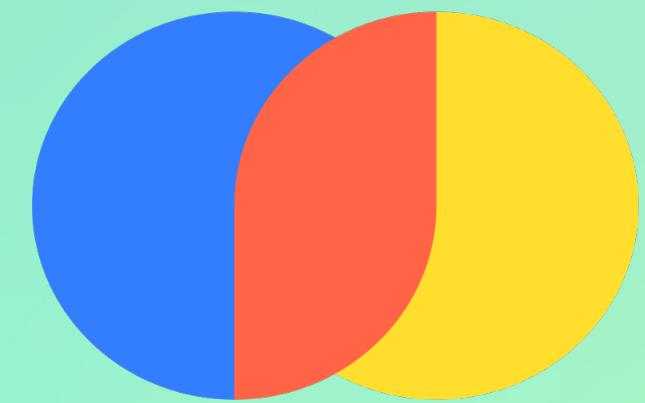
Embeddings



Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	PairClassification Average (3 datasets)	Rera Avg (4 data)
1	NV-Embed-v1	7851	29.25	4096	32768	69.32	87.35	52.8	86.91	60.5
2	voyage-large-2-instruct			1024	16000	68.28	81.49	53.35	89.24	60.0
3	Linq-Embed-Mistral	7111	26.49	4096	32768	68.17	80.2	51.42	88.35	60.2
4	SFR-Embedding-Mistral	7111	26.49	4096	32768	67.56	78.33	51.67	88.54	60.6
5	gte-Qwen1.5-7B-instruct	7099	26.45	4096	32768	67.34	79.6	55.83	87.38	60.1
6	voyage-lite-02-instruct	1220	4.54	1024	4000	67.13	79.25	52.42	86.87	58.2
7	GritLM-7B	7242	26.98	4096	32768	66.76	79.46	50.61	87.16	60.4
8	e5-mistral-7b-instruct	7111	26.49	4096	32768	66.63	78.47	50.26	88.34	60.2

<https://huggingface.co/spaces/mteb/leaderboard>

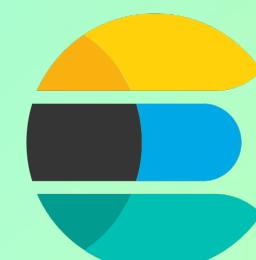
Vector Stores



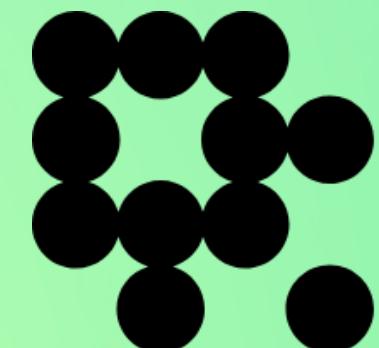
ChromaDB



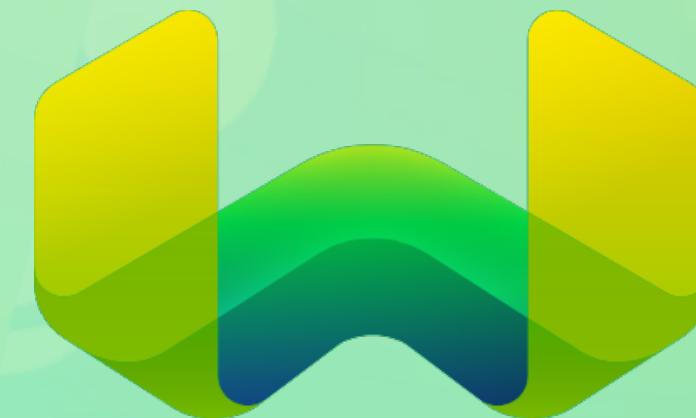
vespa



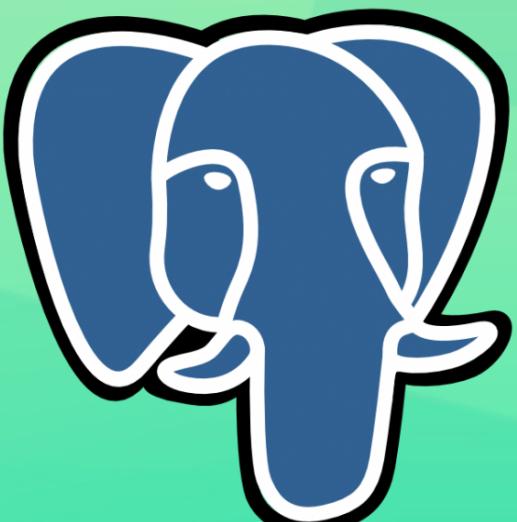
elasticsearch



LanceDB



Milvus



PostgreSQL

Weaviate

Vector Stores

Data connection

Source



Load

XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX
XXXXXXXXXX

Transform

XXXXXXXXXX
XXXXXXXXXX

XXXXXXXXXX
XXXXXXXXXX

XXXXXXXXXX
XXXXXXXXXX

Embed

0.5, 0.2...0.1, 0.9
-0.1, 0.4...1.4, 5.9

0.2, 0.7...2.1, -1.2
4.1, 3.4...-1.5, 2.5

5.5, -0.3...0.8, 2.3
2.1, 0.1...-1.7, 0.9

Store

0.5, 0.2...0.1, 0.9
;
2.1, 0.1...-1.7, 0.9

Retrieve

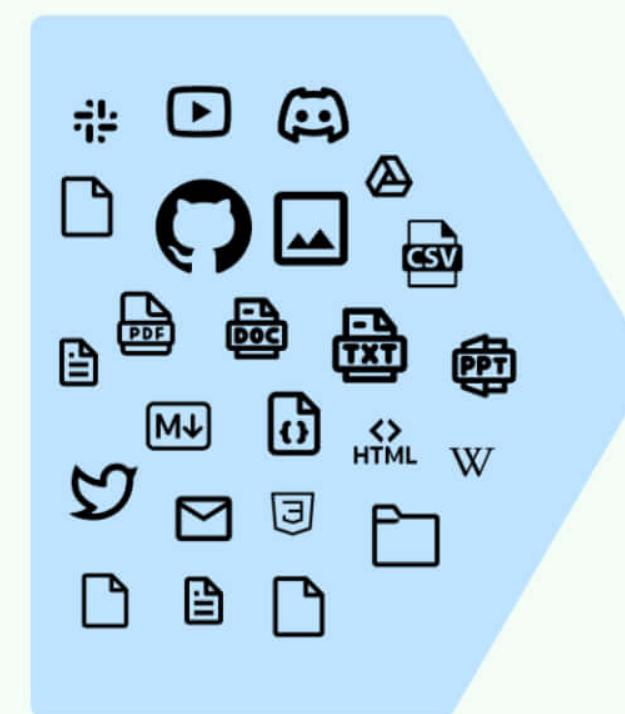
XXXXXXXXXX
XXXXXXXXXX

https://python.langchain.com/v0.1/docs/modules/data_connection/

Retrievers

Vector Stores

1. Load Source Data



Load, Transform, Embed

Vector Store

0.5, 0.2....0.1, 0.9
:
2.1, 0.1....-1.7, 0.9

Embed

5.5, -0.3...
2.1, 0.1

2. Query Vector Store

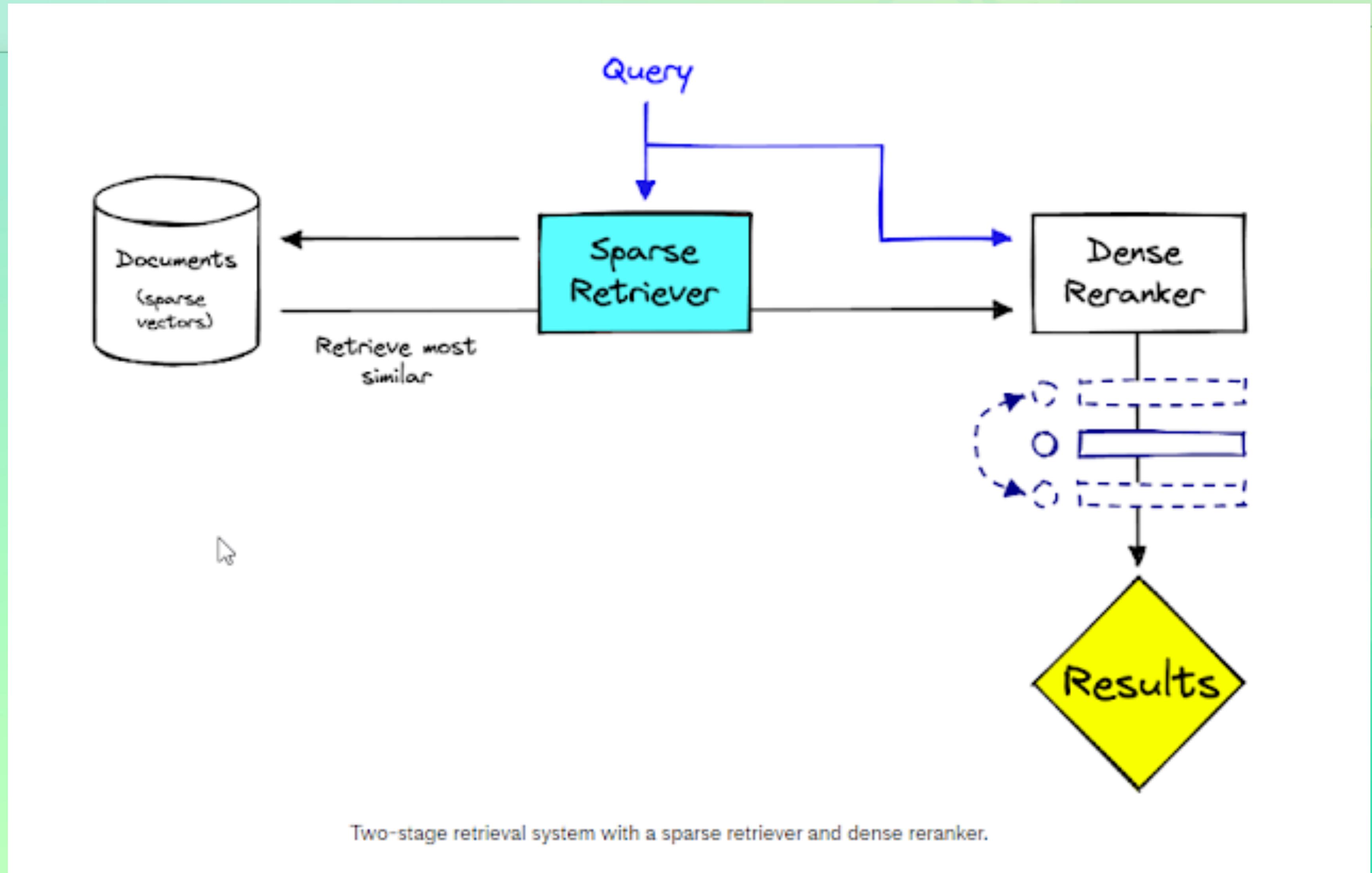
XXXXXXXXXXXXXX
XXXXXXXXXXXXXX

XXXXXXXXXXXXXX
XXXXXXXXXXXXXX

3. Retrieve 'most similar'

https://python.langchain.com/v0.1/docs/modules/data_connection/vectorstores/

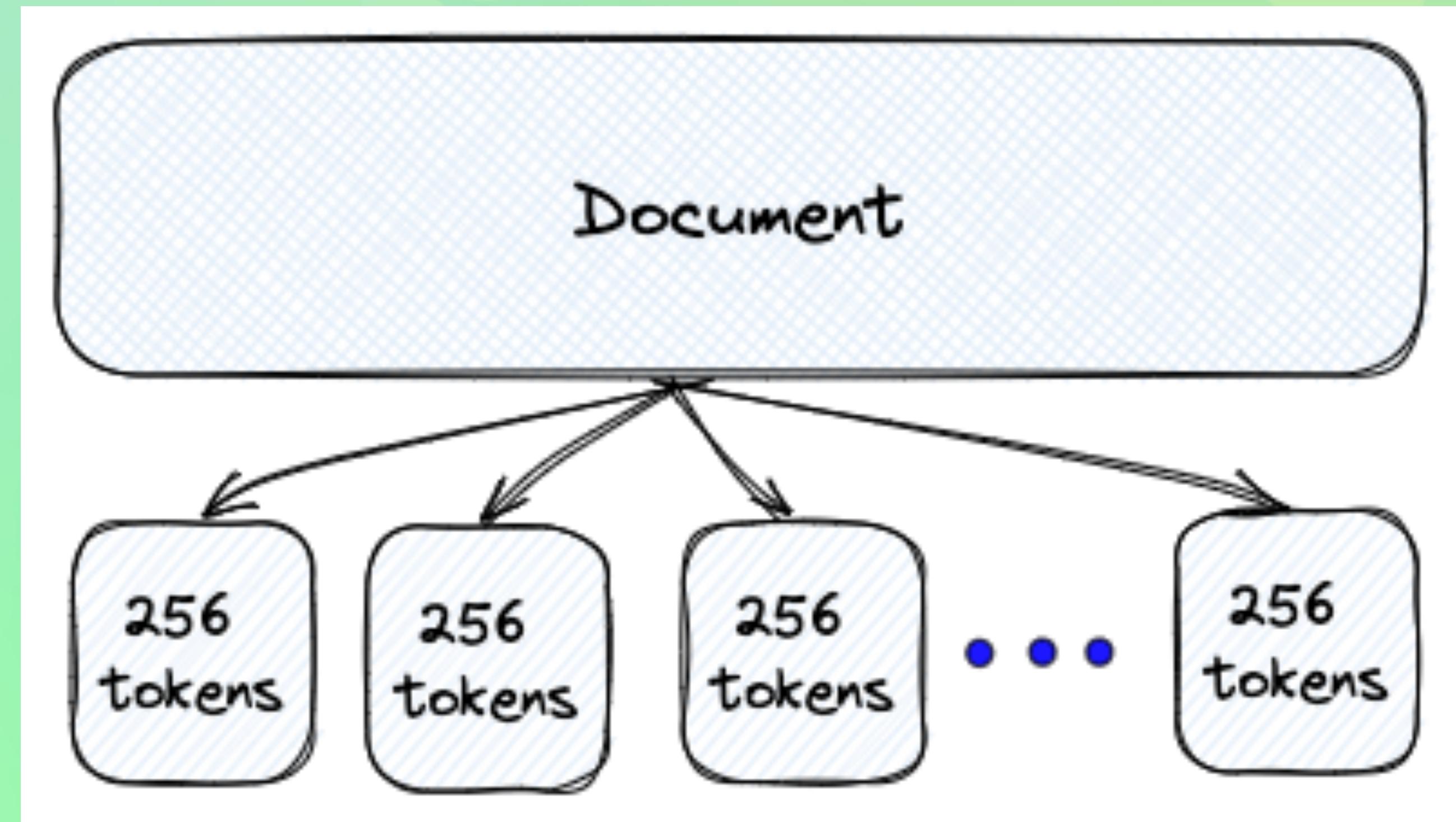
Retrievers



<https://aarontay.medium.com/boolean-vs-keyword-lexical-search-vs-semantic-keeping-things-straight-95eb503b48f5>

Chunking and Why

1. LLMs have limited input size
(Context Window)
2. We only want the Relevant
Pieces of Information



<https://app.iki.ai/content/link/1705592>

What are LLMs Good and Bad At?

LLMs are Good at

Text Generation

Reasoning

Summarization

Paraphrasing

Data Extraction from Context

LLMs are Bad at

Knowing about New Data

Knowing about Our Data

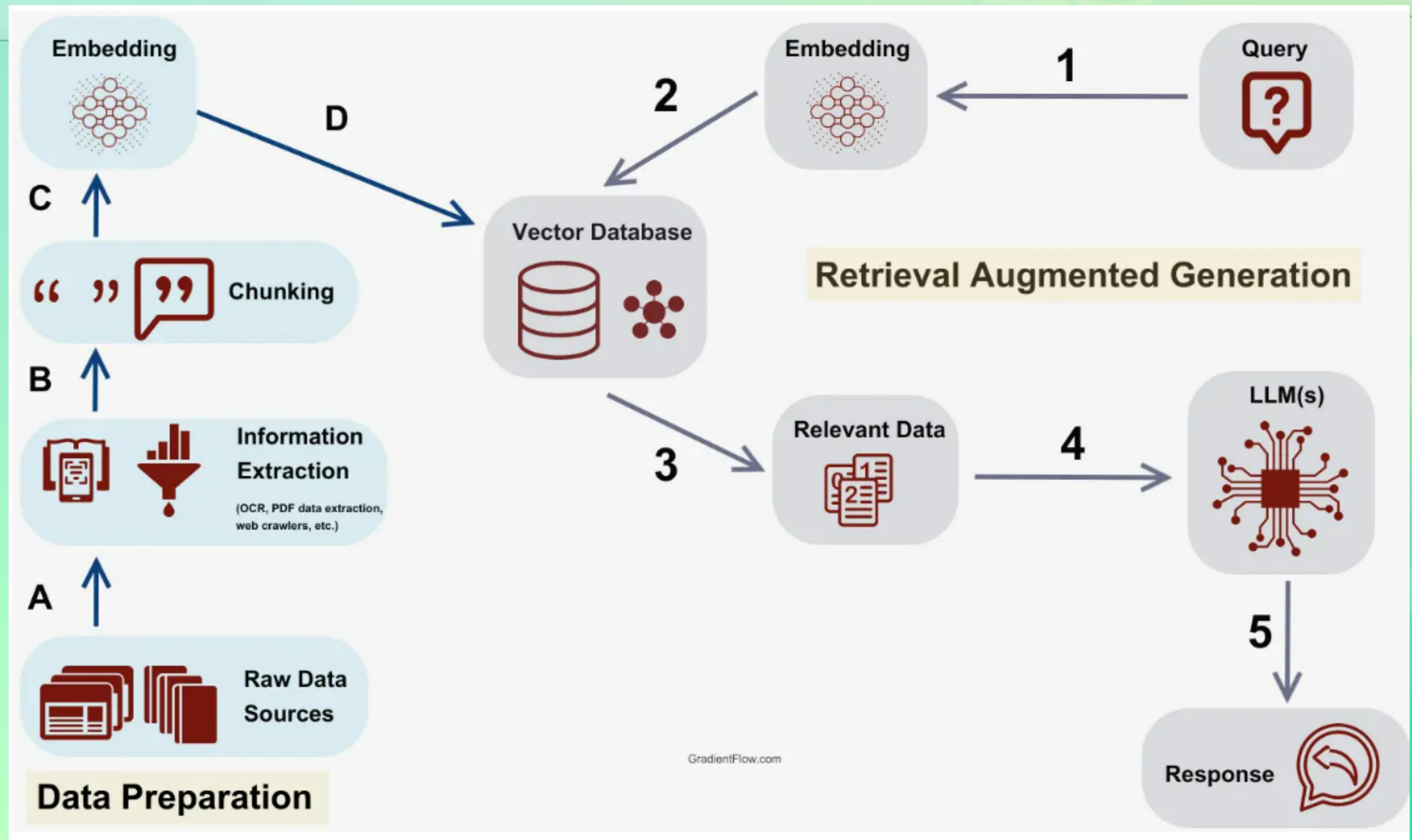
Hallucination

Remembering Long-tail Information

What are LLMs Good and Bad At?

1. Our Chatbot must have the latest information
2. Our Chatbot must know about our Business Specific Data
3. Our Chatbot should be trustworthy

What is RAG?



<https://www.linkedin.com/pulse/what-rag-why-should-you-care-aymen-noor-eyymf/>

What is LangChain?



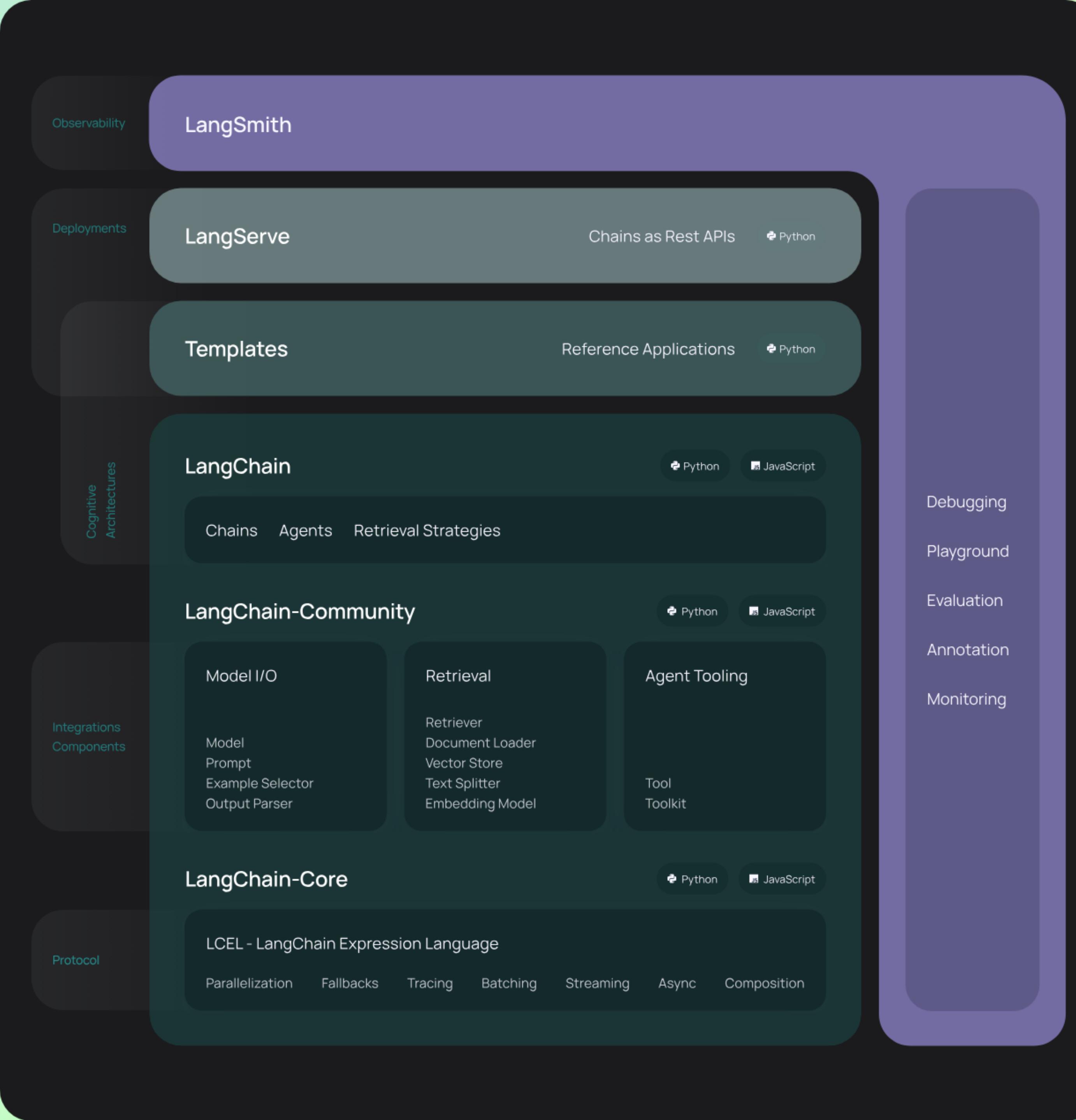
<https://pub.dev/documentation/langchain/latest/>

What is LangChain?

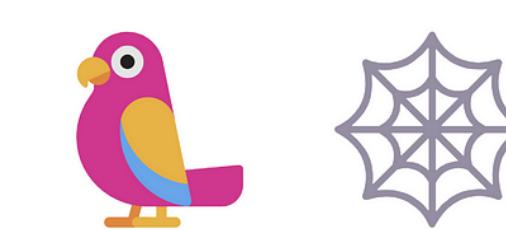


LangChain

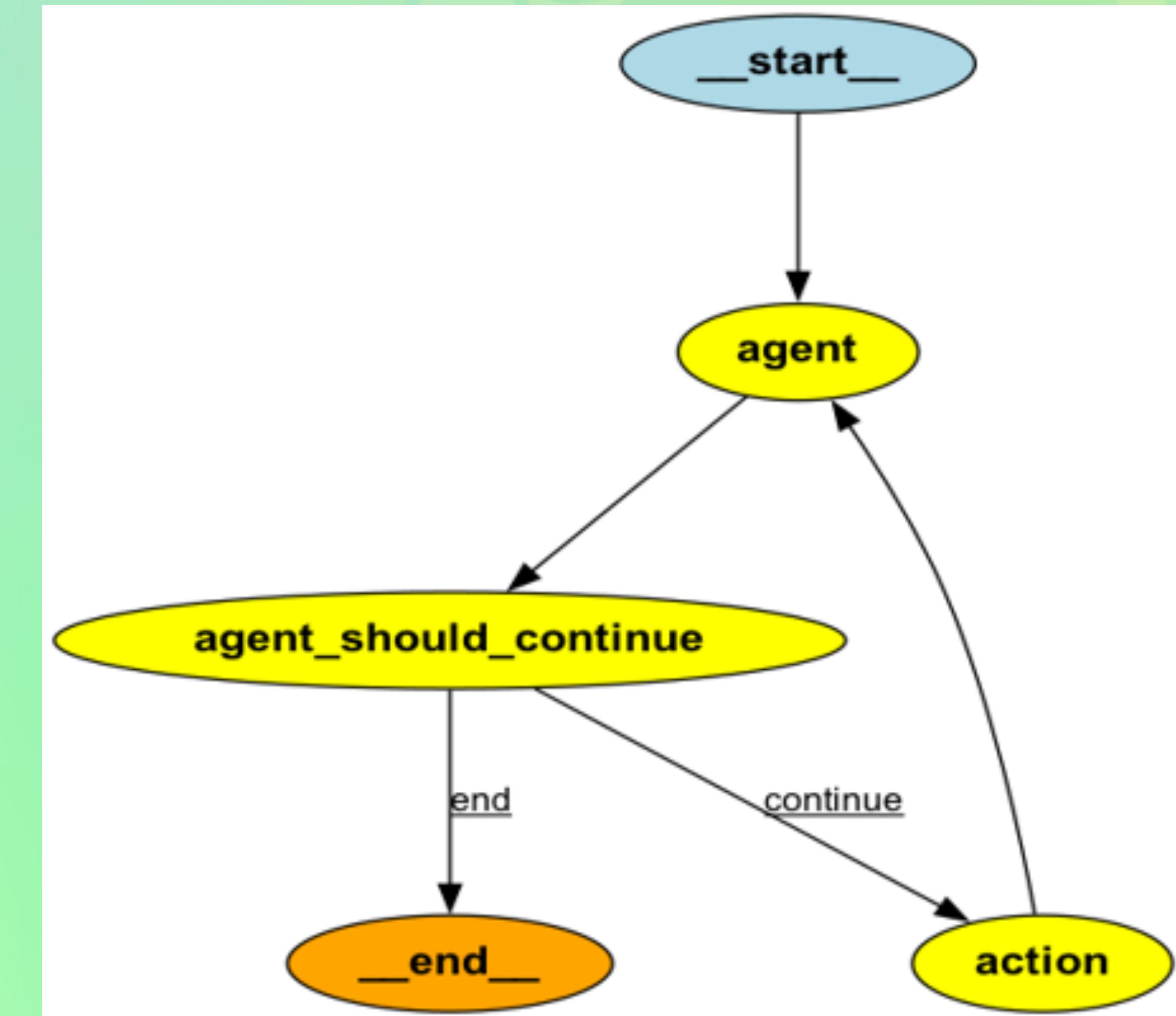
<https://python.langchain.com/v0.2/docs/introduction/>



What is LangGraph?



LangGraph



<https://blog.gopenai.com/conditional-edge-and-cycle-in-langgraph-explained-da4a112bf1ea>

What is LangGraph?

```
class AgentState(TypedDict):
    """The dictionary keeps track of the data required by the various nodes in the graph"""

    query: str
    chat_history: list[BaseMessage]
    generation: str
    documents: list[Document]
```

