

Scan me

Developing GenAI Applications with LangChain

chukmunlee@gmail.com



Updated the workshop template

- Either re clone
- Copy the files to your existing repository
 - `2_vector_store/2_llm_with_retrievers.py`
 - `ebook/*.pdf`

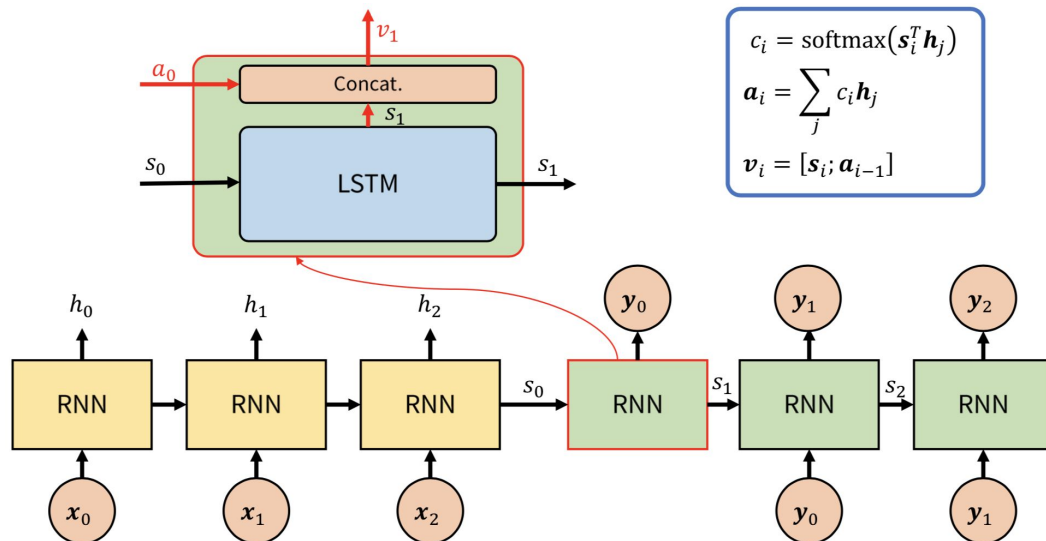
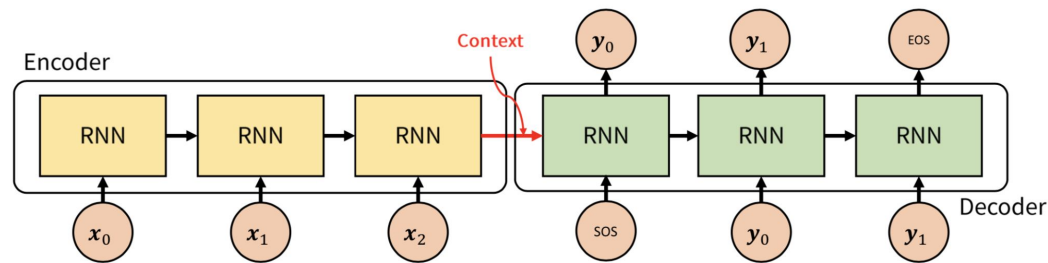
What are LLMs?

Model that is trained on text

- Eg books, code, chat conversations
- Learn statistical relationships between words and phrases
- Applications include code generation, text summary, translation, etc

Key to LLMs is transformers with attention

- Attention mechanism allows the model to focus on key and relevant information



What is LangChain?

Framework for building LLM applications

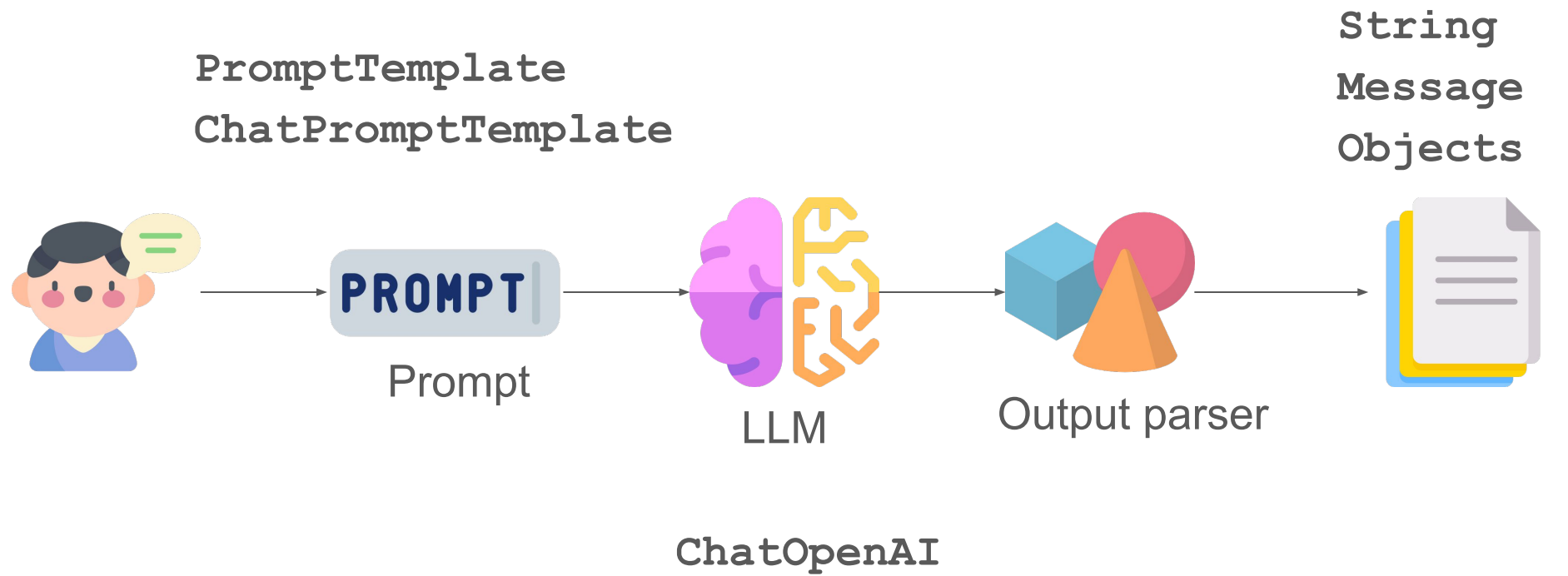
- Applications can be written in Python, JavaScript, Java, C# (not official)

Single API for different LLMs - GPT, Anthropic, Ollama, Vertex, etc

Alternative to LangChain

- LlamaIndex
- Flowise

Basic LLM Application Structure

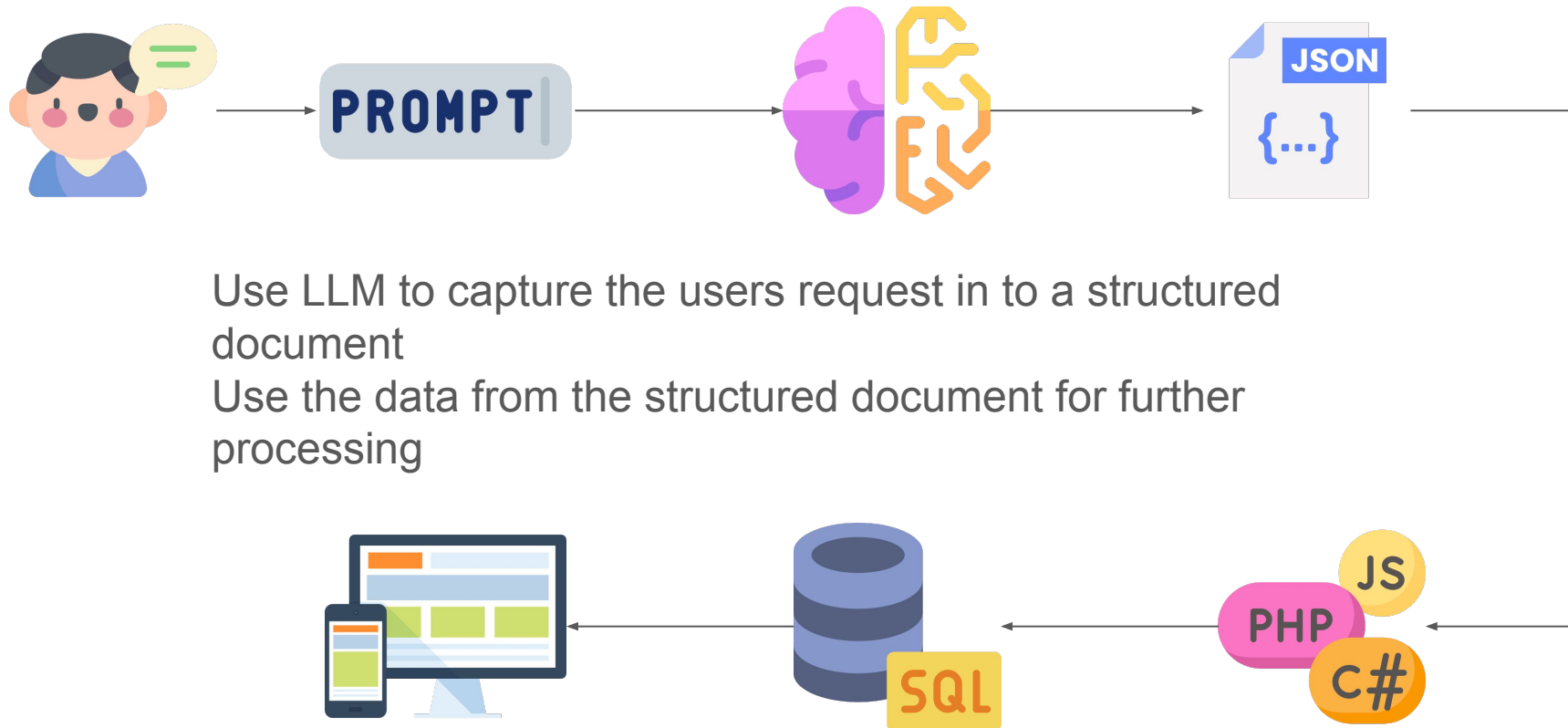




Prompts and LLM

1. Your first prompt and LLM
 - `0_basic_prompt.py`
2. Cost of answering a question
3. Conditioning and alignment
 - `1_chat_messages.py`
4. Structuring outputs (illustration next slide)
 - `2_generate_questions_and_answers.py`

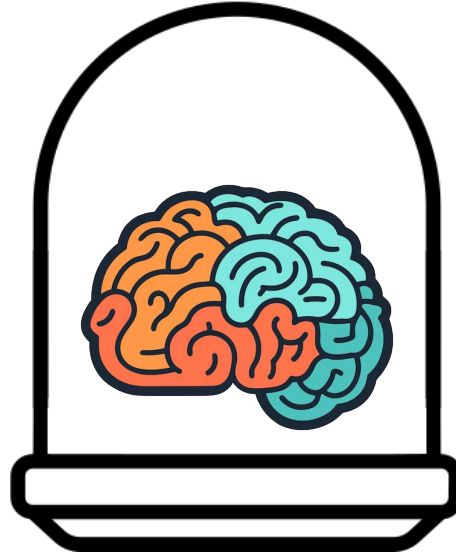
Structure Output



Language Models are Self Contained

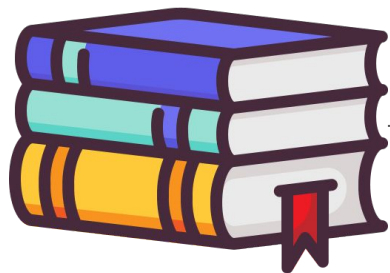


I don't know



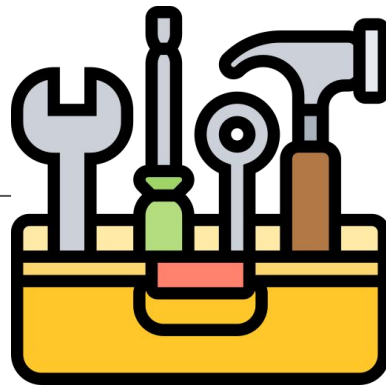
Make up an answer

LLM with External Knowledge



Documents

Lookup proprietary information



Tools

Interact with the real-world to get information

Embeddings

RGBA embedding describe the intensity of each of the component



[199, 21, 133]



[55, 192, 203]



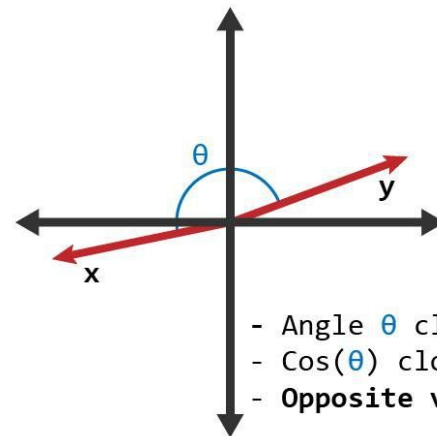
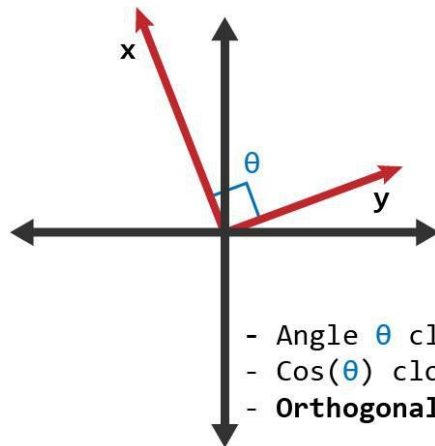
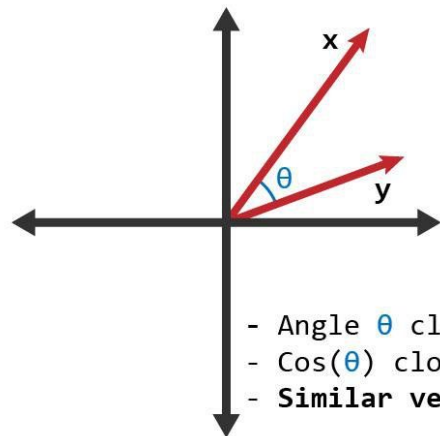
[0, 255, 0]



[199, 21, 133, 20]

More dimension, more information

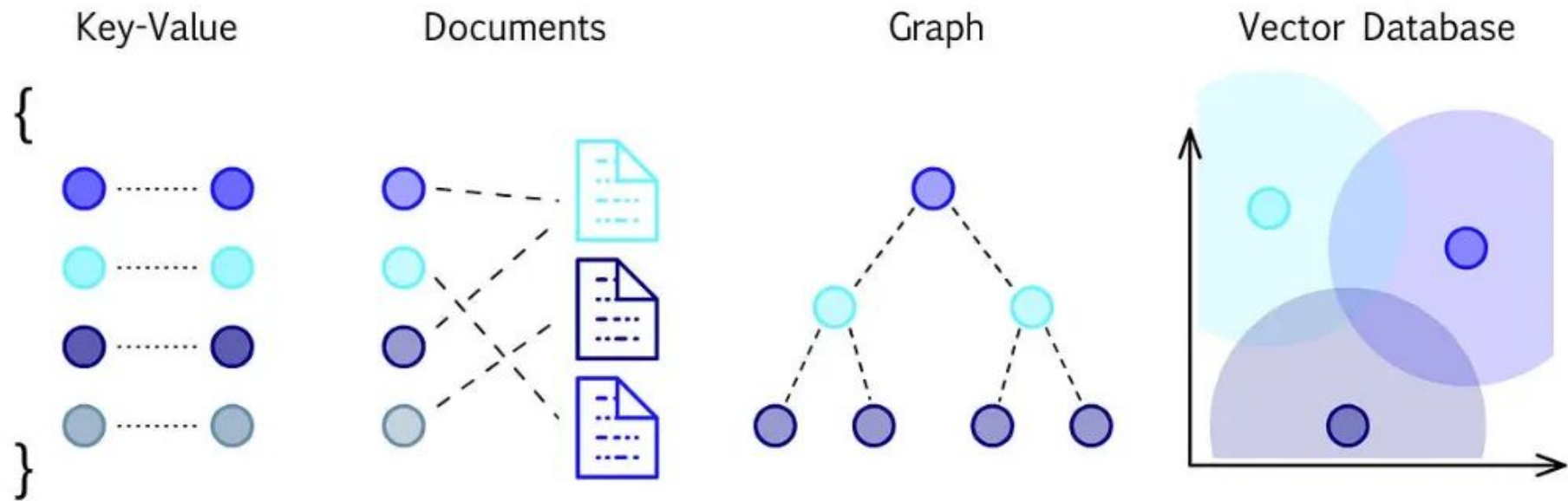
Cosine Similarity






$$\text{cosine}(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



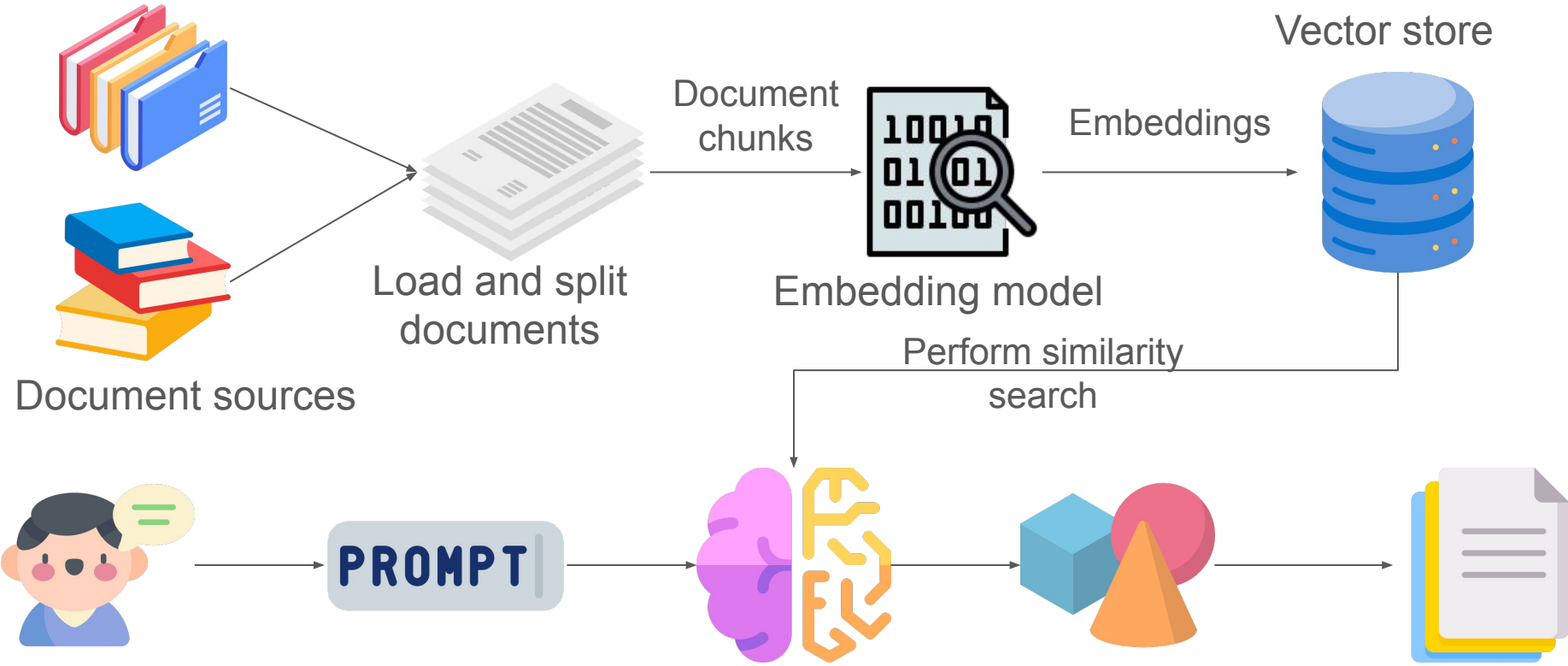
Database Structure and Use Cases



Some Vector Databases

	Dedicated	Supports
Memory		
Disk		

Querying from Vector Store

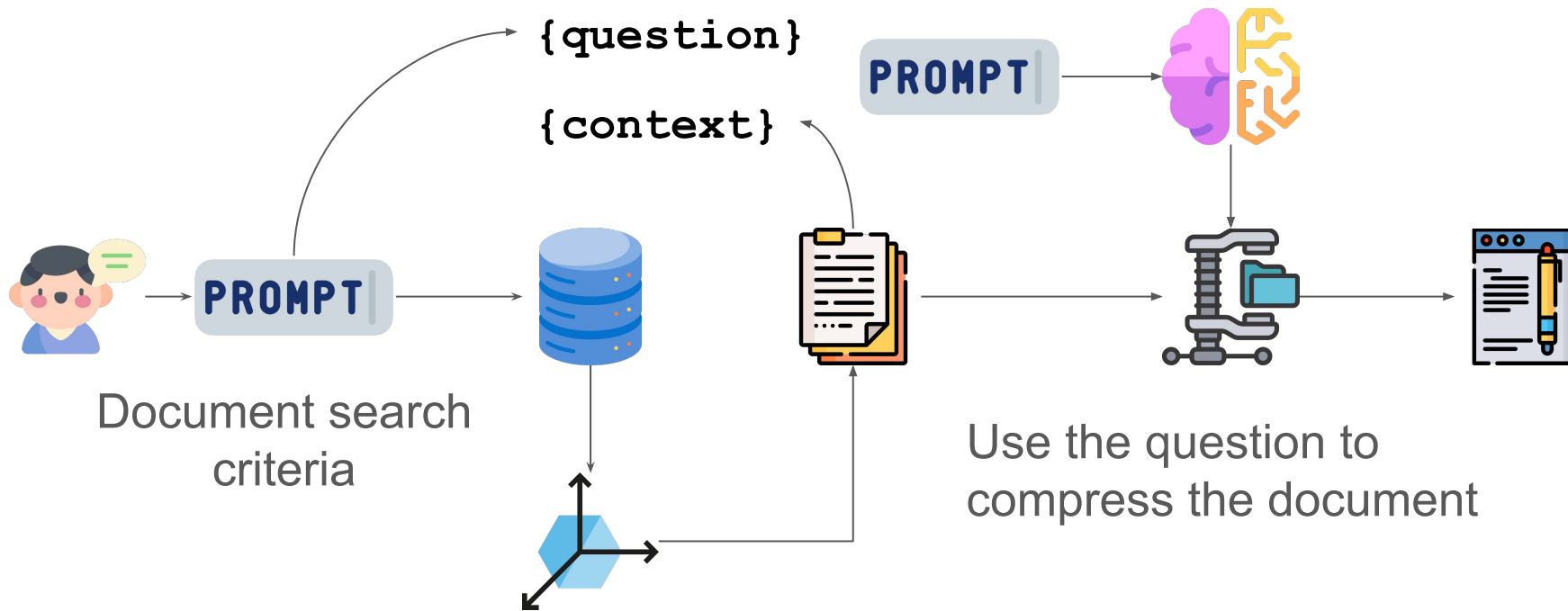




Vector Store

1. Loading and querying documents
 - `0_load_and_query.py`
2. Contextual compression (illustration next slide)
 - `1_retriever_with_compressor.py`
 - `load_vector_db.py`
3. Retrieving from vector store (illustration next next slide)
 - `2_llm_with_retrievers.py`

Contextual Compression

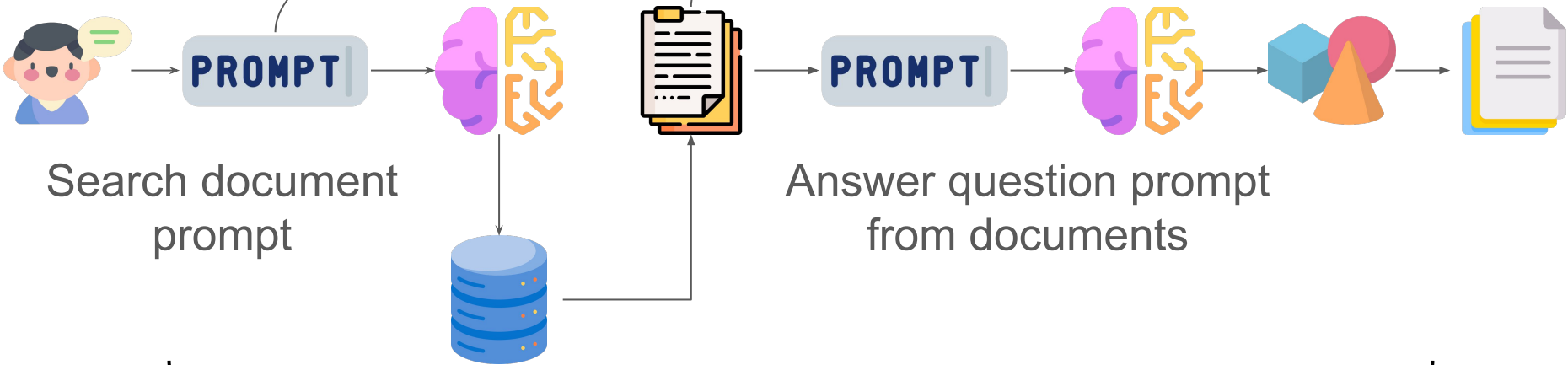


Augmenting LLM with Vector Store

Retrieval-Augmented Generation → {question}

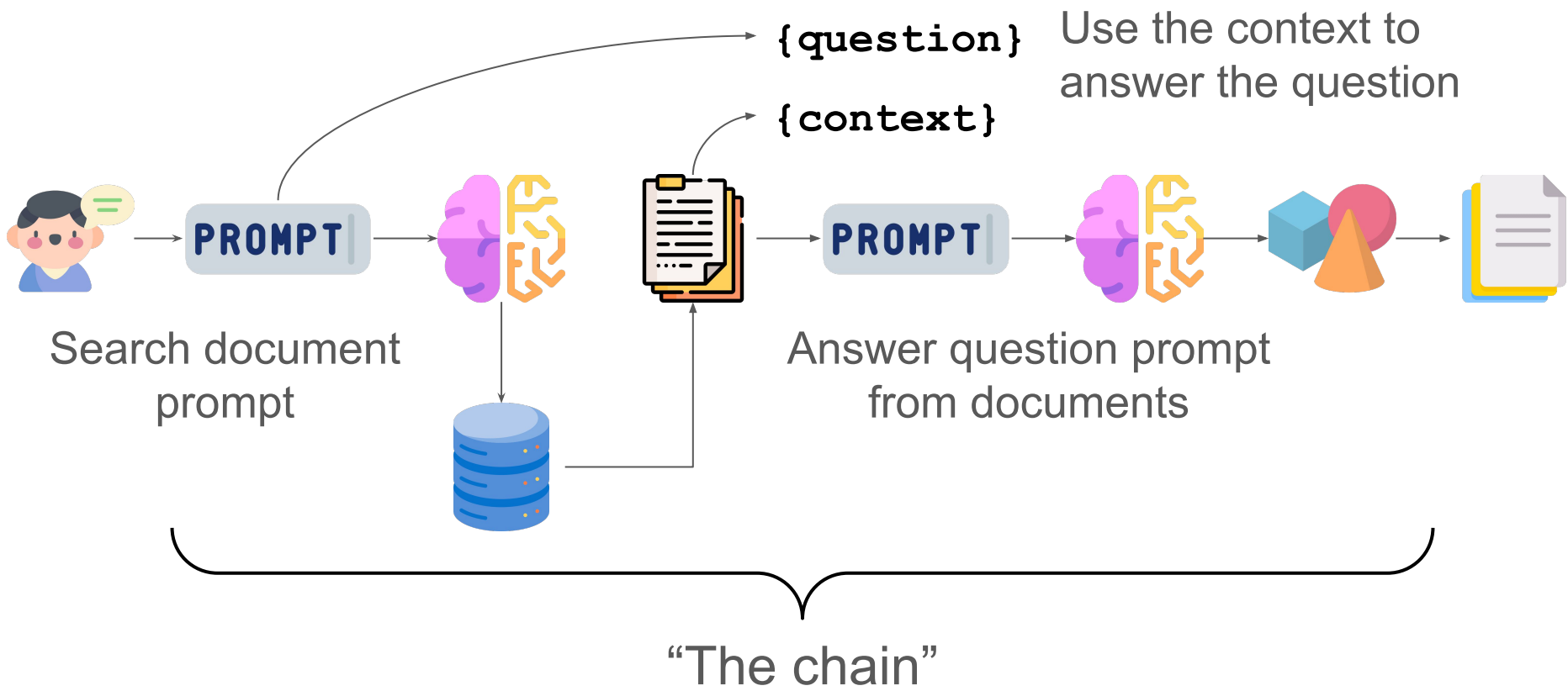
Use the context to answer the question

{context}



“The chain”

Augmenting LLM with Vector Store



Agents

Agents determines how best to answer a question with a given set of tools

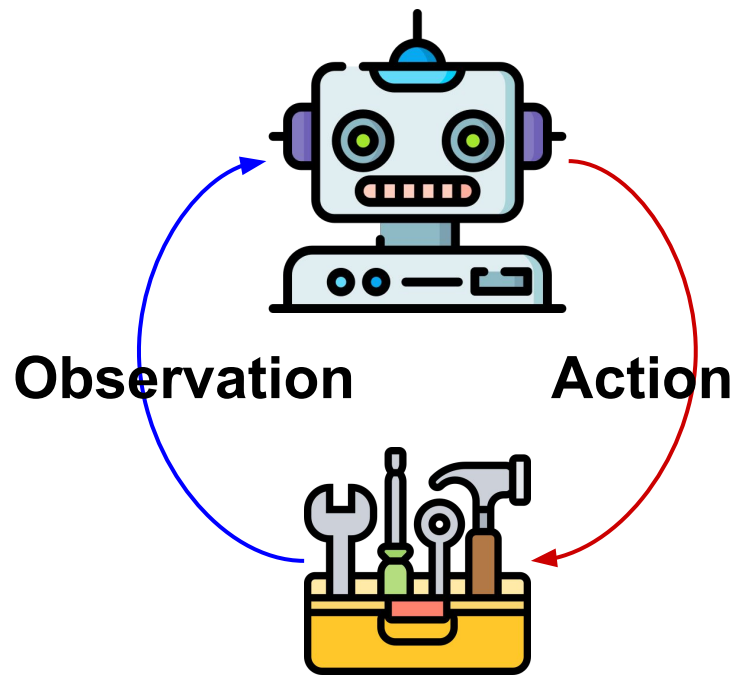
- Tools include build-in eg wikipedia, duck duck go search
- Custom tools

May iterate to refine the answer

-

Memory allow agents to recall previous interactions

- Infer from past interactions eg 'the problem'



Example of Agent's Prompt

<https://smith.langchain.com/hub/hwchase17/react-chat>

Key variables

- `{tool_names}`, `{tools}` - tool names and custom tools
- `{chat_history}` - past conversation from memory
- `{agent_scratchpad}` - used by the agent
- `{input}` - question the agent is answering



Agents

1. Agents and tools
 - `0_tools.py`
2. Conversation agents
 - `1_conversation_agent.py`

Capstone Example

More advance example in 4_capstone

- Retrievers as tools
- Preloading the conversation history

Workshop solution

<https://github.com/chukmunlee/jul13-2024-geekshacking-genai-workshop.git>

Thank you
for
attending!