# Evaluating Agents and Assistants

Jason Lopatecki | Cofounder & CEO, Arize AI

# 2024 - "The Year of Agents"

MultiOn
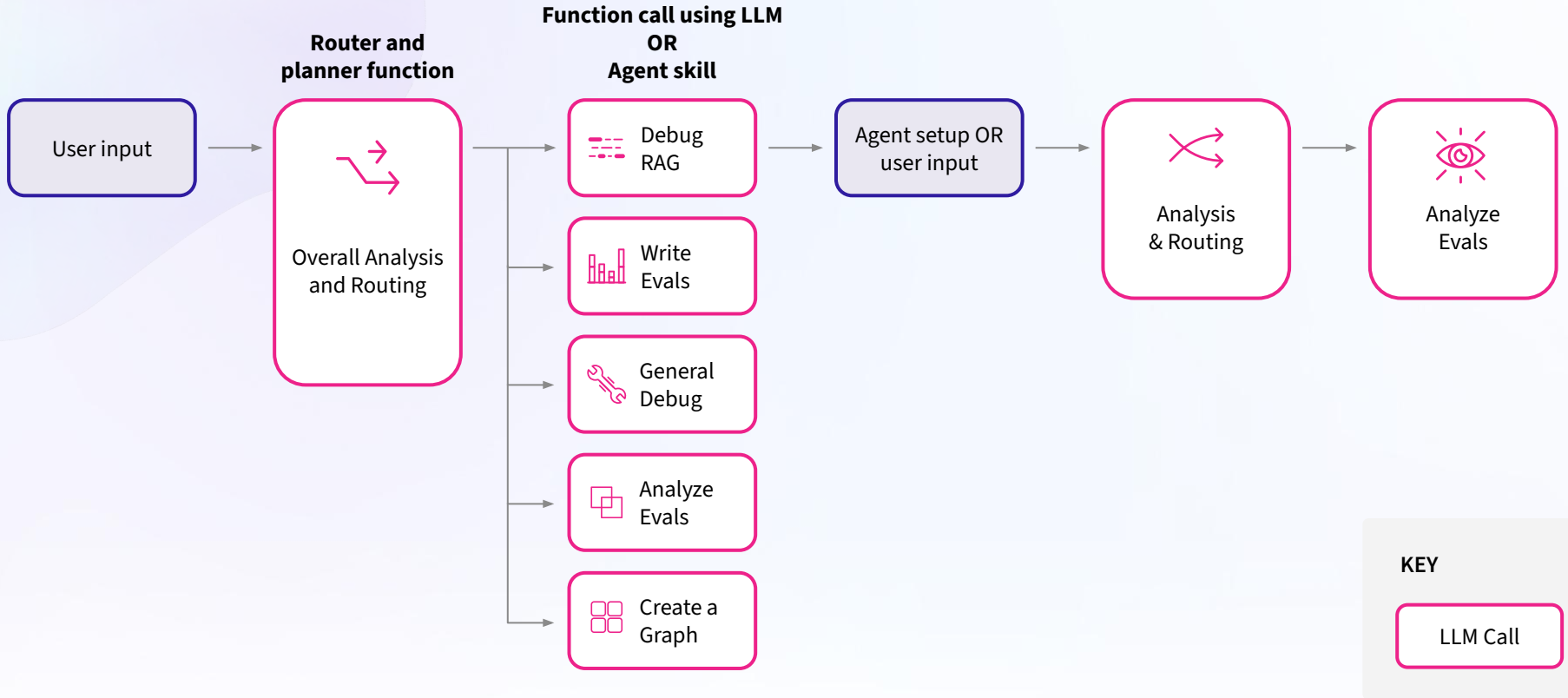
crewai

LlamaIndex

LangChain
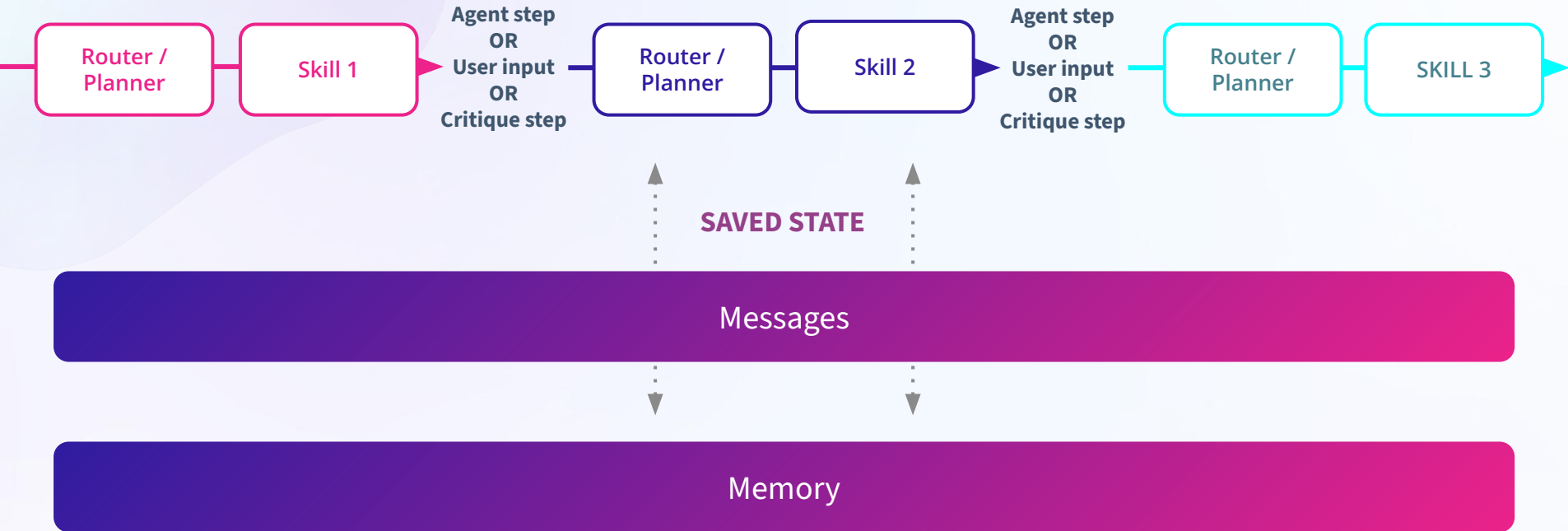
Vertex.ai

OpenAI

ANTHROP\C

MISTRAL AI_

# What We See in the Ecosystem



**Router and planner function**

**Function call using LLM OR Agent skill**

User input → Overall Analysis and Routing →

- Debug RAG
- Write Evals
- General Debug
- Analyze Evals
- Create a Graph

Debug RAG → Agent setup OR user input → Analysis & Routing → Analyze Evals

**KEY**

LLM Call

# Architecture: Agent Process



Router / Planner → Skill 1 → **Agent step OR User input OR Critique step** → Router / Planner → Skill 2 → **Agent step OR User input OR Critique step** → Router / Planner → SKILL 3
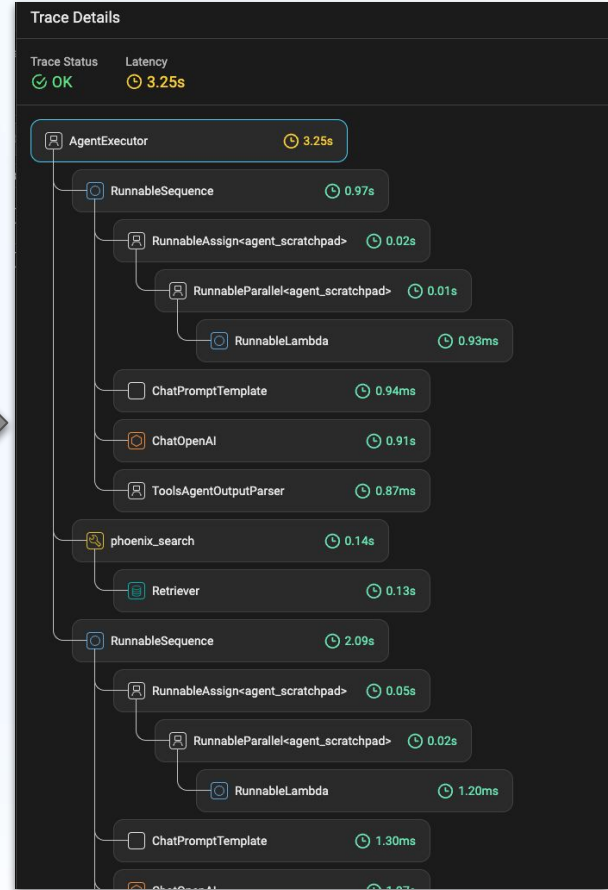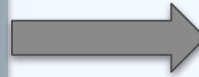
**SAVED STATE**

Messages

Memory

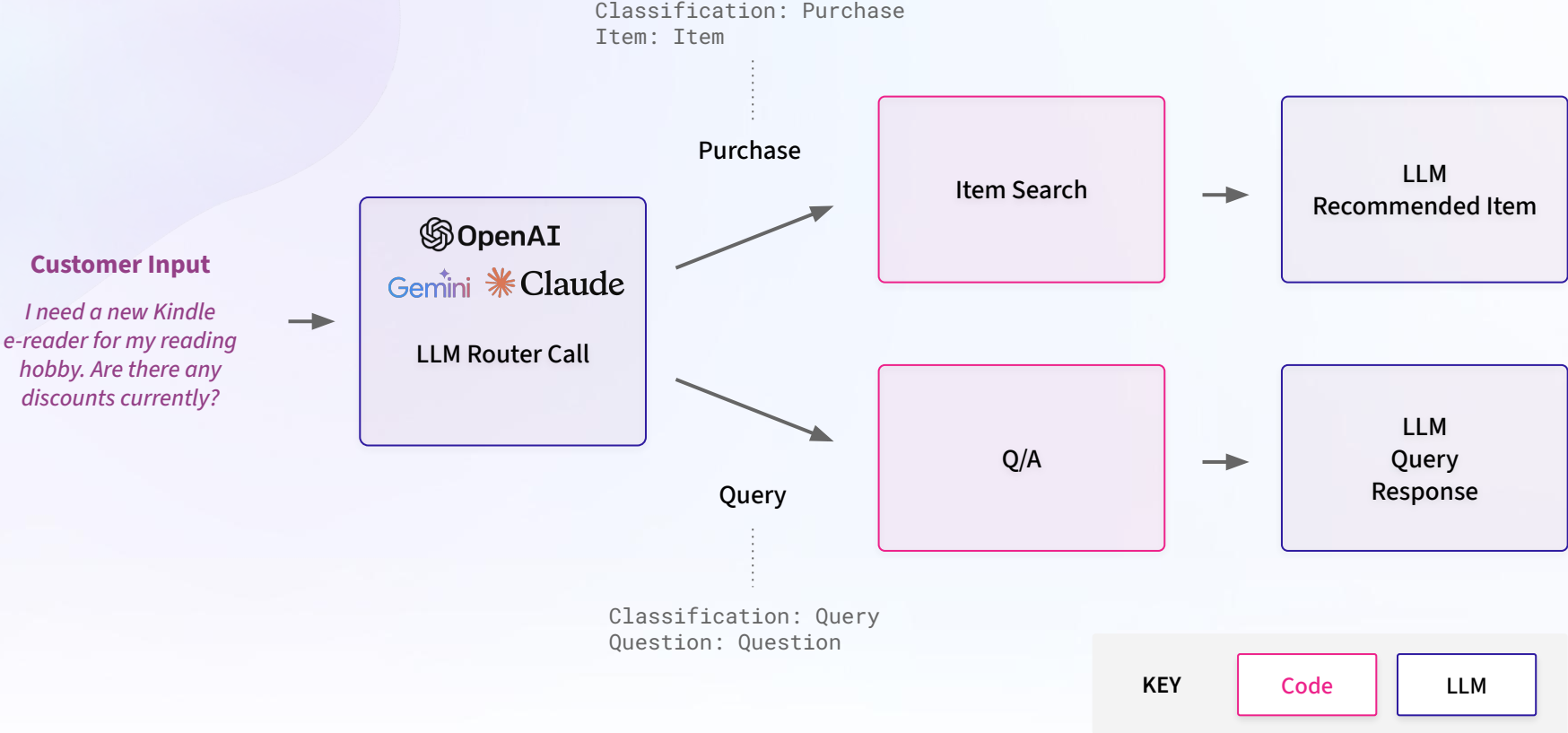# Mixing Branches of LLM Calls with Tools & State
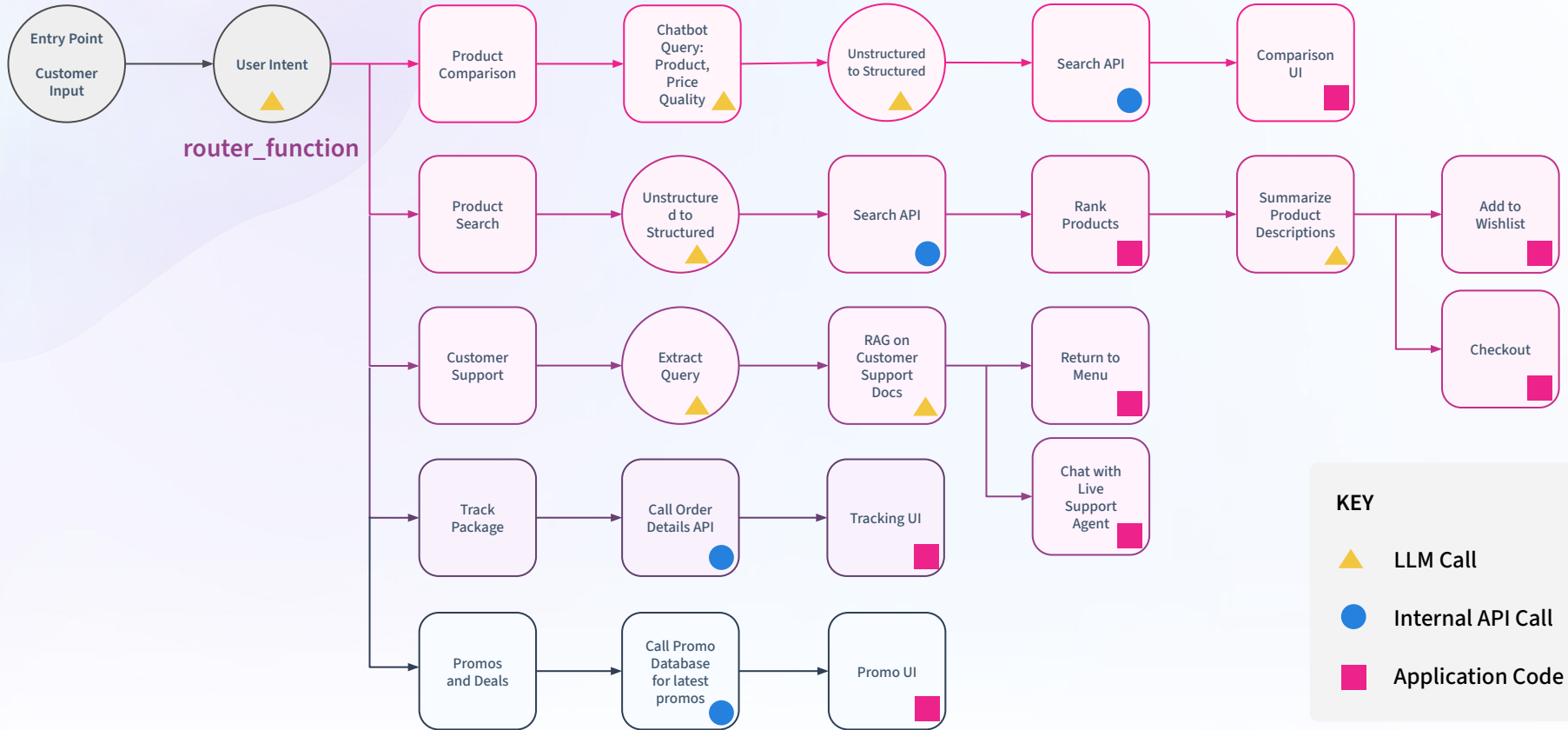
# Difficulties of Evaluating Agents

```python
1   search = TavilySearchResults()
2   loader = WebBaseLoader("https://docs.arize.com/phoenix")
3   docs = loader.load()
4   documents = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200).split_documents(docs)
5   vector = FAISS.from_documents(documents, OpenAIEmbeddings())
6   retriever = vector.as_retriever()
7
8   retriever_tool = create_retriever_tool(
9       retriever,
10      "phoenix_search",
11      "Search for information about Phoenix. For any questions about Phoenix, you must use this tool!",
12  )
13
14  tools = [search, retriever_tool]
15  llm = ChatOpenAI(model="gpt-3.5-turbo-0125", temperature=0)
16  prompt = hub.pull("hwchase17/openai-functions-agent")
17  agent = create_tool_calling_agent(llm, tools, prompt)
18  agent_executor = AgentExecutor(agent=agent, tools=tools, verbose=True)
19
20  agent_executor.invoke({"input": "What kinds of evaluators does Phoenix have?"})
```

**Trace Details**

Trace Status    Latency
⊘ OK            🕑 3.25s

| | |
|---|---|
| 🔲 AgentExecutor | 🕑 3.25s |
| ◎ RunnableSequence | 🕑 0.97s |
| 🔲 RunnableAssign<agent_scratchpad> | 🕑 0.02s |
| 🔲 RunnableParallel<agent_scratchpad> | 🕑 0.01s |
| ◎ RunnableLambda | 🕑 0.93ms |
| 🔲 ChatPromptTemplate | 🕑 0.94ms |
| ◎ ChatOpenAI | 🕑 0.91s |
| 🔲 ToolsAgentOutputParser | 🕑 0.87ms |
| 🔍 phoenix_search | 🕑 0.14s |
| ▤ Retriever | 🕑 0.13s |
| ◎ RunnableSequence | 🕑 2.09s |
| 🔲 RunnableAssign<agent_scratchpad> | 🕑 0.05s |
| 🔲 RunnableParallel<agent_scratchpad> | 🕑 0.02s |
| ◎ RunnableLambda | 🕑 1.20ms |
| 🔲 ChatPromptTemplate | 🕑 1.30ms |

# Example Router: Chat to Purchase App



Classification: Purchase
Item: Item

Purchase

Customer Input

*I need a new Kindle e-reader for my reading hobby. Are there any discounts currently?*

OpenAI
Gemini ❋ Claude

LLM Router Call

Item Search

LLM Recommended Item

Q/A

LLM Query Response

Query

Classification: Query
Question: Question

KEY    Code    LLM

# Under the Hood: Chat-to-purchase Router



router_function

**KEY**

▲ LLM Call

● Internal API Call

■ Application Code

Entry Point — Customer Input → User Intent ▲

Product Comparison → Chatbot Query: Product, Price Quality ▲ → Unstructured to Structured ▲ → Search API ● → Comparison UI ■

Product Search → Unstructured to Structured ▲ → Search API ● → Rank Products ■ → Summarize Product Descriptions ▲ → Add to Wishlist ■ / Checkout ■

Customer Support → Extract Query ▲ → RAG on Customer Support Docs ▲ → Return to Menu ■ / Chat with Live Support Agent ■

Track Package → Call Order Details API ● → Tracking UI ■

Promos and Deals → Call Promo Database for latest promos ● → Promo UI ■

# Agent Frameworks



**LlamaIndex Workflows**

- RetryQueryEvent
- StartEvent → Answer Query
- QueryFailedEvent → Improve Query
- QuerySucceededEvent → done

**LangGraph**

- input
- classify_input
- is_greeting
  - true → handle_greeting
  - false → handle_RAG
- check_rag_response
  - fail → handle_RAG
  - pass → _end_
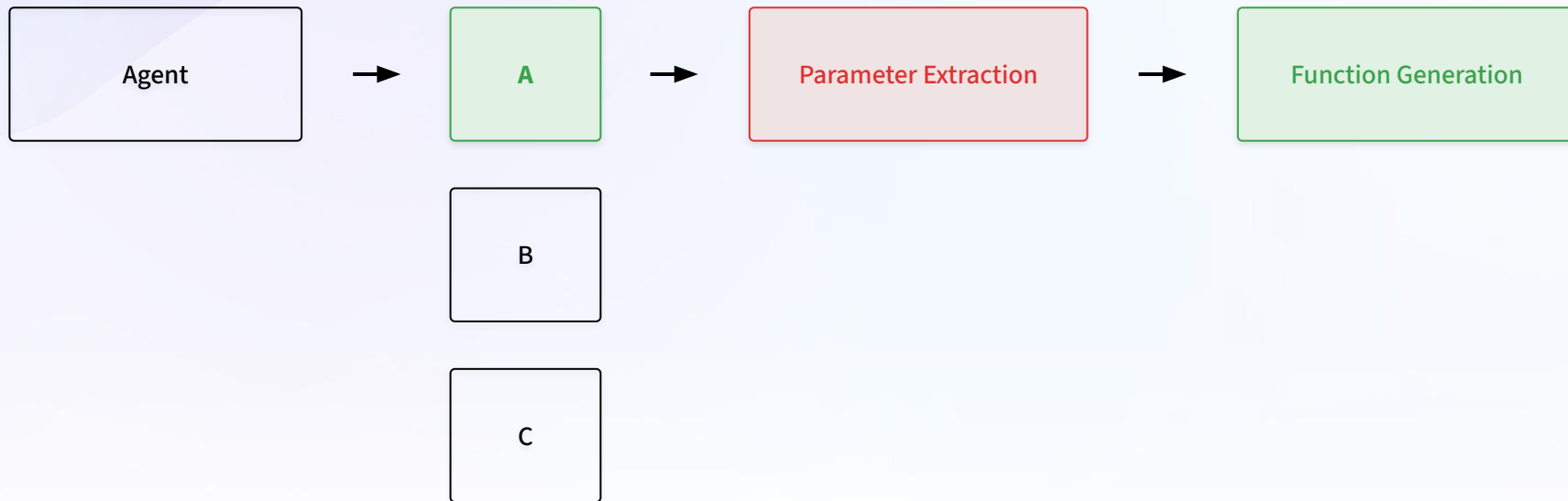
# "Working" ≠ Performing Well

# So how do you evaluate these agents?

0.  (Add observability to your application!)

1.  Build a set of test cases

2.  Breakdown individual agent steps

3.  Create evaluators for each step

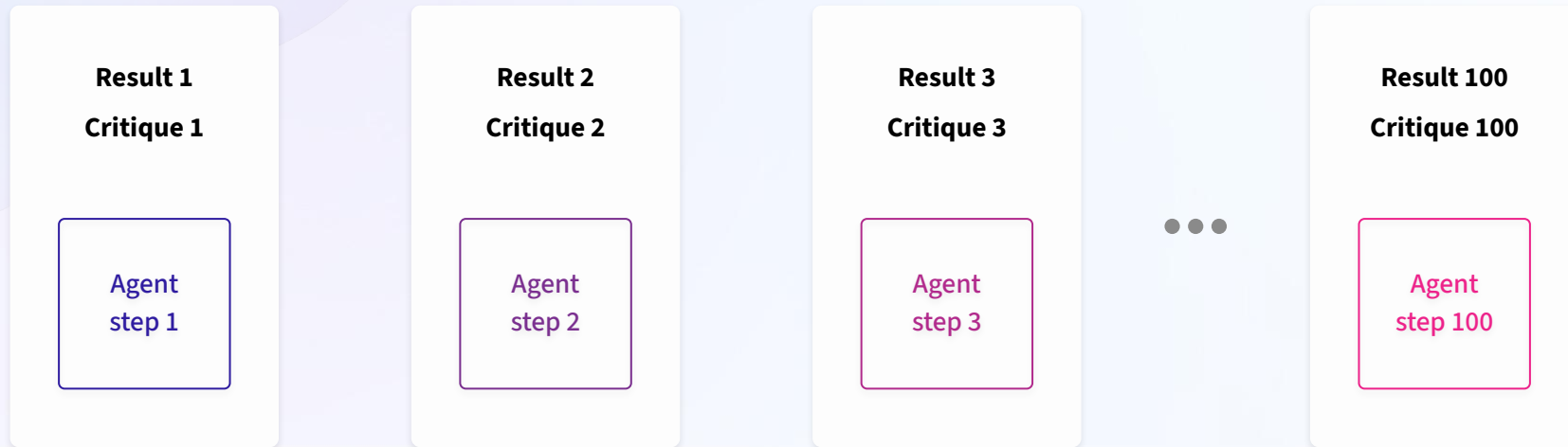4.  Experiment and iterate

# Breakdown of "Steps"

Am I using the right skill correctly?

**Skills**

# Am I Converging?

## LLM as a judge eval: Agent convergence

**Result 1**

**Critique 1**

Agent
step 1

**Result 2**

**Critique 2**

Agent
step 2

**Result 3**

**Critique 3**

Agent
step 3

● ● ●

**Result 100**

**Critique 100**

Agent
step 100

# Teaser: Convergence Evals

```
You are an AI assitant evaluting AI Agent execution. Your overall goal is
  to asses Agent convergence. AI Agents can sometimes get stuck in loops
  and not accomplish or coverge on the overall task.

  Is it moving forward in its goals, is it stuck repeating the same tasks, is
  it likely to converge to an outcome.

    [BEGIN DATA]
    ************
    [Agent Previous Step Results]: {agent_results}
    ************
    [Agent Previous Step Critques]: {agent_step_critques}
    ************
    [END DATA]
```

| Total Traces | Total Tokens | Latency P50 | Latency P99 | Function Calls | Parameter Extraction JSON |
|---|---|---|---|---|---|
| 49 | 33,149 | ⏱ 1.22s | ⏱ 5.42s | 🟢 0.00 | 🟢 0.33 |

Stream 🔵   📅 All Time

**Traces**   Spans

🔍 filter condition (e.x. span_kind == 'LLM')   ⊕   | ▦ Columns ▾

| ☐ | > | kind | name | input | output | evaluations | start time | latency |
|---|---|---|---|---|---|---|---|---|
| ☐ | > | agent | AgentExecutor | Can you help me decide whether to upgrade my current product or wait for a new one to arrive? | Of course! To help you make an informed decision, I can compare the features of your current prod... | – | 6/26/2024, 10:44 PM | ⏱ 1.23s |
| ☐ | > | agent | AgentExecutor | I'm not sure if I should request a refund or just track the delivery status of my order. What do ... | To provide the best recommendation, I need a bit more information: 1. **Current Status of Your O... | – | 6/26/2024, 10:44 PM | ⏱ 2.12s |
| ☐ | > | agent | AgentExecutor | Could you tell me if there are any current promotions for the Samsung 106i smartphone? | It looks like there are no current promotions or listings for the Samsung 106i smartphone at the ... | Function Calls incorrect  Parameter Extraction JS... incorrect | 6/26/2024, 10:44 PM | ⏱ 1.46s |
| ☐ | > | agent | AgentExecutor | I need help with an issue related to an iPhone 16H I bought, but I'm not sure when I purchased it. | I can assist you with that. Since you are unsure about the purchase date, it would be helpful to ... | – | 6/26/2024, 10:44 PM | ⏱ 1.27s |
| ☐ | > | agent | AgentExecutor | Is there a way to check if the Samsung 15H adapter in my last order is still under warranty? | I can't directly check the warranty status of a specific product from your order history. However... | – | 6/26/2024, 10:44 PM | ⏱ 1.30s |
| ☐ | > | agent | AgentExecutor | I recently purchased a product, but I'm not sure if it was the Vizio 14Y TV or another model. Can... | Sure, I can help you track your package. Could you please provide the tracking number for the pac... | – | 6/26/2024, 10:44 PM | ⏱ 0.73s |
| ☐ | > | agent | AgentExecutor | Perhaps you could assist me in understanding the difference between | I can help with that! To get started, I'll need the unique identifiers (product IDs) for the | – | 6/26/2024, 10:44 PM | ⏱ 1.19s |

# Q/A

---

✨ [Star us on Github](github.com/Arize-ai/phoenix) ✨

github.com/Arize-ai/phoenix