# LLM Observability
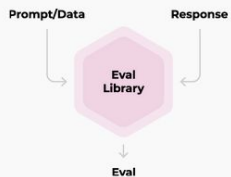
Amber Roberts, ML Growth Lead @ Arize AI

# LLM Observability Pillars
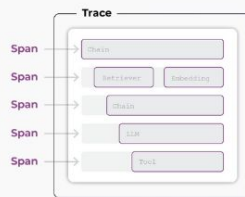
## Evaluation
*Evaluations of LLM outputs by using a separate evaluation LLM*



## Traces & Spans
*Visibility into where the agentic workflow broke*



## Prompt Engineering
*Iterating on prompt templates for improved results*



## Search & Retrieval
*Locate and improve retrieved context*



## Fine-tuning
*Re-train LLM on use case / company data*

# LLM Observability Pillars

## Evaluation
*Evaluations of LLM outputs by using a separate evaluation LLM*



## Traces & Spans
*Visibility into where the agentic workflow broke*



## Prompt Engineering
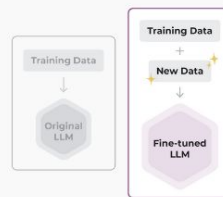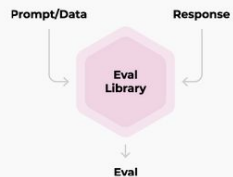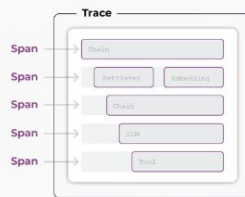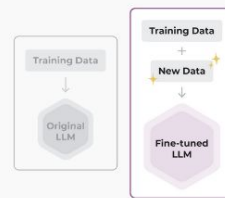*Iterating on prompt templates for improved results*



## Search & Retrieval
*Locate and improve retrieved context*



## Fine-tuning
*Re-train LLM on use case / company data*

# How does RAG (Retrieval Augmented Generation) work?

**Knowledge base of articles**

**User query**

*Do you support international calling?*

**Query embedding**

**<1, 2, 3, 4>**

**Cosine Similarity From Lookup**

0.8

0.4

0.1

| Document Chunk Embedding | Document Chunk ID | Document Chunk |
|---|---|---|
| **<1, 1, 2, 4>** | **1** | **What countries support International Calling**<br>• The International Calling feature is available to all countries when it is enabled. See **How to Enable International Calling** for more information.<br>• Some countries may not be available when International Calling is enabled. This means that RingCentral has restricted International Calling to those countries. **Contact Support** if you need to reach any of the restricted countries. See **How to Open a RingCentral Tech Support Case** for more information. |
| <100, 309, 4, 7> | 2 | **Configuring Outbound Call Prefix**<br>1. Contact Support to enable Outbound Call Prefix on your account.<br>2. Sign in to the RingCentral admin portal.<br>3. Navigate to **Admin Portal > More > Account Settings > Outbound Call Prefix**.<br>4. Toggle to enable **Use outgoing call prefix**.<br>5. Enter the single-digit that users will dial before dialing an external number.<br>6. Click **Validate & Save**. |
| <59, 71, 73, 95> | 3 | |

**Prompt**

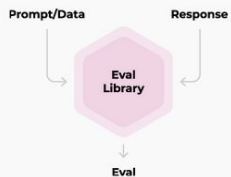User is asking "*Do you support international calling?*"

Here's relevant content. Can you answer?

**What countries support International Calling**
• The International Calling feature is available to all countries when it is enabled. See How to Enable International Calling for more information.
• Some countries may not be available when International Calling is enabled. This means that RingCentral has restricted International Calling to those countries. Contact Support if you need to reach any of the restricted countries. See How to Open a RingCentral Tech Support Case for more information.
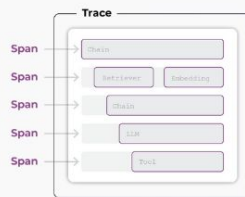
# LLM Observability Pillars

## Evaluation
*Evaluations of LLM outputs by using a separate evaluation LLM*



## Traces & Spans
*Visibility into where the agentic workflow broke*



## Prompt Engineering
*Iterating on prompt templates for improved results*



## Search & Retrieval
*Locate and improve retrieved context*



## Fine-tuning
*Re-train LLM on use case / company data*

# Traces and Spans

- **LLM usage is documented in a callback system by a trace.**

- **In this trace a span can refer to any unit of execution**

- **You may annotate a span with a specific name or a general term like a chain.**



**Chains/Pipelines**

Chain

Retriever — Emb

Chain

LLM

Tool

# Span Types

**LLM**
Call to a LLM for completion or chat

**Chain**
Link between application steps

**Tool**
API or Function invoked on behalf of an LLM

**Agent**
Root of a set of LLM and Tool invocations

**Embedding**
Encoding of unstructured data

**Retriever**
Query for context from a data store

**Reranker**
Relevance-based re-ordering of documents

# LLM Observability Pillars

**Evaluation**
*Evaluations of LLM outputs by using a separate evaluation LLM*

**Traces & Spans**
*Visibility into where the agentic workflow broke*

**Prompt Engineering**
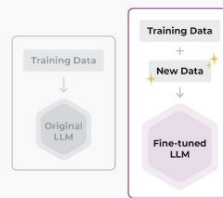*Iterating on prompt templates for improved results*

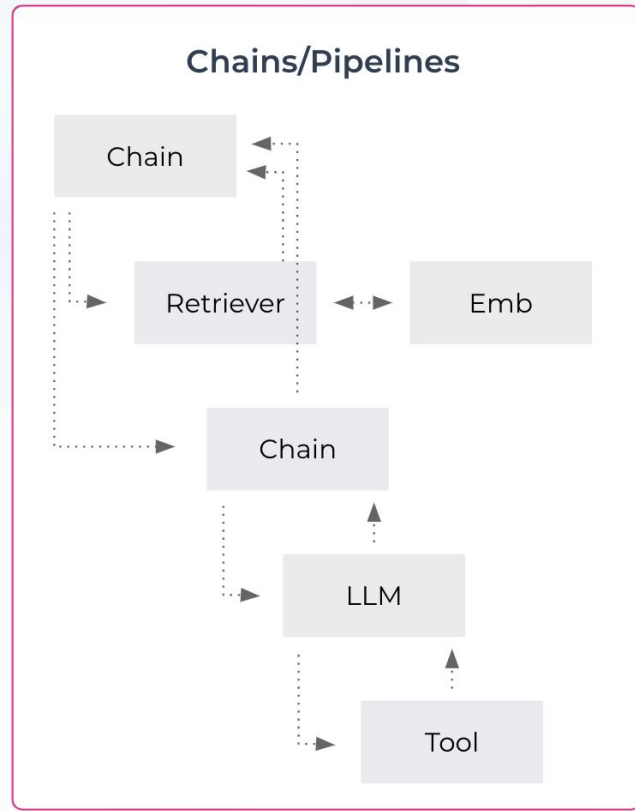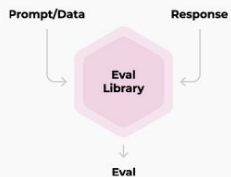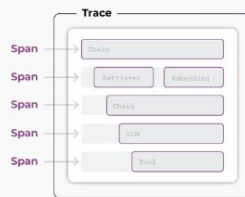**Search & Retrieval**
*Locate and improve retrieved context*

**Fine-tuning**
*Re-train LLM on use case / company data*

# Model vs. System Evals

# LLM Model Evals

| INPUT DATA | PROMPT TEMPLATE | MODEL TESTED | OUTPUT |
|---|---|---|---|
| Prompt 1 | | → **Llama** → | Response 1 |
| Prompt 2 | "You are an AI assistant. Respond succinctly." | | Response 2 |
| Prompt 3 | | | Response 3 |

| INPUT DATA | | MODEL TESTED | OUTPUT |
|---|---|---|---|
| Prompt 1 | | → **Vicuna** → | Response 1 |
| Prompt 2 | | | Response 2 |
| Prompt 3 | | | Response 3 |

# OpenAI Eval library

- **The OpenAI Eval library and LLM leaderboards look at how well the various LLMs stack up against each other.**
- **Using LLM evaluation metrics like:**
  - HellaSwag – how well an LLM can complete a sentence
  - TruthfulQA - truthfulness of responses
  - MMLU - how well the LLM can multitask.

# LLM System Evals

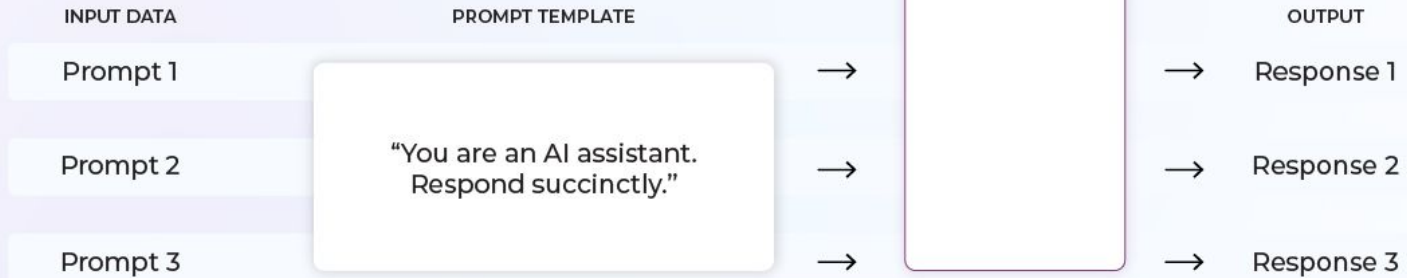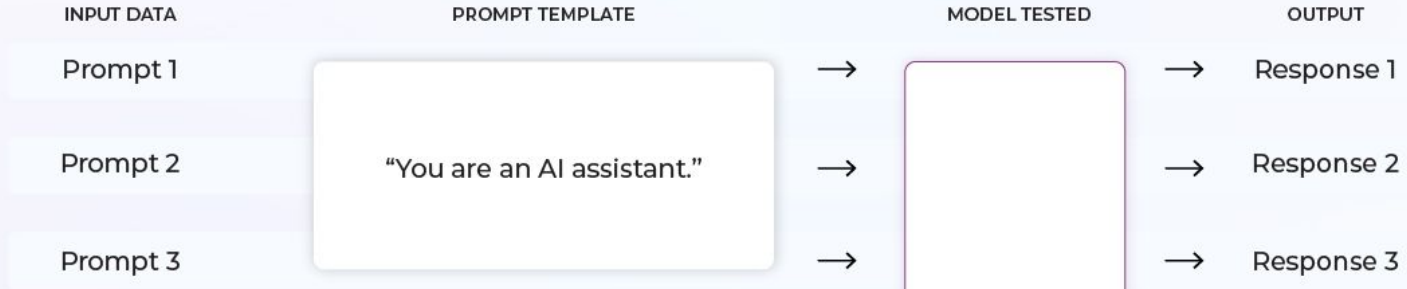| INPUT DATA | PROMPT TEMPLATE | MODEL TESTED | OUTPUT |
|------------|-----------------|--------------|--------|
| Prompt 1 | | | Response 1 |
| Prompt 2 | "You are an AI assistant." → | | Response 2 |
| Prompt 3 | | | Response 3 |

**Llama**

| INPUT DATA | PROMPT TEMPLATE | | OUTPUT |
|------------|-----------------|--------------|--------|
| Prompt 1 | | | Response 1 |
| Prompt 2 | "You are an AI assistant. Respond succinctly." → | | Response 2 |
| Prompt 3 | | | Response 3 |

# Response vs. Retrieval Evals

# Retrieval Evals vs Response Evals



Knowledge base
Pinecone
Milvus
Chroma

User query

Query embedding

Search & Retrieval

Vector store → Prompt With Context and User Query

LLM → Response → User feedback

Sample of the vector store

Prompt & prompt template

Retrieved context

**Precision @4 = 80%
NDCG = 75%
Hit = True**

LLM response

**Hallucinations = False
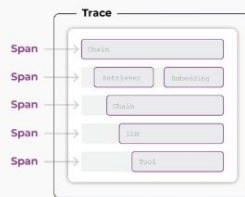Correctness = True**

# Putting it all together

### Evaluation
*Evaluations of LLM outputs by using a separate evaluation LLM*



### Traces & Spans
*Visibility into where the agentic workflow broke*



### Prompt Engineering
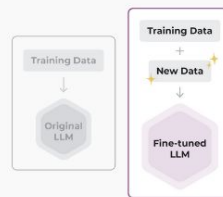*Iterating on prompt templates for improved results*



### Search & Retrieval
*Locate and improve retrieved context*



### Fine-tuning
*Re-train LLM on use case / company data*

# Arize Phoenix Overview

4–5 lines of code  ⟶  10+ LLM calls

```python
from langchain.chains import RetrievalQA
from langchain.chat_models import ChatOpenAI
from langchain.embeddings import OpenAIEmbeddings
from langchain.retrievers import KNNRetriever

embeddings = OpenAIEmbeddings(model="text-embedding-ada-
002")
knn_retriever = KNNRetriever(
    index=vectors,
    texts=texts,
    embeddings=OpenAIEmbeddings(),
)

llm = ChatOpenAI(model_name="gpt-3.5-turbo")
chain = RetrievalQA.from_chain_type(
    llm=llm,
    chain_type="map_reduce",
    retriever=knn_retriever,
)
```

Making sense of a large number of distributed system calls is what Phoenix is designed to do.

chatbot

embedding

retriever

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce

LLM call - map reduce
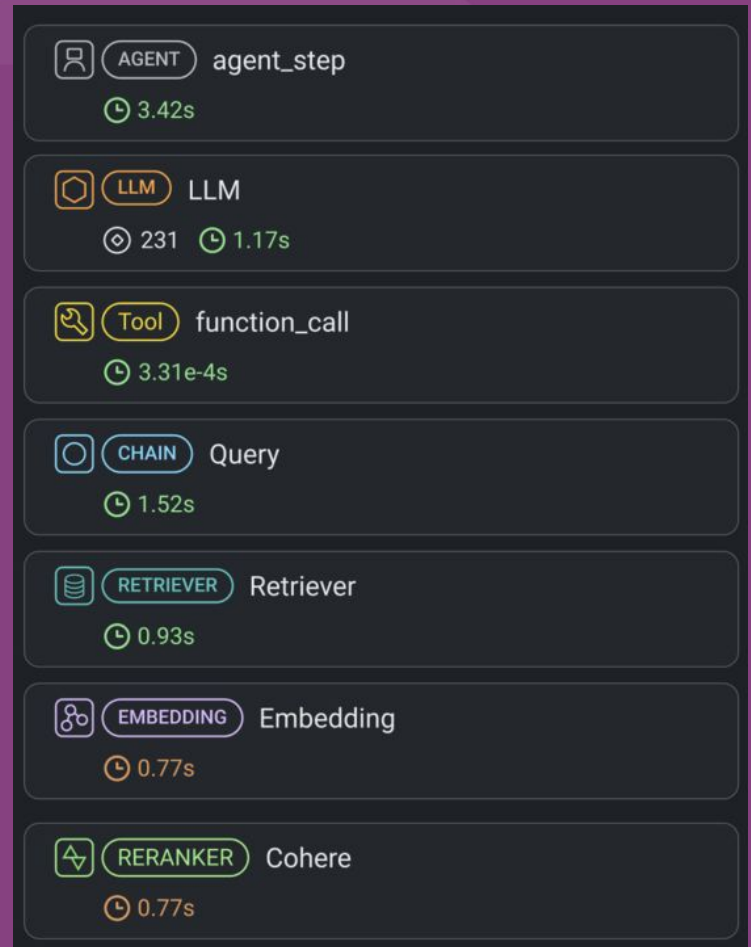
...

# Arize Phoenix Overview

4–5 lines of code $\longrightarrow$ 10+ LLM calls

```python
from langchain.chains import RetrievalQA
from langchain.chat_models import ChatOpenAI
from langchain.embeddings import OpenAIEmbeddings
from langchain.retrievers import KNNRetriever

embeddings = OpenAIEmbeddings(model="text-embedding-ada-
002")
knn_retriever = KNNRetriever(
    index=vectors,
    texts=texts,
    embeddings=OpenAIEmbeddings(),
)

llm = ChatOpenAI(model_name="gpt-3.5-turbo")
chain = RetrievalQA.from_chain_type(
    llm=llm,
    chain_type="map_reduce",
    retriever=knn_retriever,
)
```

Making sense of a large number of distributed system calls is what Phoenix is designed to do.

# Phoenix Demo

projects > **default**

| Total Traces | Total Tokens | Latency P50 | Latency P99 | Hallucination | QA Correctness | Relevance | Stream |
|---|---|---|---|---|---|---|---|
| 4 | 1,671 | ⏱ 2.60s | ⏱ 3.31s | ◐ 0.75 | ○ 0.00 | ndcg 0.50  precision 0.38  hit rate 0.50 | |

**Traces**    Spans

🔍 filter condition (e.x. span_kind == 'LLM')                                              ⊞ Columns ▾

| ˃ | kind | name | input | output | evaluations | start time | latency | total tokens | status |
|---|---|---|---|---|---|---|---|---|---|
| ˃ | chain | query | How do I log a prediction using the python SDK? | To log a prediction using the Python SDK, you can use the `arize.log()` function. You need to provid... | Hallucination factual  QA Correctness correct | 12/11/2023, 11:57 AM | ⏱ 2.60s | 👁 653 | ✓ |
| ˃ | chain | query | How much does an enterprise license of Arize cost? | I'm sorry, but I don't have access to pricing information for Arize. For detailed pricing informatio... | Hallucination factual  QA Correctness correct | 12/11/2023, 11:57 AM | ⏱ 3.31s | 👁 293 | ✓ |
| ˃ | chain | query | How do I delete a model? | To delete a model, you would need to access the model management or administration section of the pl... | Hallucination hallucinated  QA Correctness incorrect | 12/11/2023, 11:57 AM | ⏱ 2.82s | 👁 319 | ✓ |
| ˃ | chain | query | How can I query for a monitor's status using GraphQL? | You can query for a monitor's status using GraphQL by including the "status" field in your query. | Hallucination factual  QA Correctness correct | 12/11/2023, 11:57 AM | ⏱ 2.18s | 👁 406 | ✓ |

# Thank you.

Sign up for a free account at:

**arize.com/join**

Check out Phoenix, our OSS tool, at:

**phoenix.arize.com**