

LangSmith

How to Improve LLM Applications in Production



Presented by

*Greg Loughnane, Founder & CEO
Chris Alexiuk, Co-Founder & CTO*



ALIGNING OUR AIM

lide



BY THE END OF TODAY...

- Baseline production RAG apps w/
evaluation & E2E testing
- Improve RAG re; eval metrics and
feedback/annotation on examples

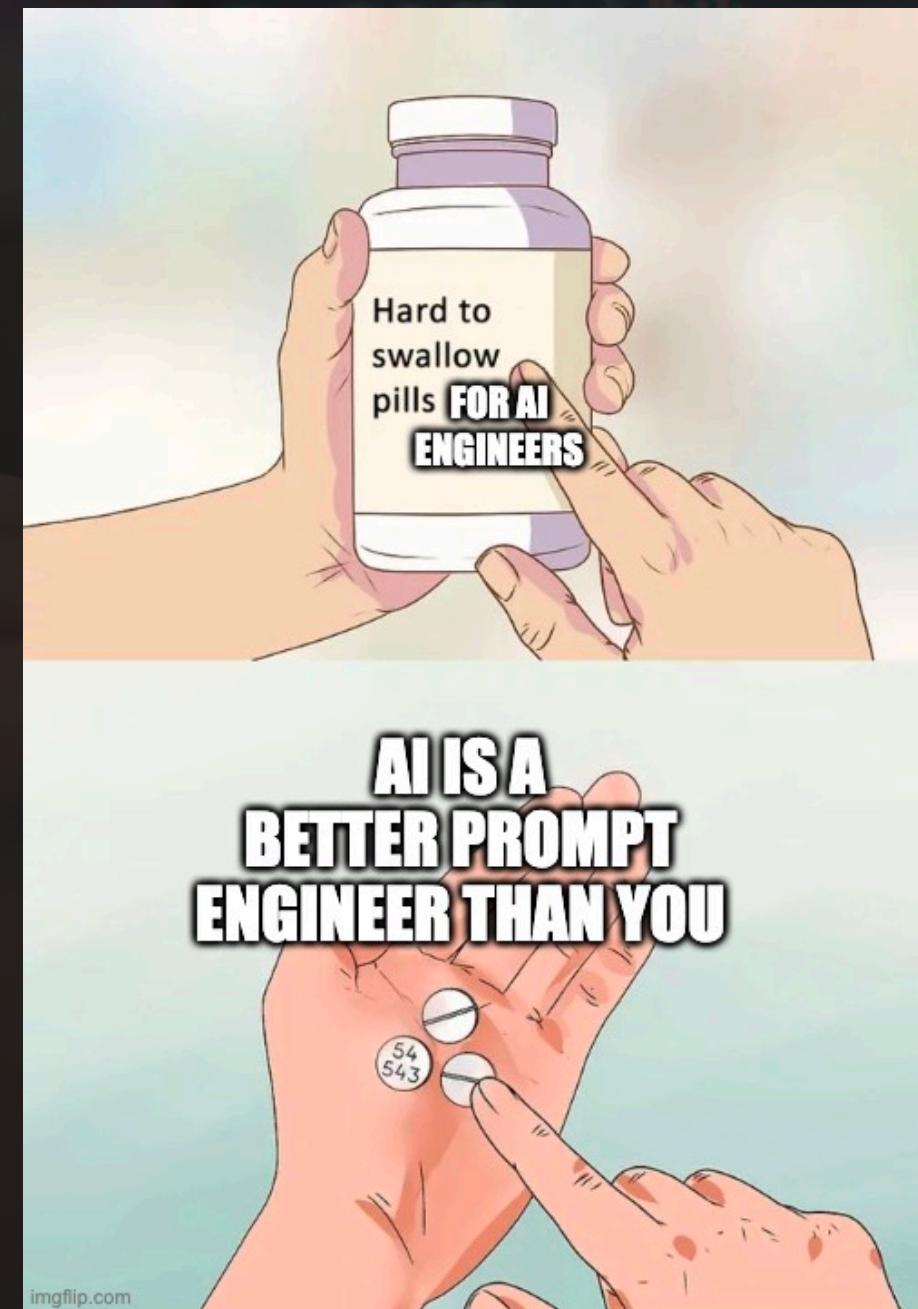
OVERVIEW

- LangChain RAG
- Monitoring/Evaluation
- Baseline RAG
- Annotation/feedback
- Advanced retrieval
- Improved RAG
- Conclusions, QA

Data Science
AI Engineering



AI ENGINEERING



AI IS A
BETTER PROMPT
ENGINEER THAN YOU



“The real power comes when you can **combine** them with other sources of computation or knowledge.”

~ Harrison Chase
Creator of LangChain

IN COLLABORATION WITH
LangChain

Functions, Tools and Agents with LangChain

Learn and apply the new capabilities of LLMs as a developer tool.

- Learn about the most recent advancements in LLM APIs.
- Use LangChain Expression Language (LCEL), a new syntax to compose and customize chains and agents faster.
- Apply these new capabilities by building up a conversational agent.

Intermediate

Harrison Chase

Prerequisite recommendation: Basic Python and familiarity with writing prompts for Large Language Models.

[Enroll For Free](#)

[Learn more](#)

IN COLLABORATION WITH
LangChain

LangChain for LLM Application Development

The framework to take LLMs out of the box. Learn to use LangChain to call LLMs into new environments, and use memories, chains, and agents to take on new and complex tasks.

- Learn LangChain directly from the creator of the framework, Harrison Chase
- Apply LLMs to proprietary data to build personal assistants and specialized chatbots
- Use agents, chained calls, and memories to expand your use of LLMs

Beginner

Harrison Chase, Andrew Ng

Prerequisite recommendation: Basic Python

[Enroll For Free](#)

[Learn more](#)

IN COLLABORATION WITH
LangChain

LangChain: Chat with Your Data

Create a chatbot to interface with your private data and documents using LangChain.

- Learn from LangChain creator, Harrison Chase
- Utilize 80+ loaders for diverse data sources in LangChain
- Create a chatbot to interact with your own documents and data

Beginner

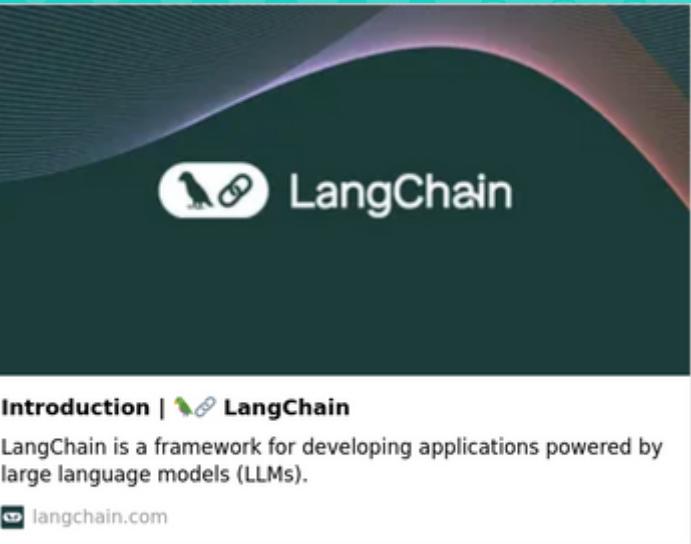
Harrison Chase

Prerequisite recommendation: Basic Python

[Enroll For Free](#)

[Learn more](#)

:WHAT IS IT?



LangChain

LangChain is a framework for developing applications powered by language models. We believe that the most powerful and differentiated applications will **not only call out to a language model via an api, but will also:**

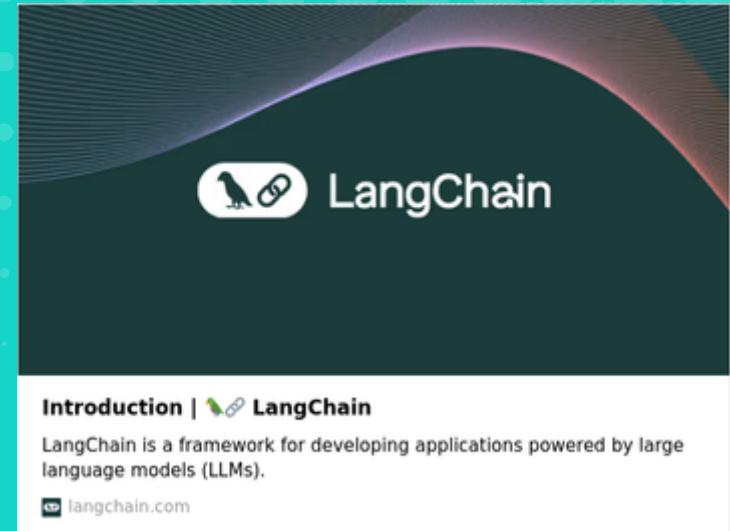
1. Be **data-aware**: connect a language model to **other sources of data**
2. Be **agentic**: Allow a language model to interact with its environment

As such, the LangChain framework is designed with the objective in mind to enable those types of applications.

:WHAT IS IT?

LangChain is a framework for developing applications powered by language models. It enables applications that:

- **Are context-aware:** connect a language model to sources of context (prompt instructions, few shot examples, content to ground its response in, etc.)
- **Reason:** rely on a language model to reason (about how to answer based on provided context, what actions to take, etc.)



OBSERVABILITY



LangSmith

DEPLOYMENT



LangServe

Chains as Rest APIs

APPLICATION



LangChain

Python

JavaScript

INTEGRATIONS
COMPONENTS

Models I/O

Model
Prompt
Example Selector
Output Parser

Retrieval

Retriever
Document Loader
Vector Store
Text Splitter
Embedding Model

Agent Tooling

Tools
Toolkits

PROTOCOL

LCEL - LangChain Expression Language

Parallelization · Fallbacks · Tracing · Batching · Streaming · Async · Composition

Monitoring

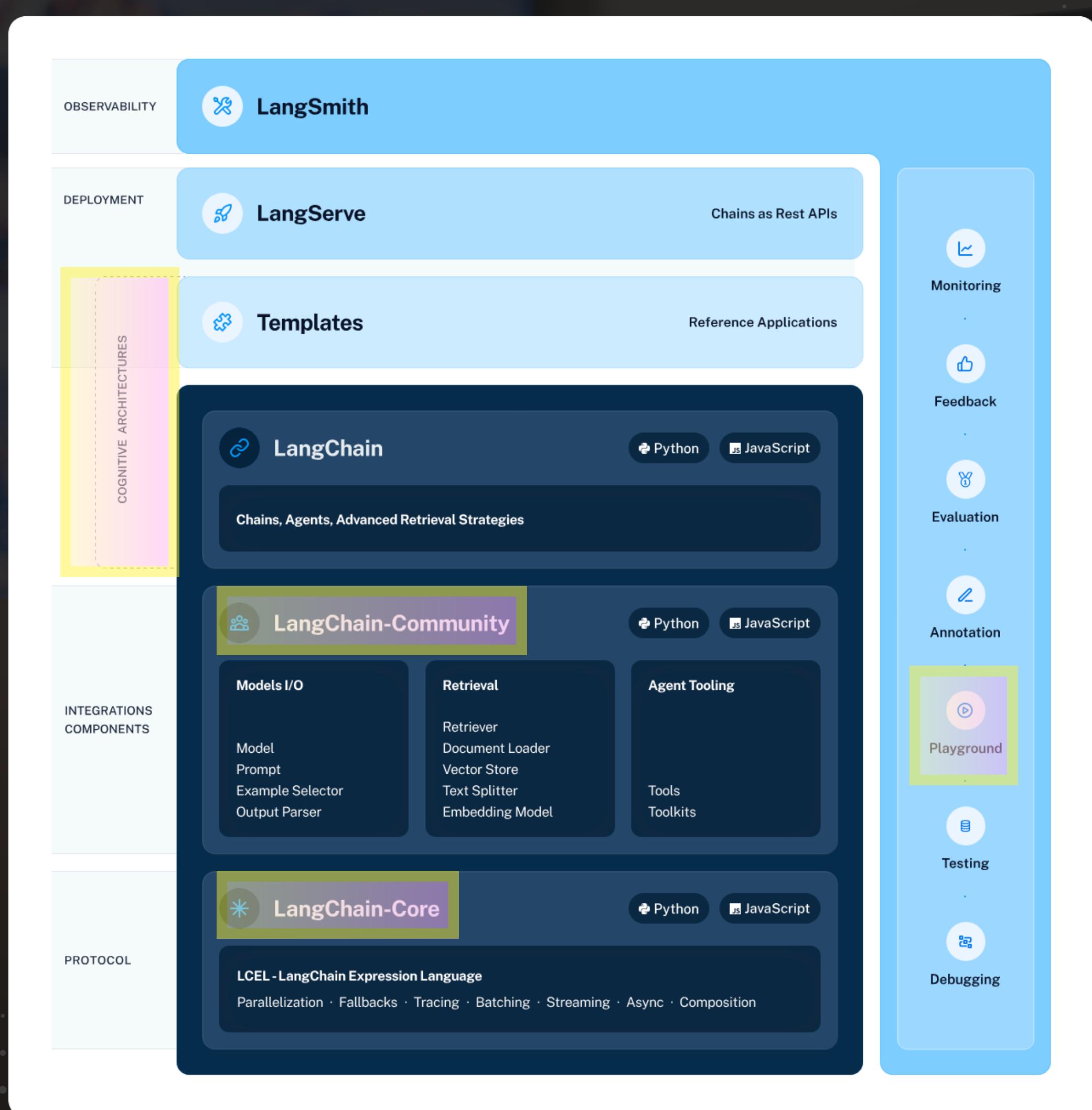
Evaluation

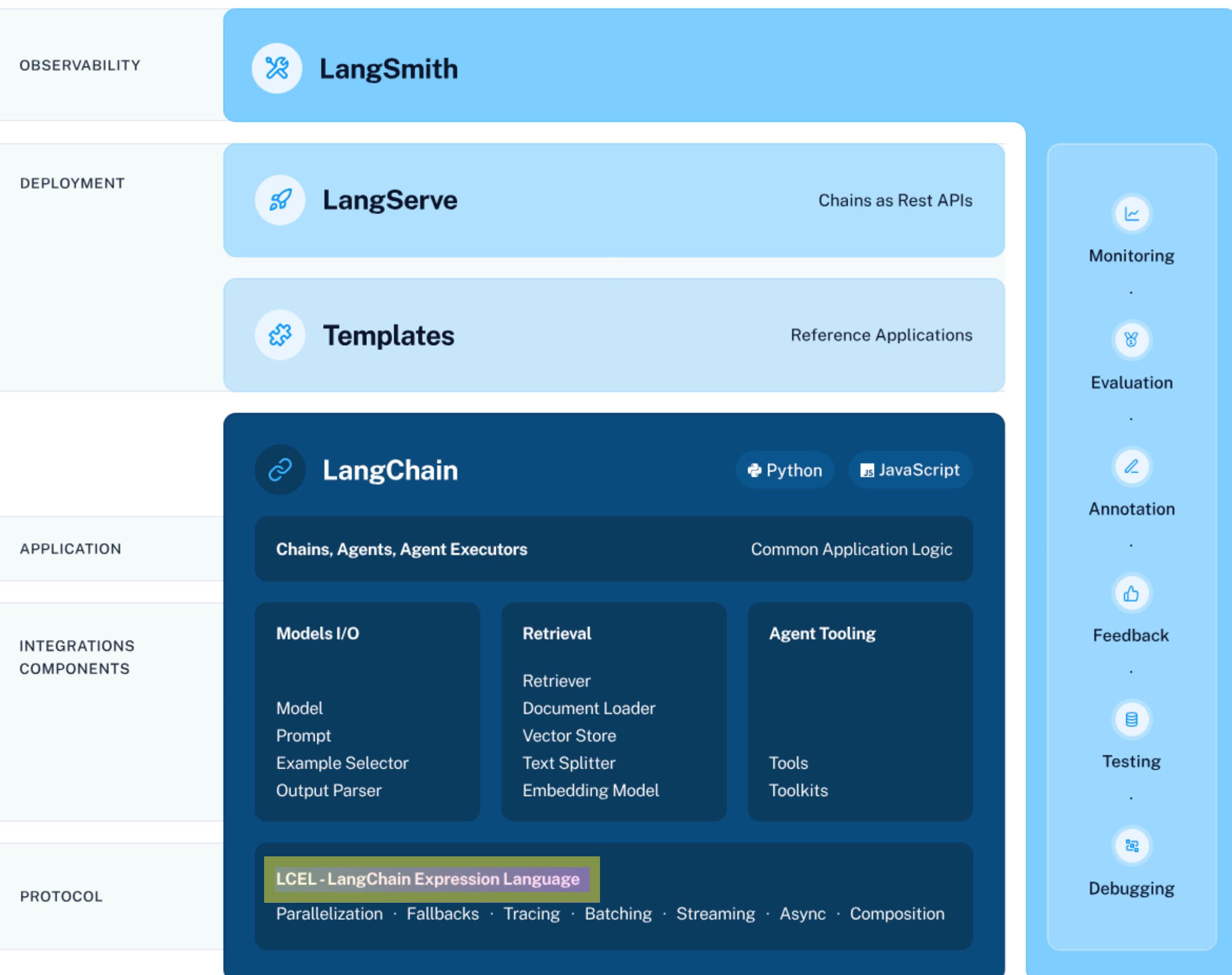
Annotation

Feedback

Testing

Debugging







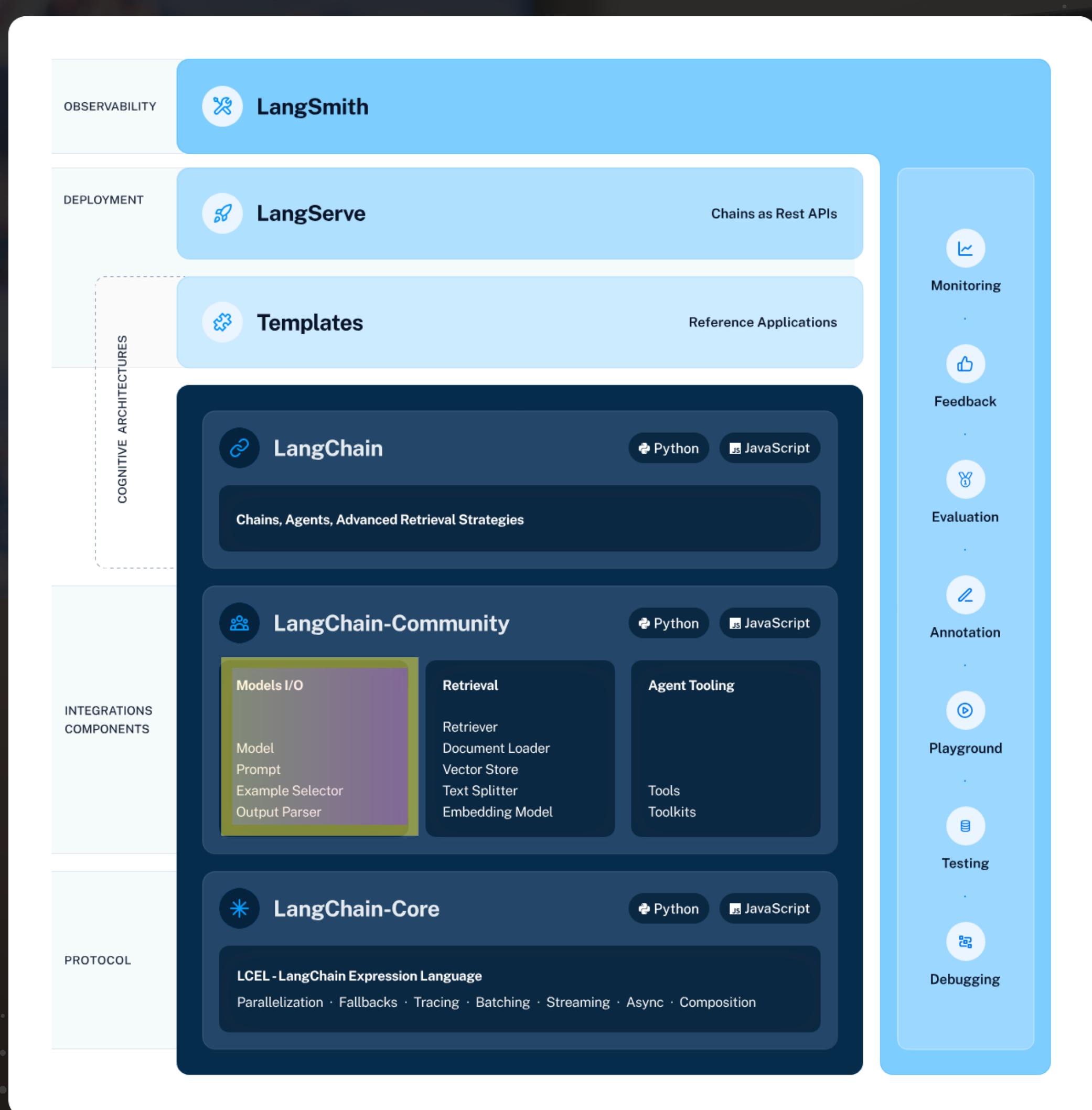
LangChain Expression Language

- Declarative way to **easily compose chains** with elegance
- Handles all the finicky bits
 - Streaming
 - Async
 - Parallel execution
 - Retries & fallbacks
 - Access intermediate results
 - I/O schemas
 - **LangSmith tracing**
 - LangServe deployment

```
llm_chain = LLMChain(prompt=prompt, llm=llm)
```

```
llm_chain = prompt | llm
```





: MODELS

- “**Chat**-Style”
 - *Fine-tuned for **chat***
 - *Often **instruction**-tuned*



System Message

User Message

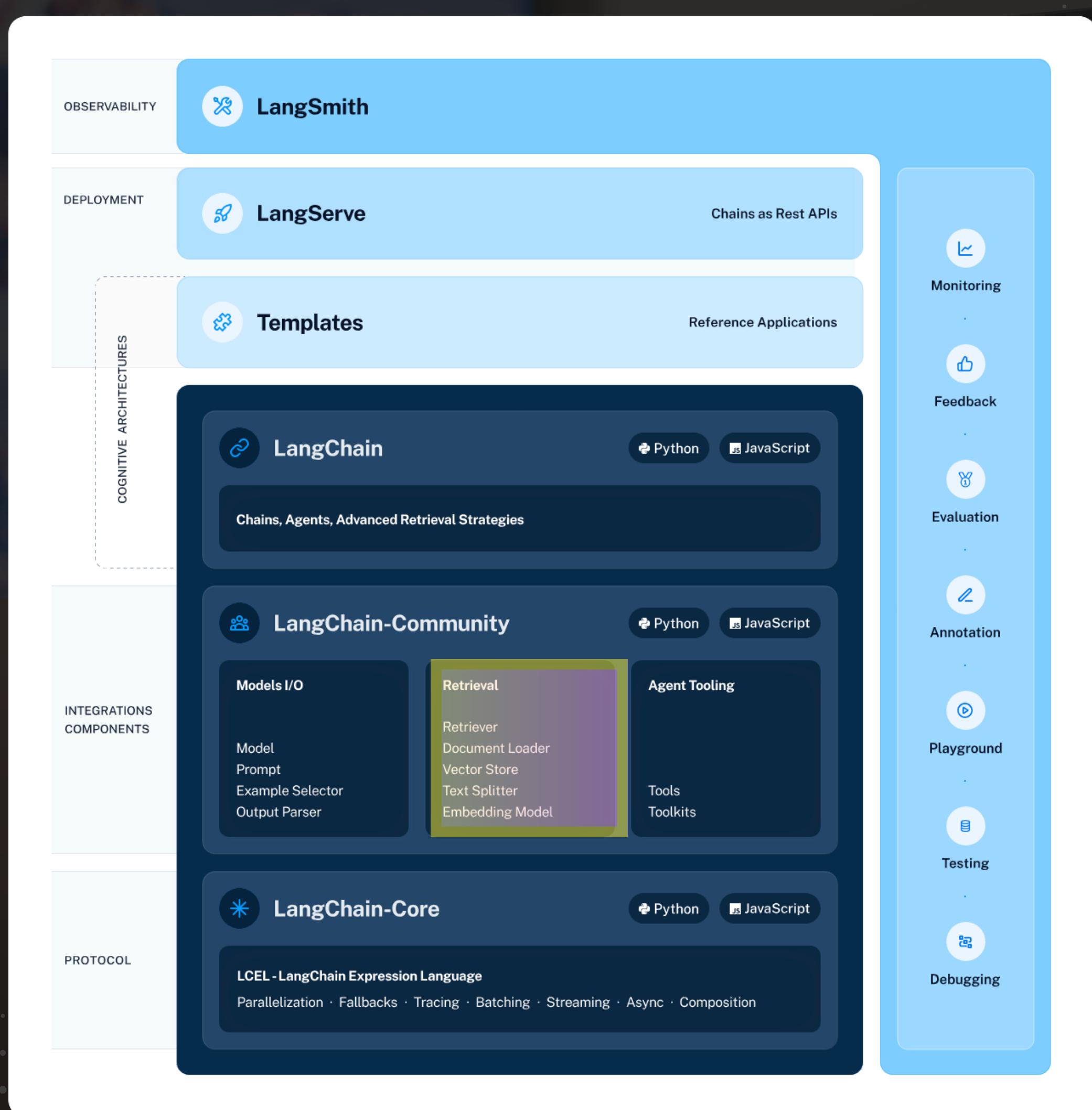
Assistant Message

System Message

Human Message

AI Message





• RETRIEVAL AUGMENTED GENERATION

Retrieval

- **Finds references** in your documents

Augmented

- **Adds references** to your prompts

Generation

- **Improves answers** to questions!

Sprinkling references into my essay that I didn't even read

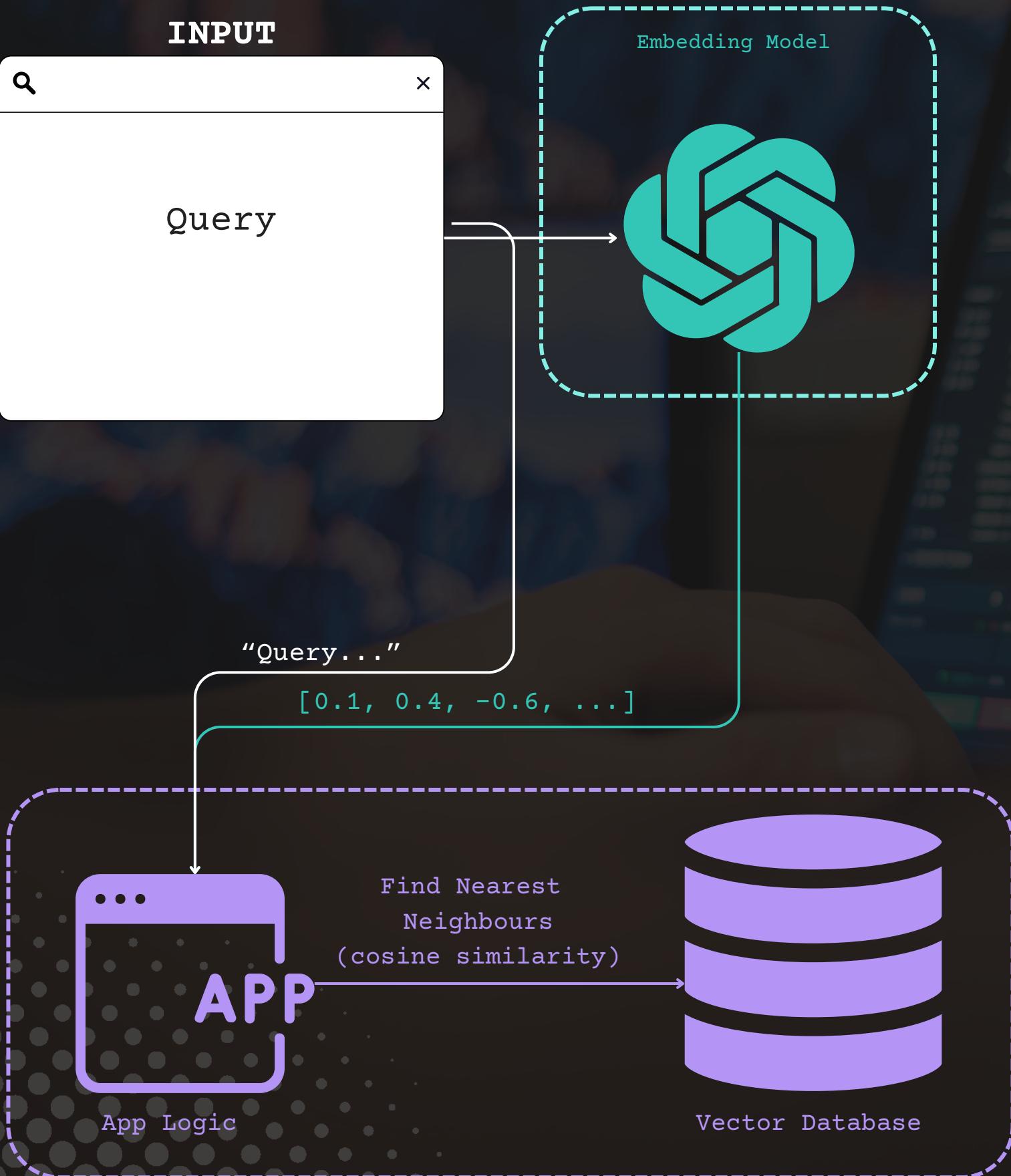


SP

Yes 😂

UNIVERSITYSTUDENT.org





Prompt Templates

Use the provided context to answer the user's query.

You may not answer the user's query unless there is specific context in the following text.

If you do not know the answer, or cannot answer, please respond with "I don't know".

Context:
{context}

User Query:
{user_query}

INPUT

Query

Embedding Model



Chat Model



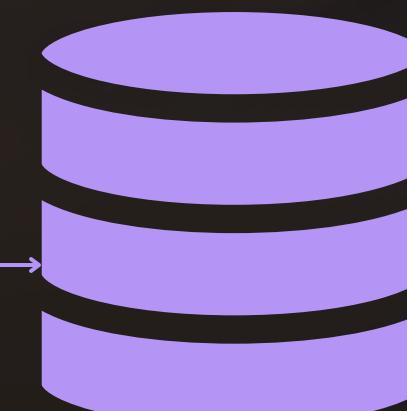
"Query..."

[0.1, 0.4, -0.6, ...]



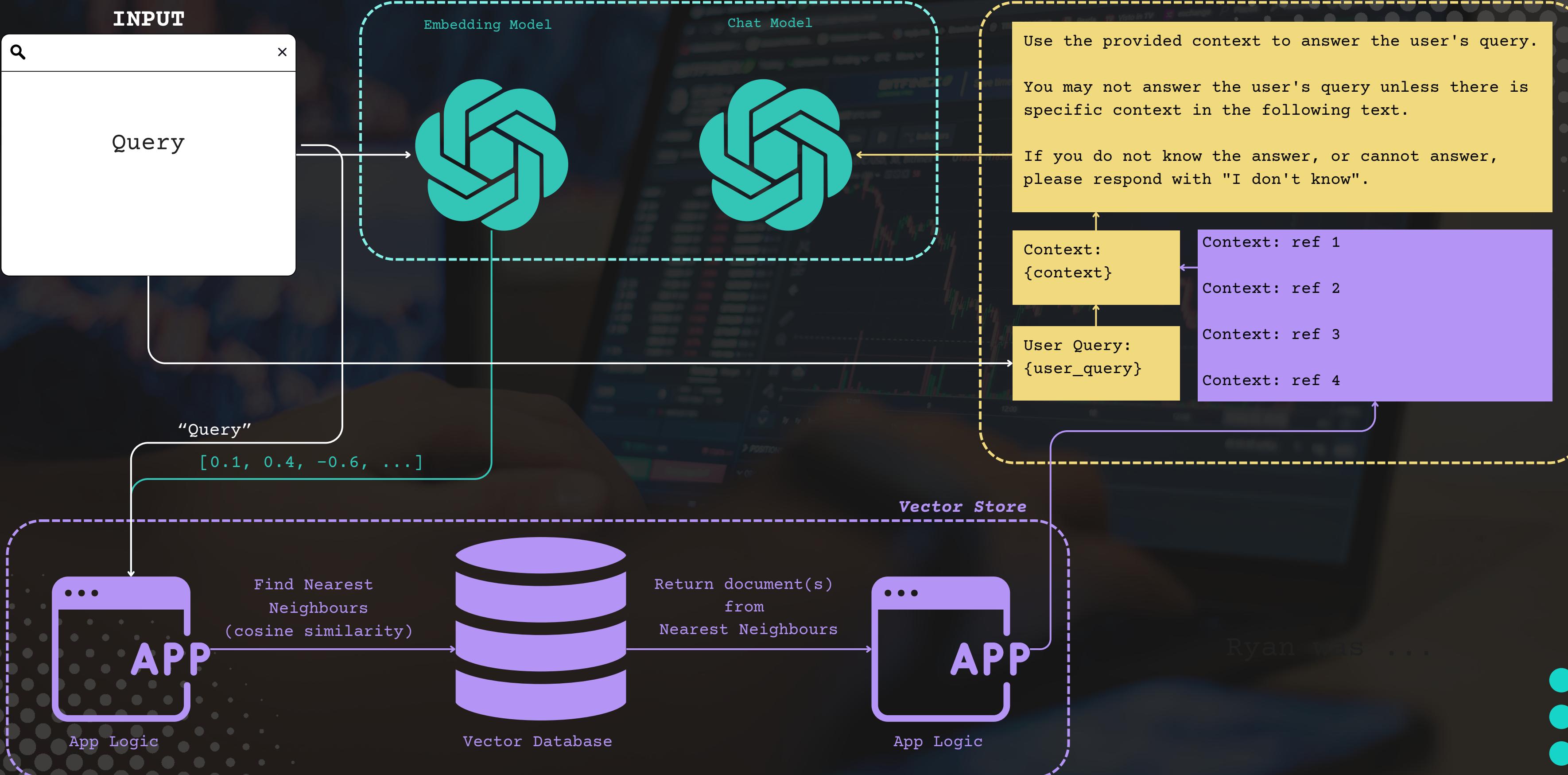
App Logic

Find Nearest
Neighbours
(cosine similarity)

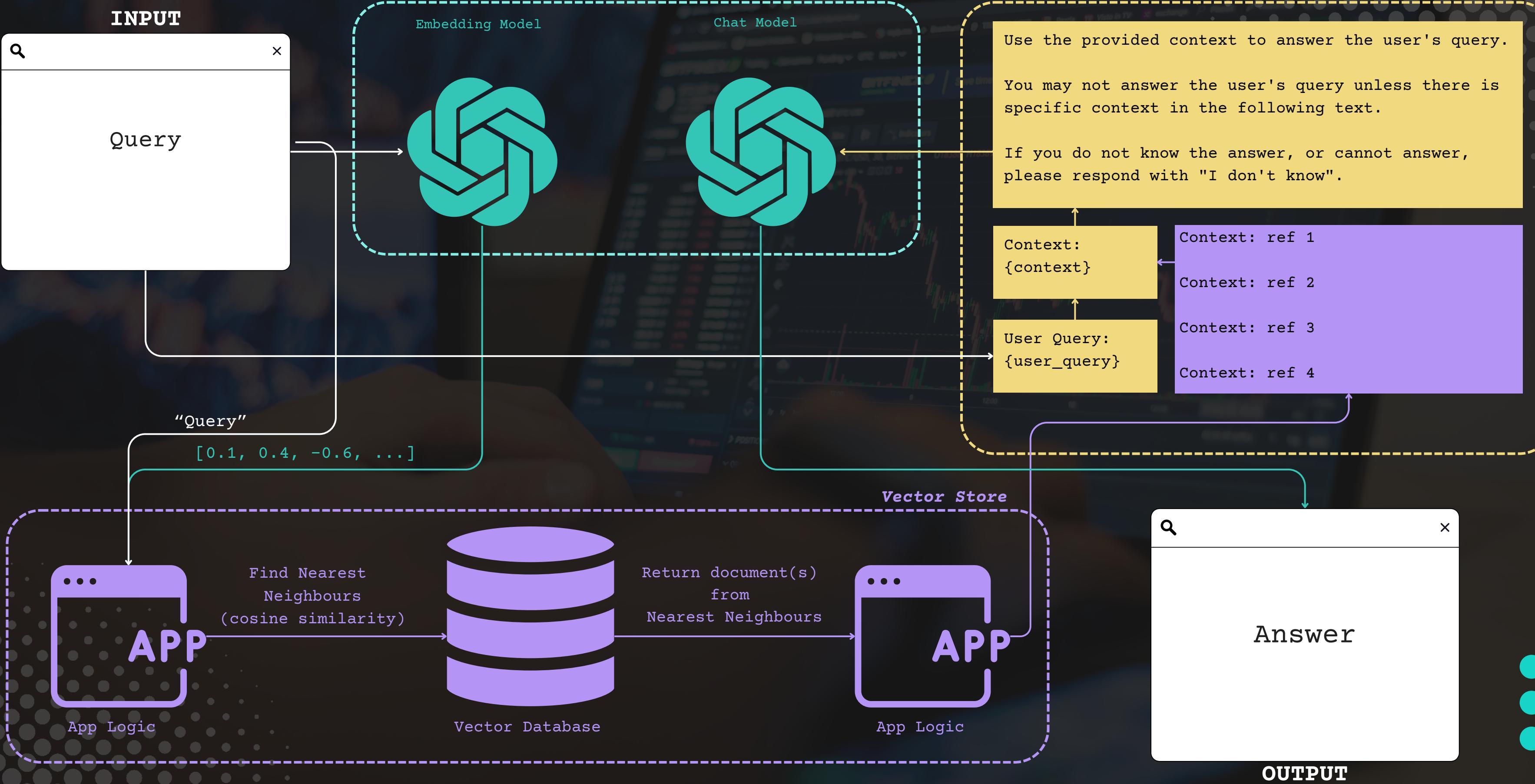


Vector Database

Prompt Templates



Prompt Templates



Prompt Templates

Use the provided context to answer the user's query.

You may not answer the user's query unless there is specific context in the following text.

If you do not know the answer, or cannot answer, please respond with "I don't know".

Context:
{context}

User Query:
{user_query}

Context: ref 1

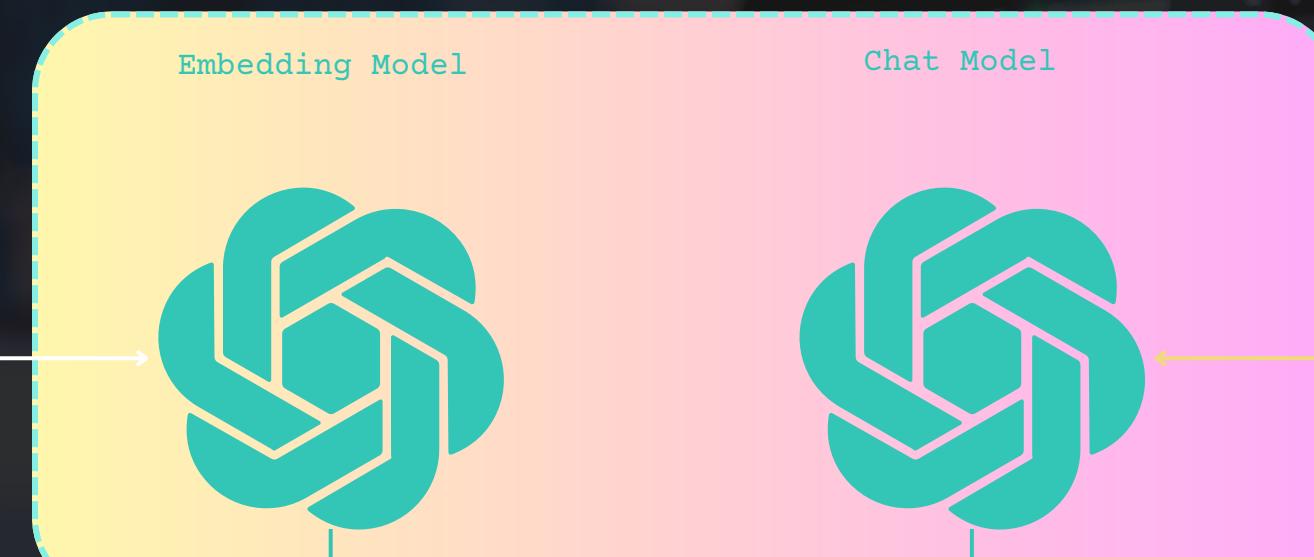
Context: ref 2

Context: ref 3

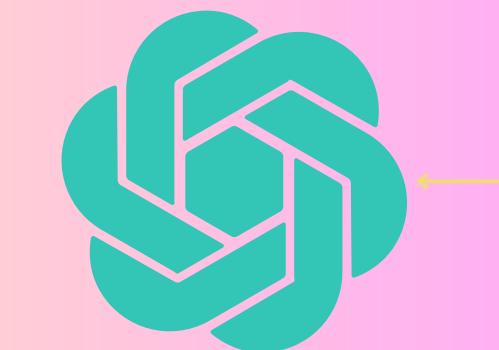
Context: ref 4

INPUT

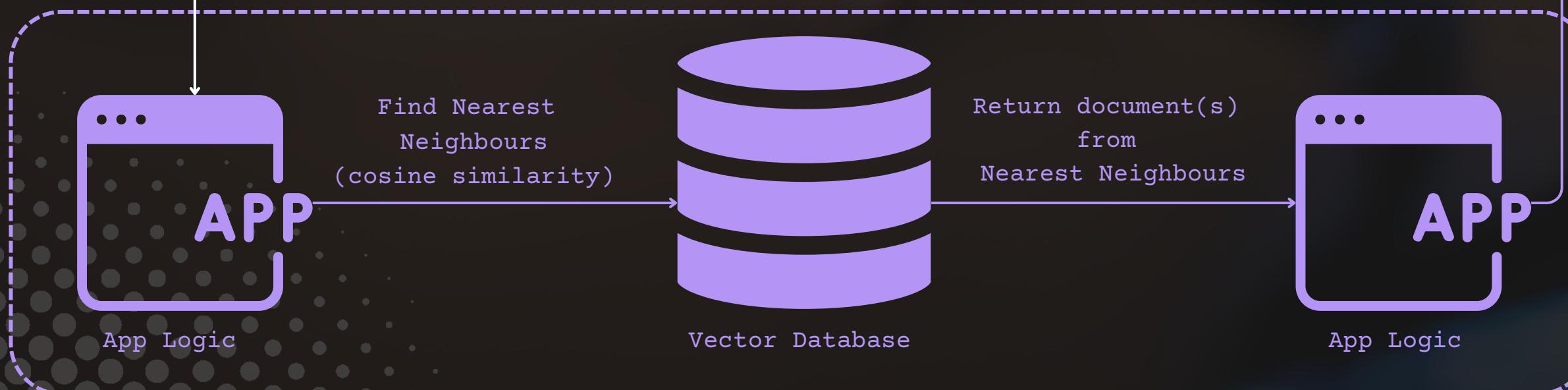
Query



Chat Model

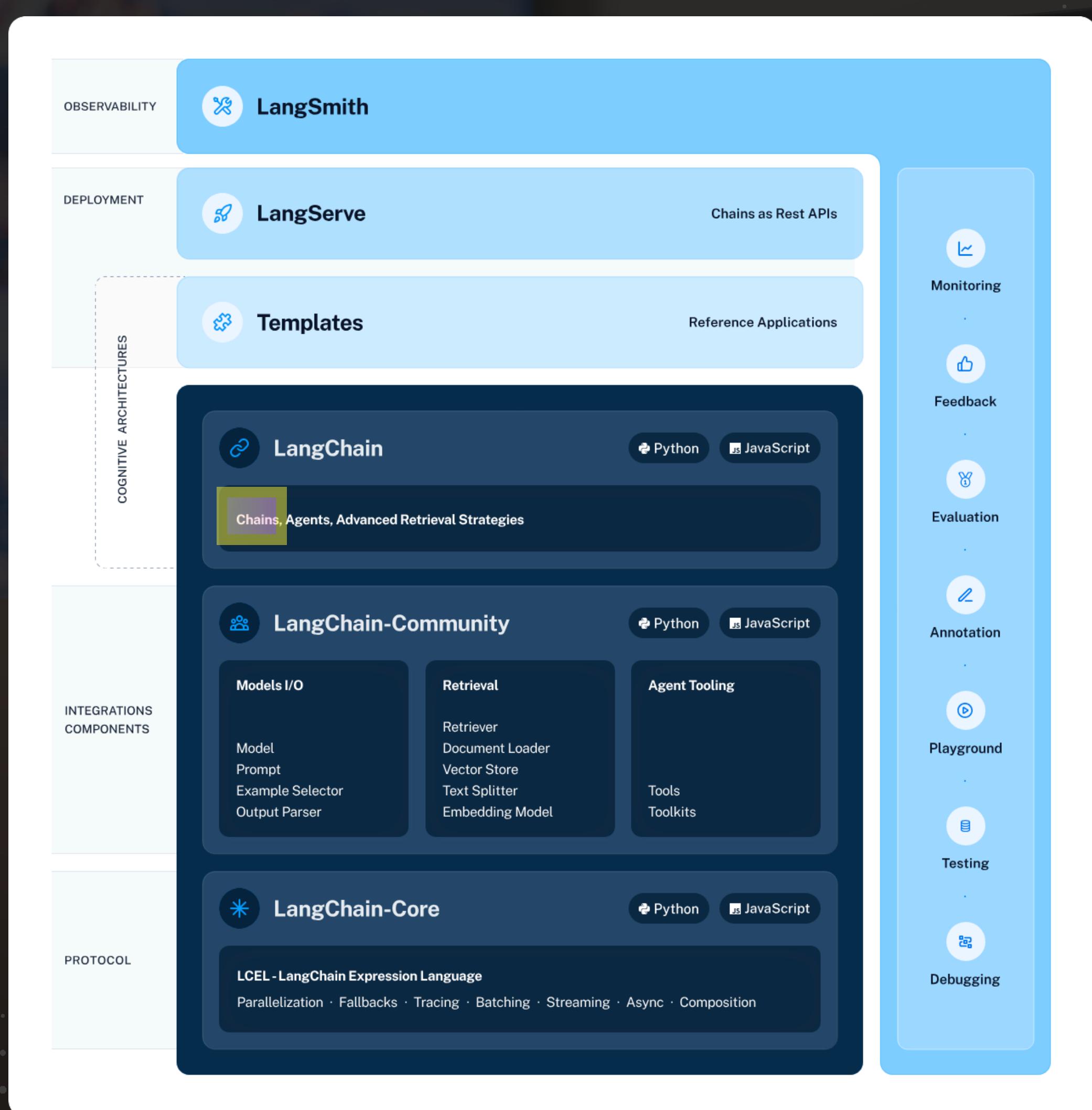


Vector Store



Answer

OUTPUT



DEF: CHAIN

A sequence of calls to other components

LLM Chain = Chat Model + Prompt Template

```
chain = LLMChain(llm=chat_model, prompt=chat_prompt)
```

```
llm_chain = prompt | llm
```



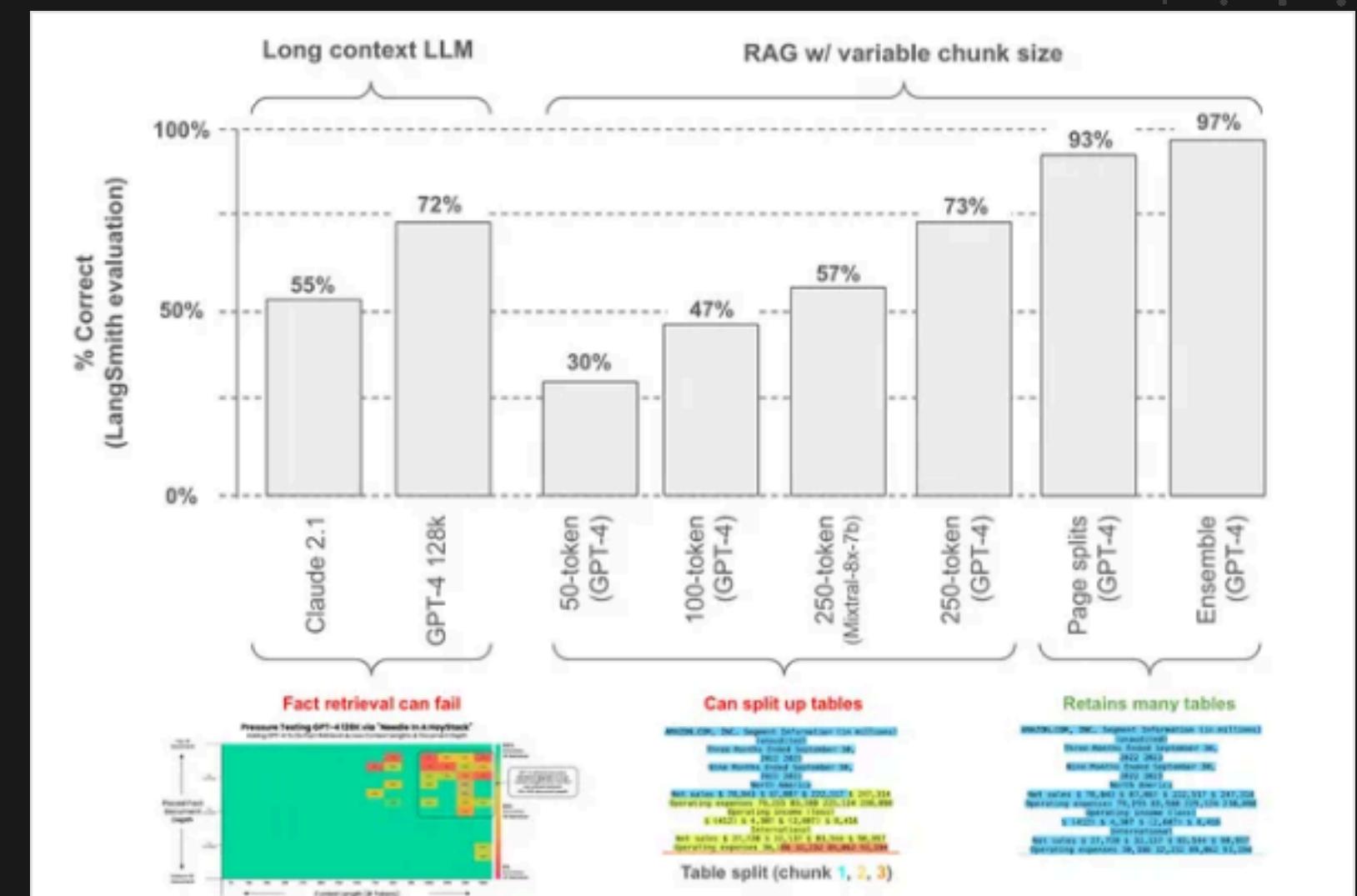


TODAY'S BUILD

LangChain Blog RAG Flow



1. **Search** LangChain blog docs for **top k** resources
2. **Ask** specific **questions** related to content
3. **Return answers** to questions with sources



Benchmarking RAG on tables

Key links LangChain public benchmark evaluation notebooks:
* Long context LLMs here
* Chunk size tuning here
* Multi vector with ensemble here
Motivation Retrieval augmented generation (RAG) is one of the mo...



OUR MODELS



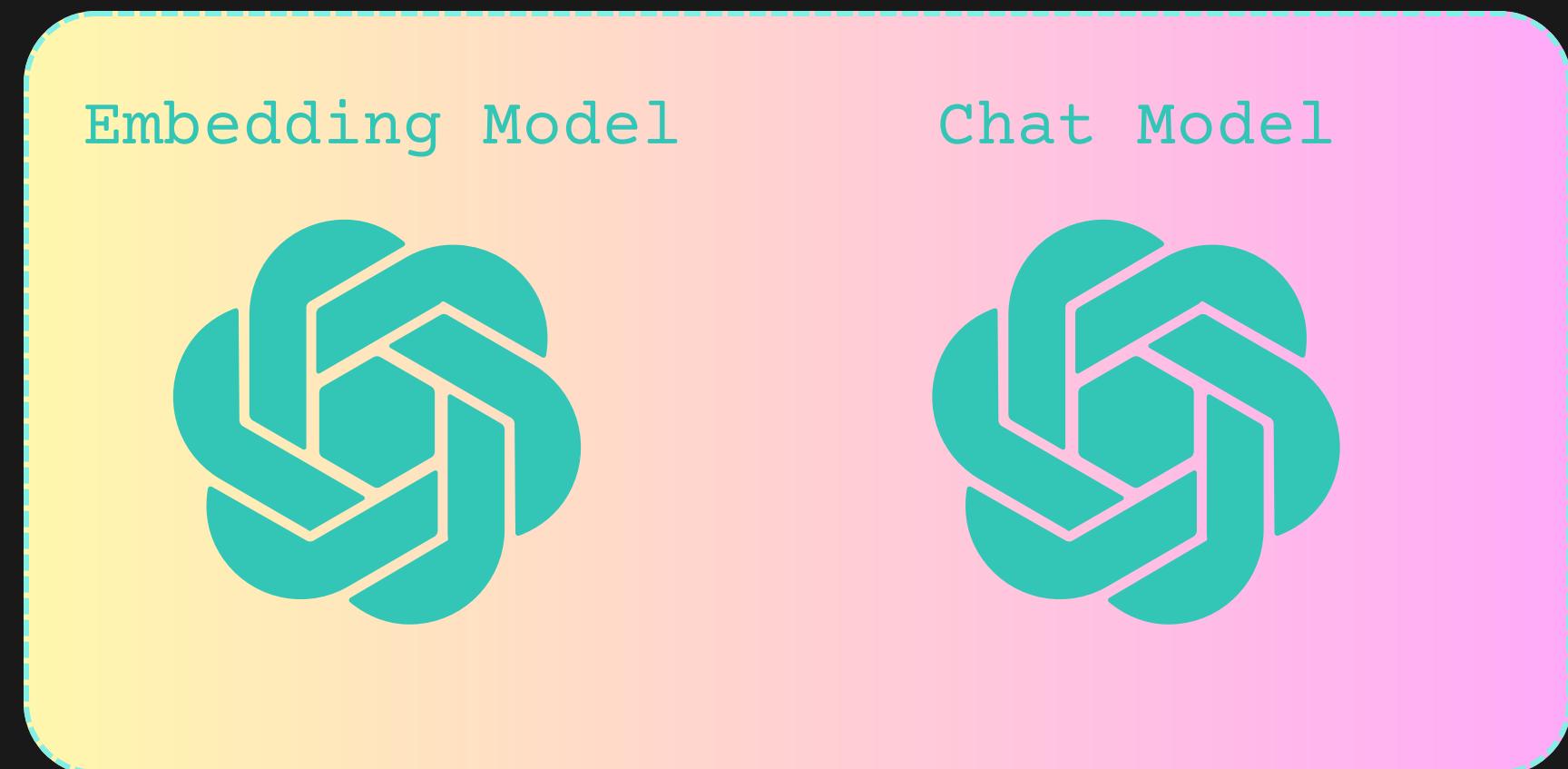
• EMBEDDING & CHAT (LLM) MODEL

Chat Model (e.g., LLM)

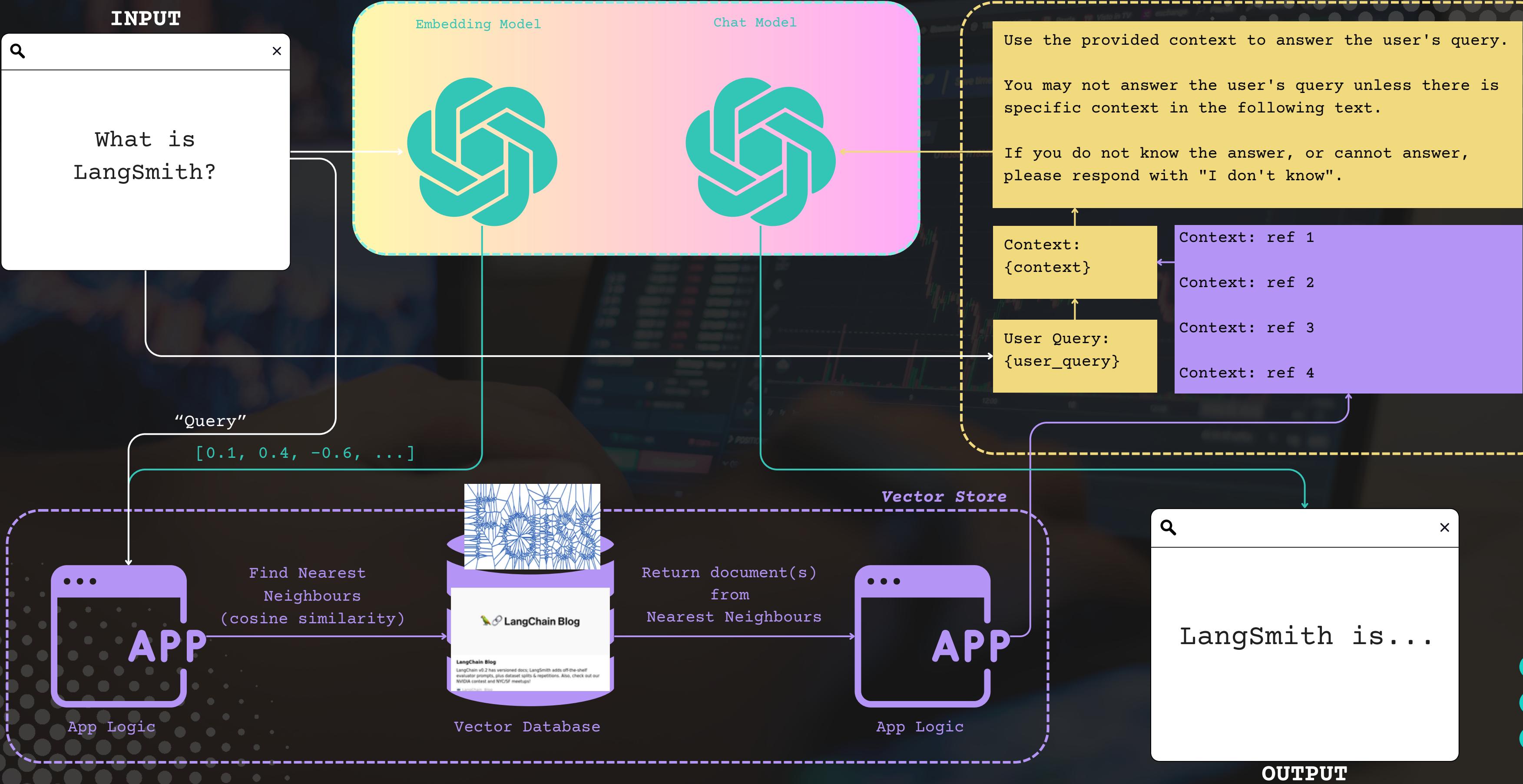
- e.g., OpenAI GPT-3.5 Turbo

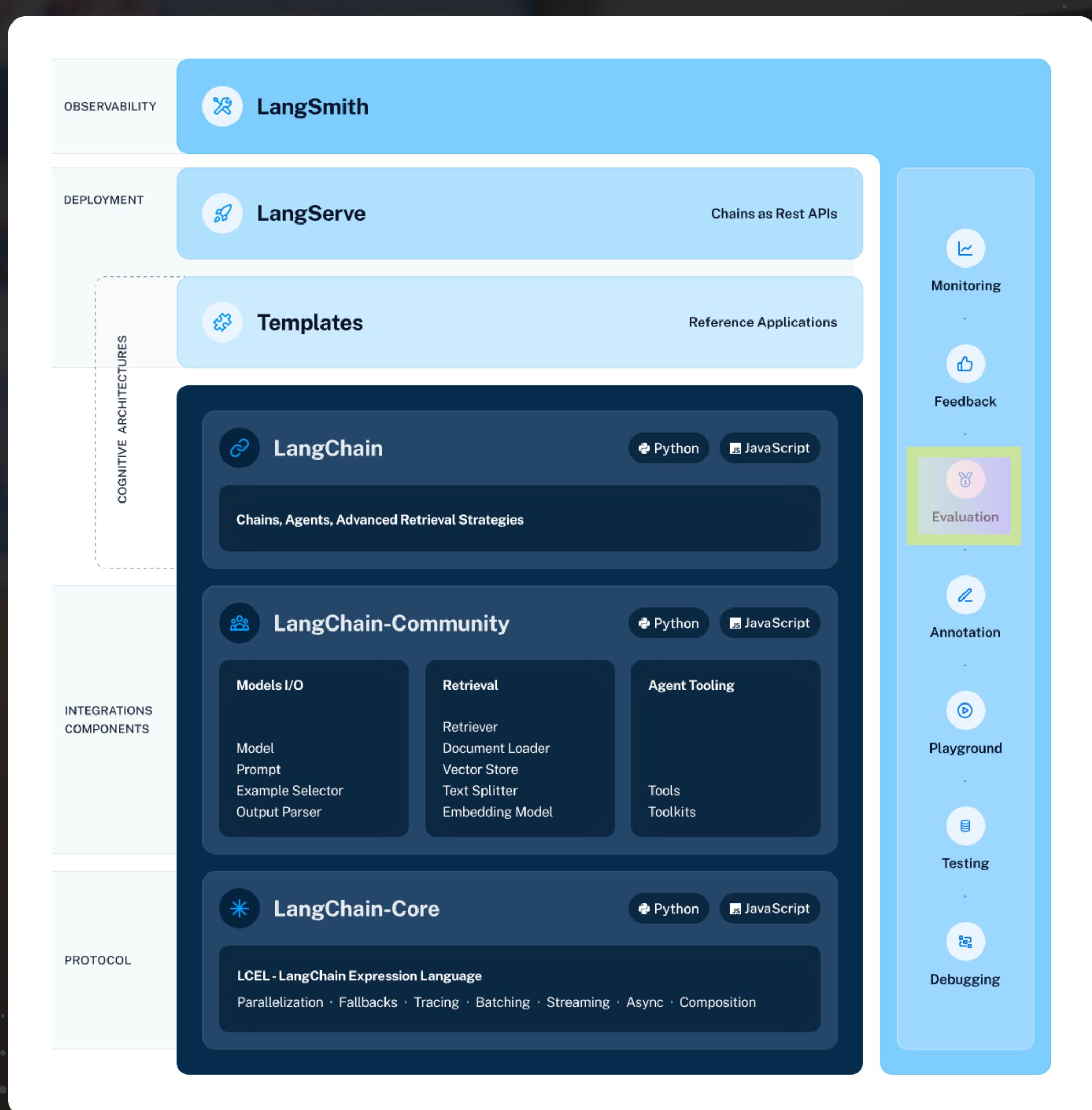
Embedding Model

- e.g., Open AI's Ada

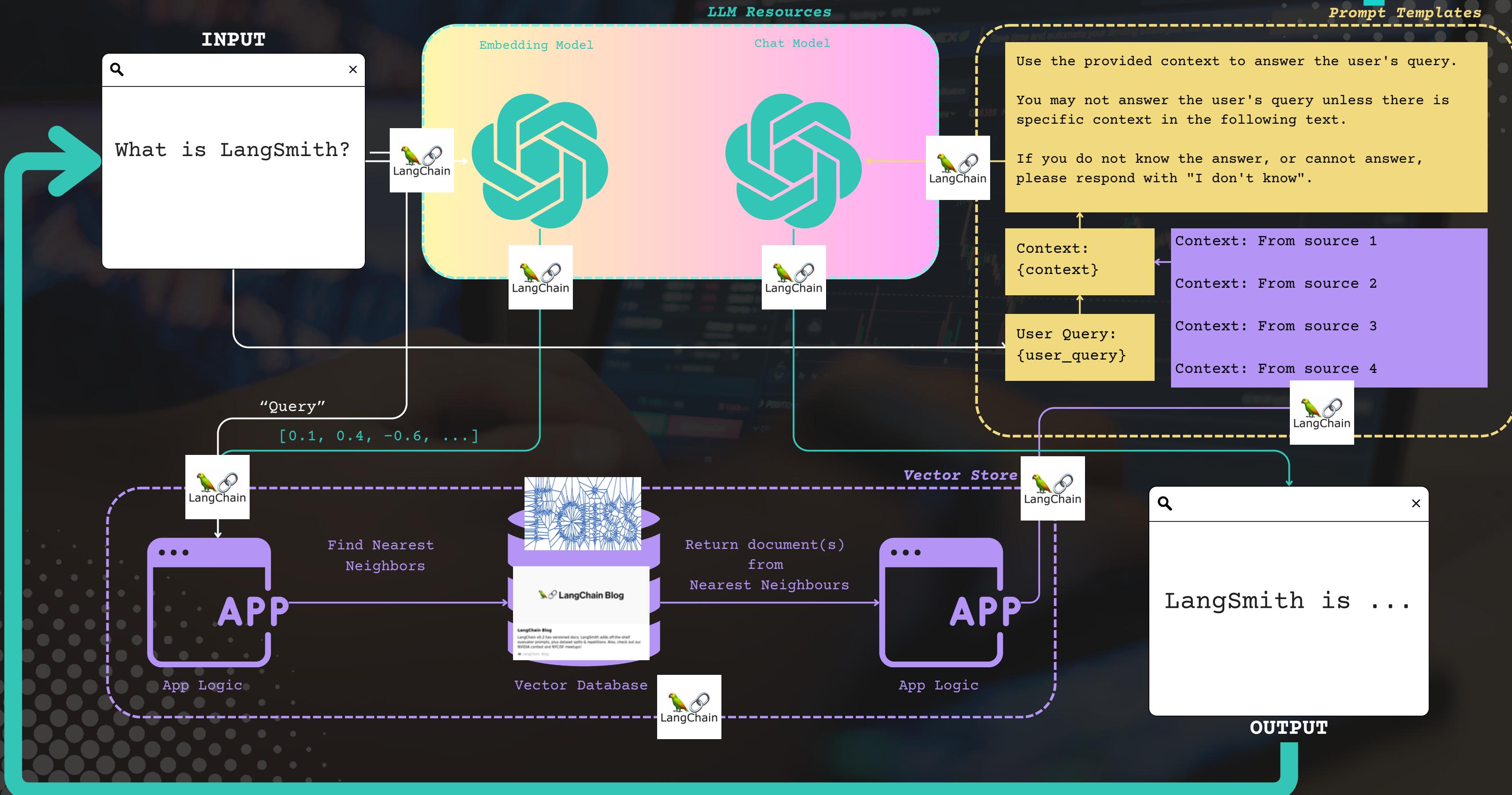


Prompt Templates





Evaluation Feedback Loop



RAG EVALUATION

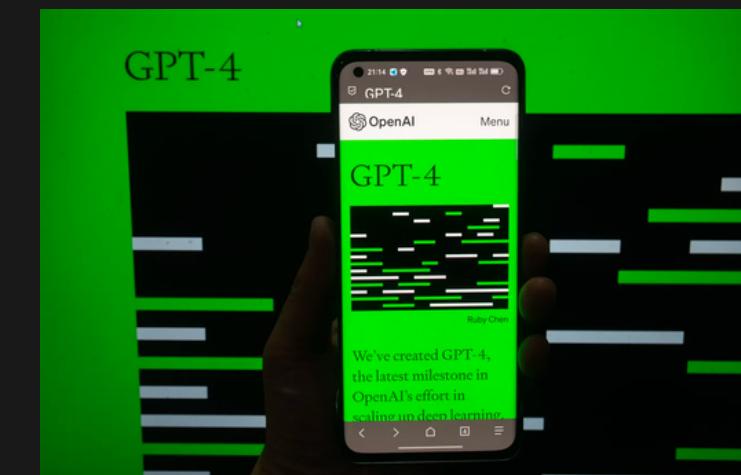
Generally we have 4 pieces of data...

Question

Result

Context

Answer



CHAIN-OF-THOUGHT QA

- **Instructs** the LLMChain to use **chain of thought** "reasoning"
- Responses **correlate better** with **human** labels
- **Higher** token and runtime **cost**.

```
_PROMPT_TEMPLATE = """You are an expert professor  
You are grading the following question:  
{query}  
Here is the real answer:  
{answer}  
You are grading the following predicted answer:  
{result}  
Respond with CORRECT or INCORRECT:  
Grade:  
"""  
  
PROMPT = PromptTemplate(  
    input_variables=["query", "answer", "result"],  
)
```

BASELINING RAG

Presented by
Chris Alexiuk, LLM Wizard ✨

DATA SCIENCE



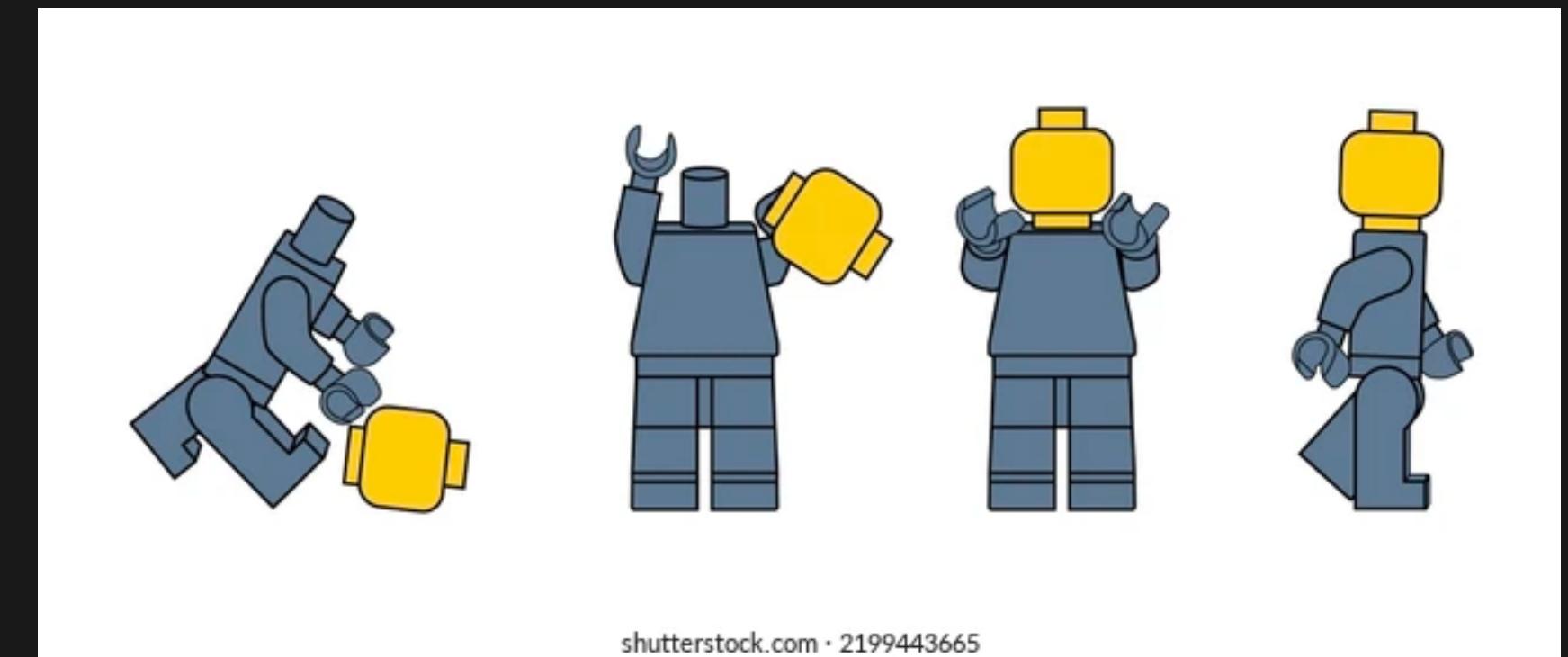
**BASIC
RETRIEVAL
AUGMENTED
GENERATION**

**CUSTOM RAG
#LLMOPS**

THE AGE OF THE AI ENGINEER

“A wide range of AI tasks that used to take 5 years and a research team to accomplish in 2013, now just require API docs and a spare afternoon in 2023.”

It is now possible to build what used to take months in a single day!

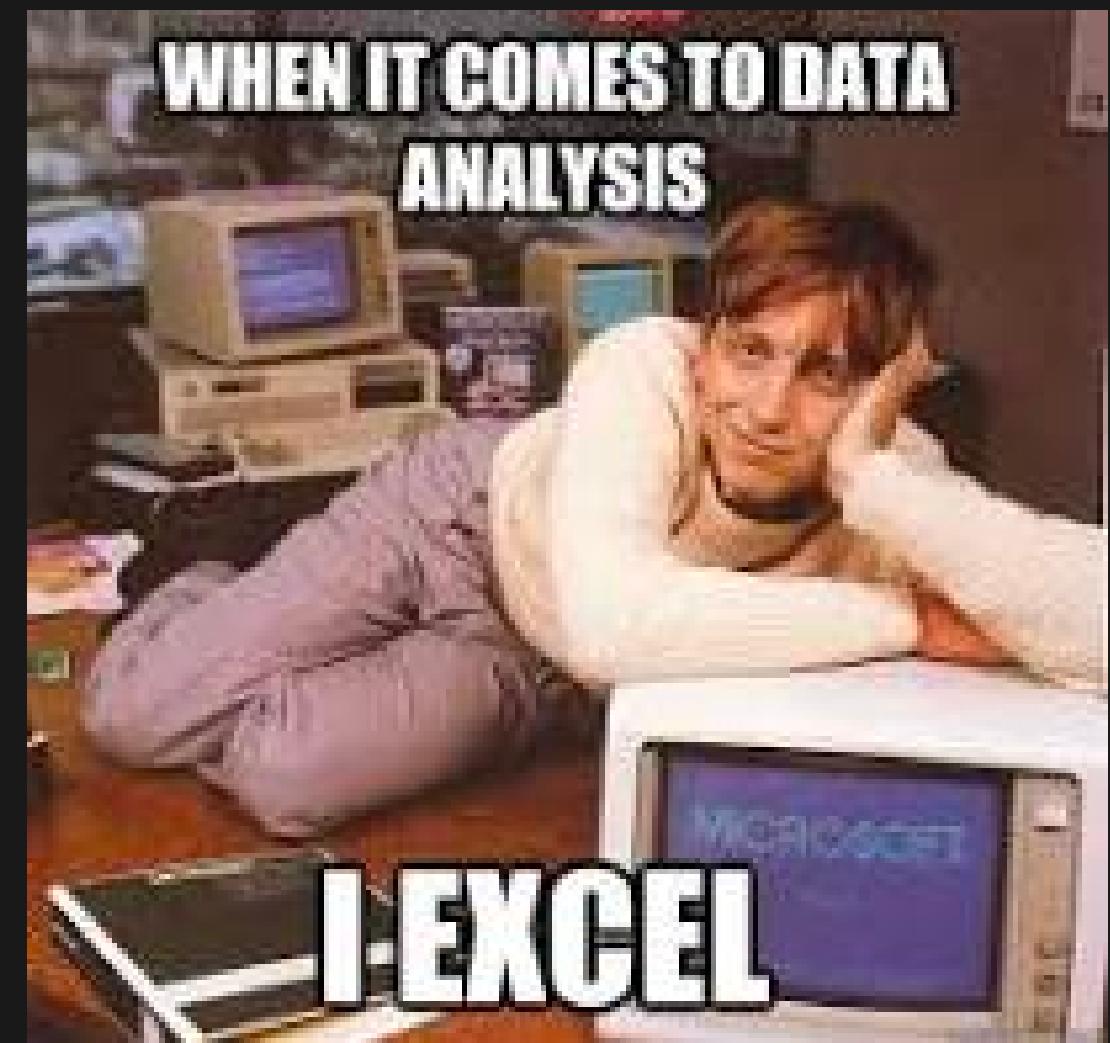


shutterstock.com · 2199443665



DATA SCIENTISTS!

- Evaluation
- Enhance Retrieval (and thus Generation!)
- Fine-Tuning
 - Embeddings
 - Chat Models



ADVANCED RAG

From Simple to Advanced RAG

Table Stakes

Better Parsers
Chunk Sizes
Hybrid Search
Metadata Filters



Less Expressive
Easier to Implement
Lower Latency/Cost

Advanced Retrieval

Reranking
Recursive Retrieval
Embedded Tables
Small-to-big Retrieval



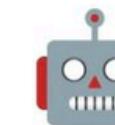
Fine-tuning

Embedding fine-tuning
LLM fine-tuning



Agentic Behavior

Routing
Query Planning
Multi-document Agents



More Expressive
Harder to Implement
Higher Latency/Cost

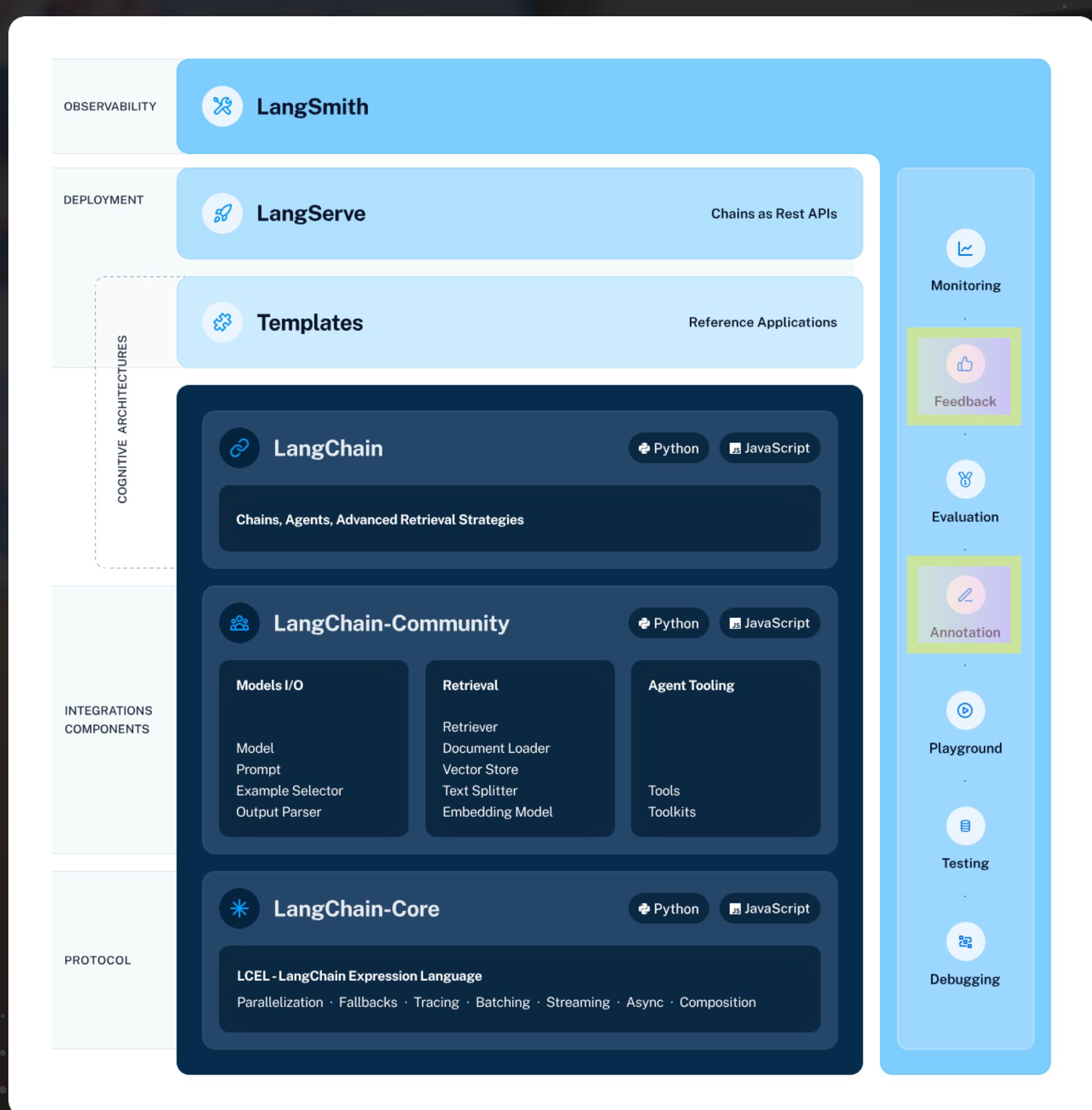


ENSEMBLE RETRIEVER

- A.k.a. “Hybrid Retrieval”
- Sparse retriever (BM25)
- Dense retriever (embedding similarity)
- Reciprocal Rank Fusion



```
from langchain.retrievers import BM25Retriever, EnsembleRetriever
```



⋮ ANNOTATION & FEEDBACK

- Feedback
 - As simple as  / 
 - As complex as engaging users for **detailed assessments**
- Annotation
 - Devs can provide **qualitative feedback** manually prior to next deployment!



MODEL ALIGNMENT WITH FINETUNING



⋮ PROTOTYPING/IMPROVING LLM APPS

1. Prompt Engineering
2. Question Answering Systems
3. Fine-Tuning Models



FEW-SHOT

- Zero-shot



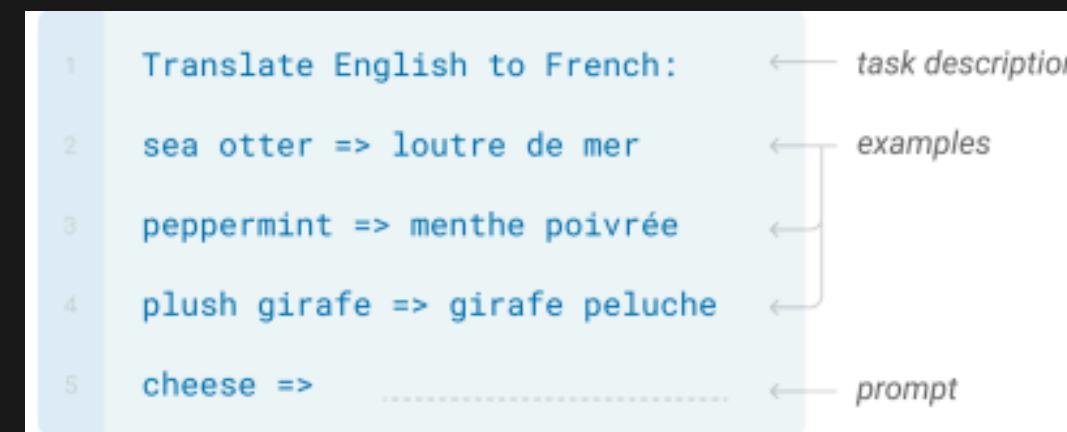
- **Instruction** only; closest to human

- One-shot

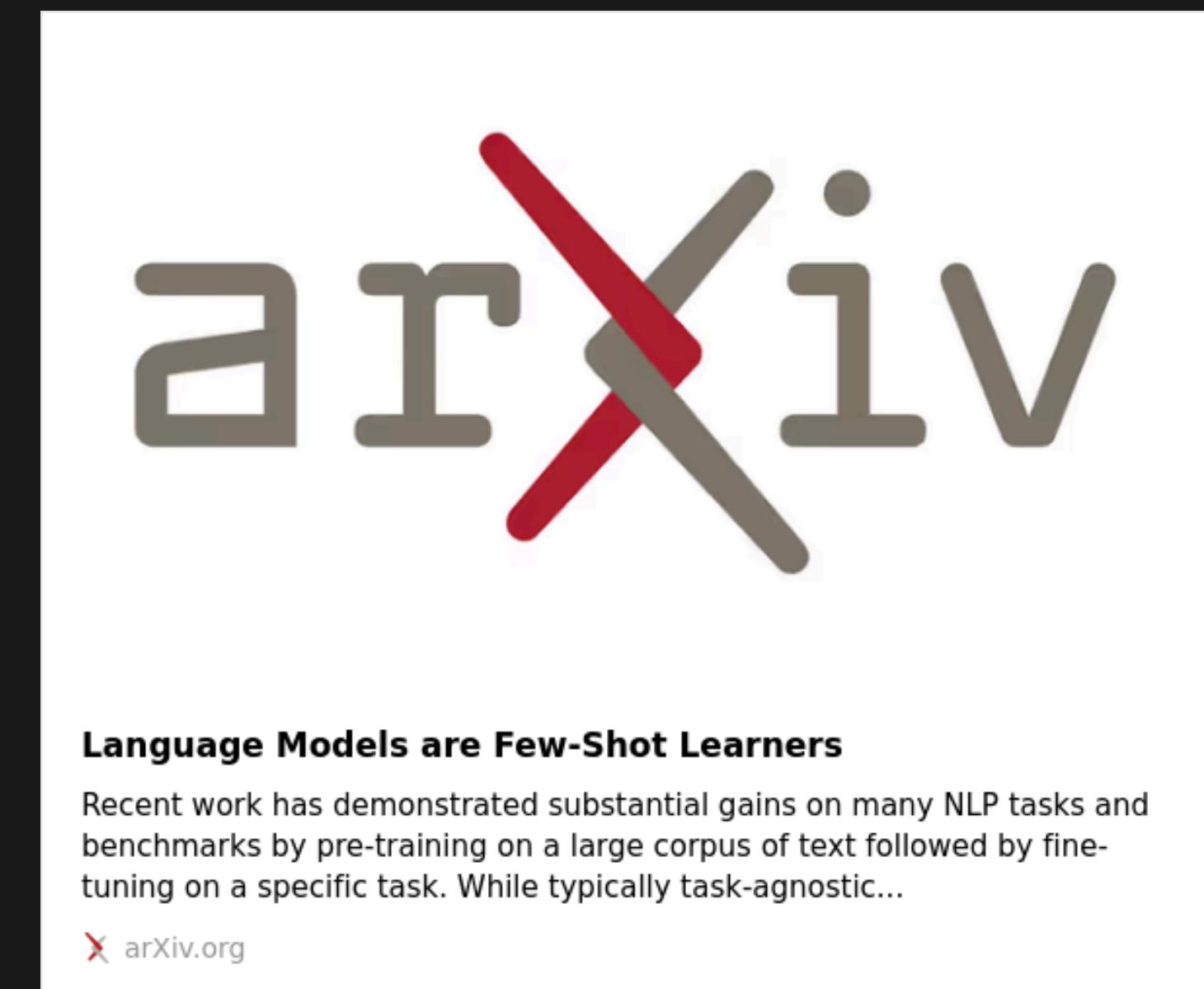


- Only **1** example

- Few-shot



- Given a “few” demonstrations (**10-100**)



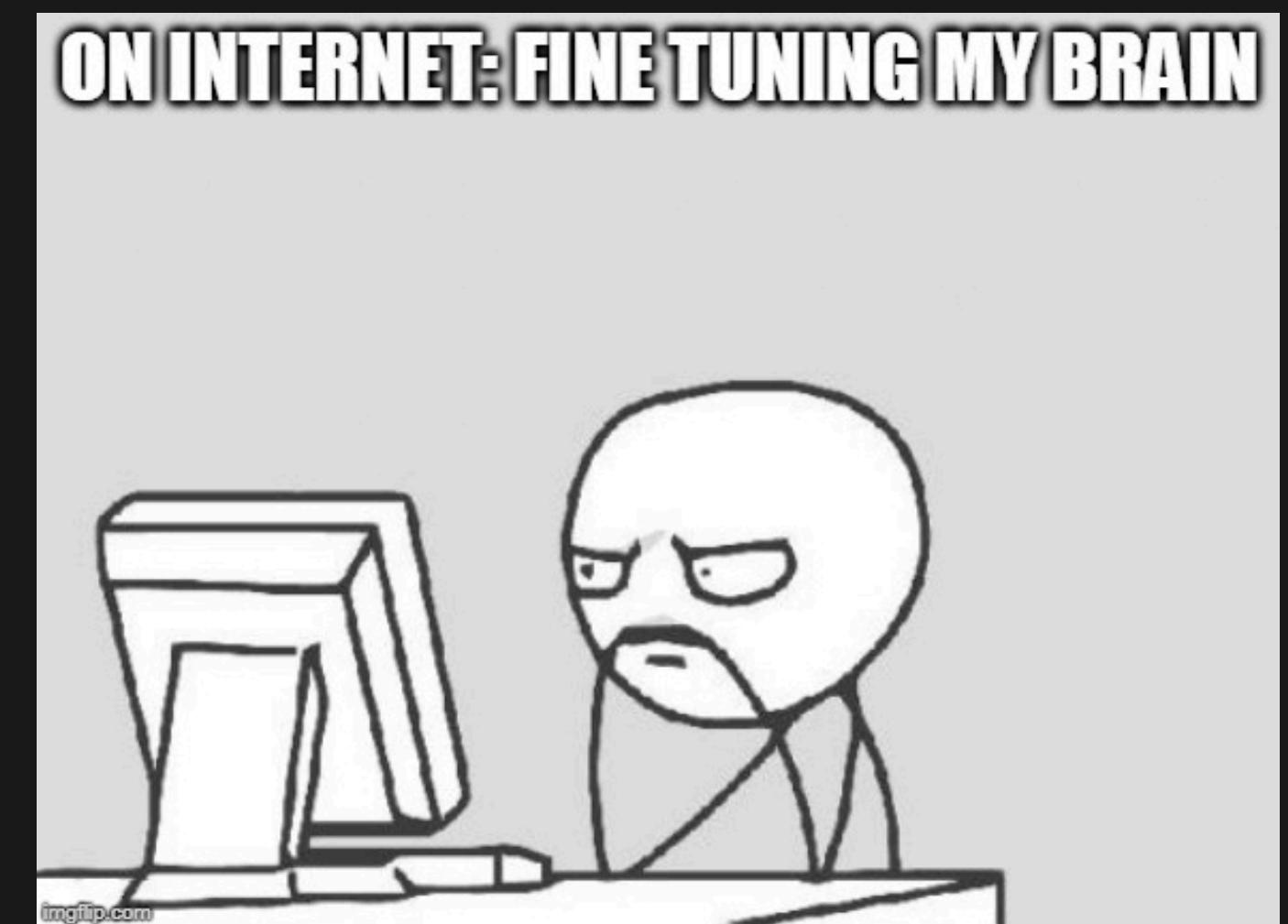
USER EXPERIENCE

- Control for **uncertainty** (PE)
- Build **guardrails** for steerability and safety (Harm/Help)



MODEL CONSISTENCY

- **Constrain** model behavior (FT)
- **Ground** the model (RAG)



The optimization flow

Context
optimization

What the model
needs to know

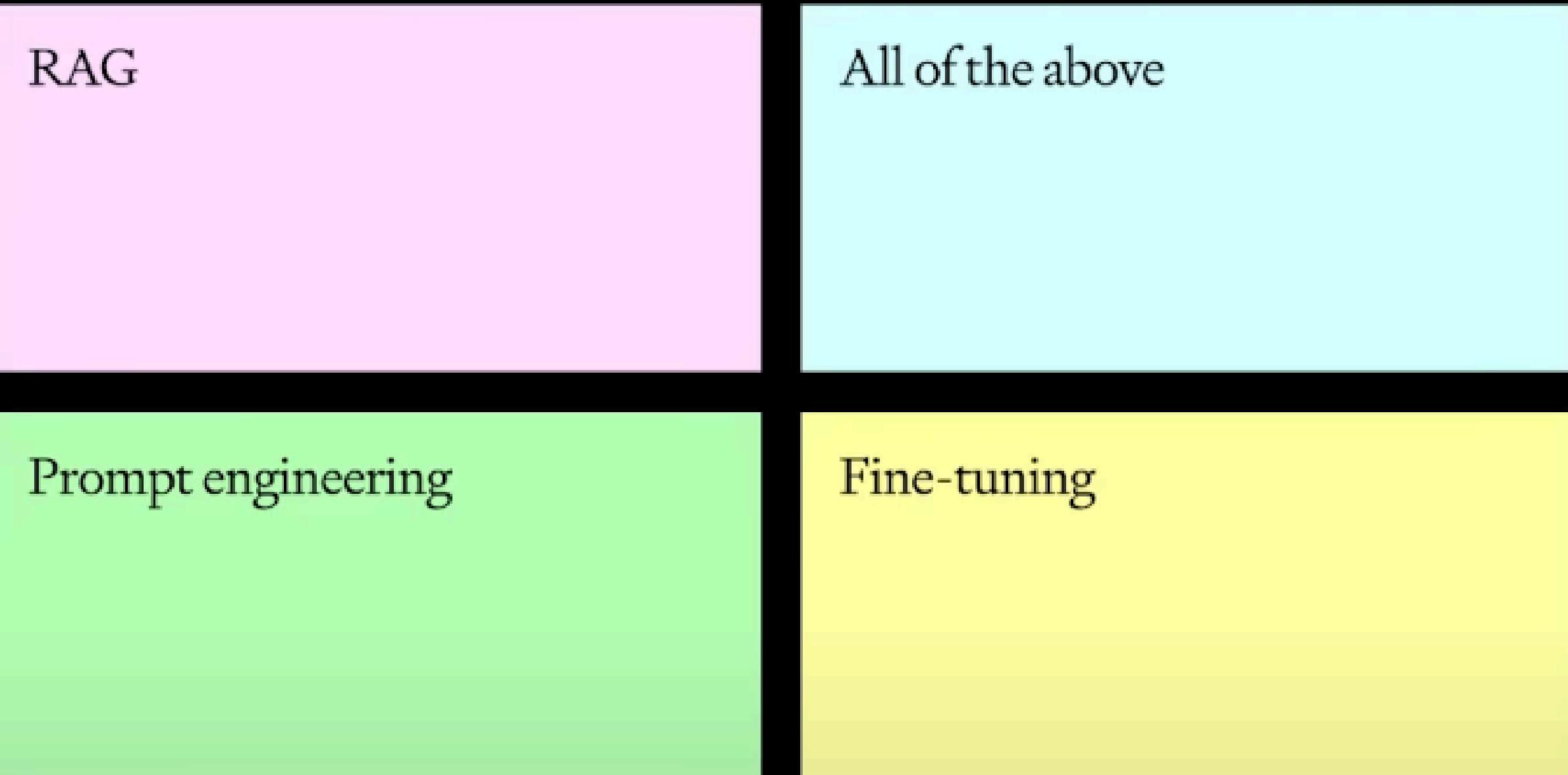
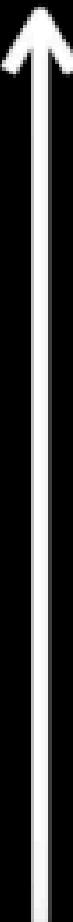


LLM optimization
How the model needs to act

The optimization flow

Context
optimization

What the model
needs to know

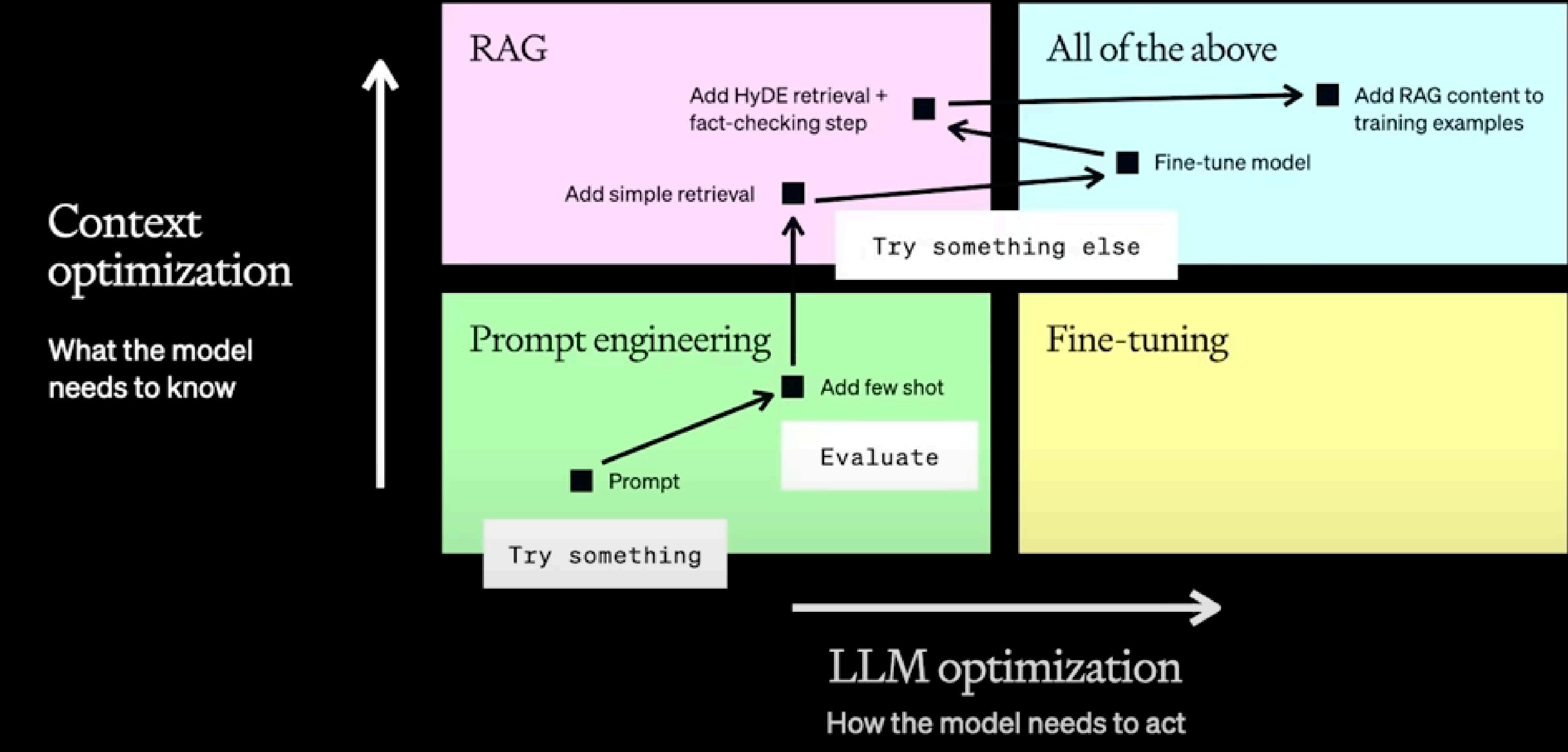


LLM optimization
How the model needs to act

The optimization flow

Context optimization

What the model needs to know



⋮ LATENCY AND COST

- Use **semantic** caching (Prompts)
- Route to **cheaper** models (FT)





IMPROVING RAG

Presented by
Chris Alexiuk, LLM Wizard ✨

CONCLUSIONS

- **LangServe** is solving end-to-end LLM application problems!
 - **Streamlines production improvements**, from Prompt Engineering to RAG to Fine-Tuning
- **Baseline evaluation** is key
 - LLMs as evaluators is an emerging space!
- From there, **improving** your system is the data science of the future!



OBSERVABILITY



LangSmith

DEPLOYMENT



LangServe

Chains as Rest APIs

APPLICATION



LangChain

Python

JavaScript

Common Application Logic

INTEGRATIONS
COMPONENTS

Models I/O

- Model
- Prompt
- Example Selector
- Output Parser

Retrieval

- Retriever
- Document Loader
- Vector Store
- Text Splitter
- Embedding Model

Agent Tooling

- Tools
- Toolkits

PROTOCOL

LCEL - LangChain Expression Language

Parallelization · Fallbacks · Tracing · Batching · Streaming · Async · Composition

Monitoring

Evaluation

Annotation

Feedback

Testing

Debugging

QUESTION?



Thank you!