

From RAG to Knowledge Assistants

Jerry Liu
Llamalndex



Llamaindex

Llamaindex is the fastest way to build production-quality **LLM applications over enterprise data**.

Products:

- **LlamaCloud**: Enterprise RAG platform
- **Open-Source**: Agent orchestration framework

Backed By: Greylock



Leading developer platform and community for GenAI:

- 200k+ followers on LinkedIn
- 70k+ followers on Twitter
- 34k+ stars
- 2M+ monthly downloads
- 20k+ Discord members
- 600+ integrations

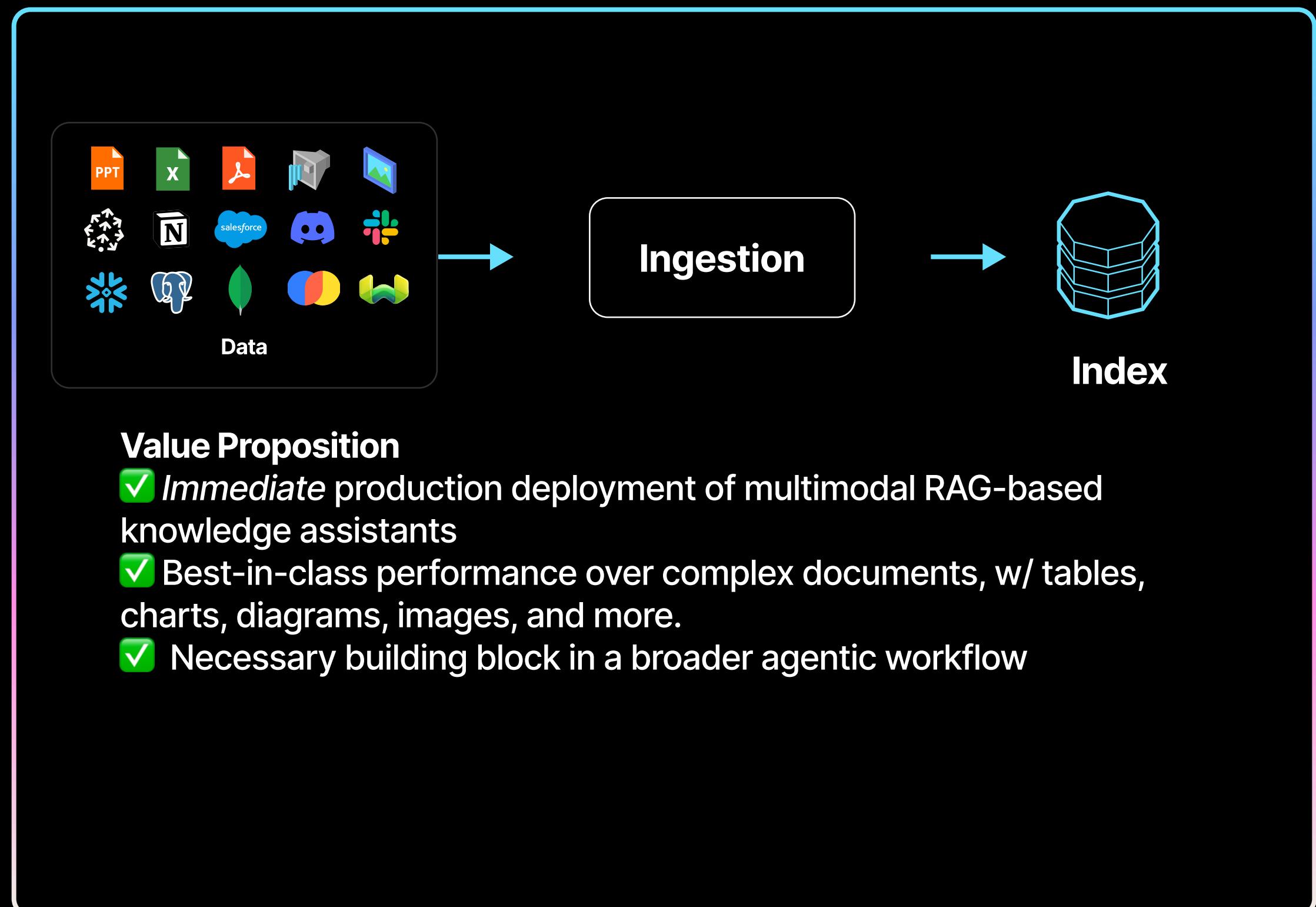
Fortune 500 Companies to Startups



LLM Agents over any Data Type

LlamaCloud

Enterprise RAG platform to rapidly connect unstructured enterprise data to LLM agents.



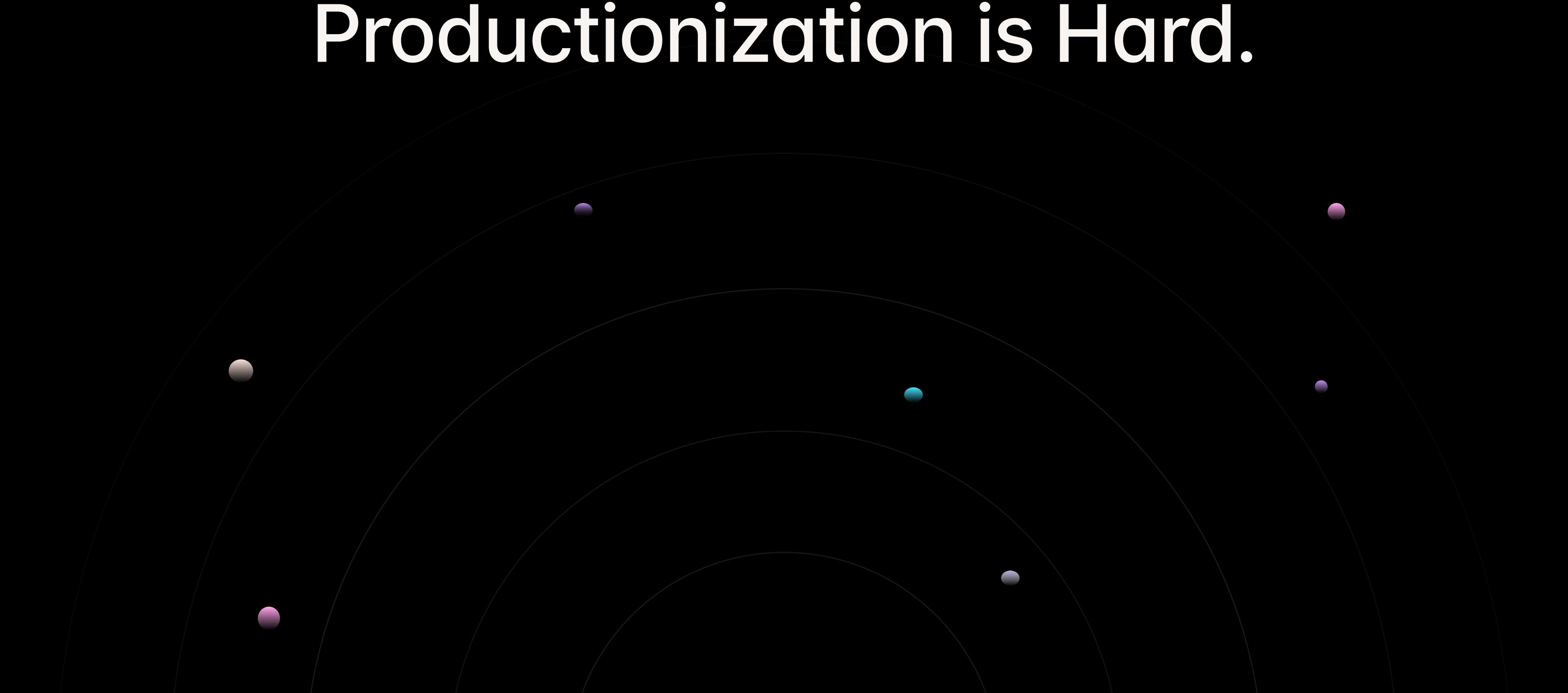
Open-Source

Leading developer framework for orchestrating **single and multi-agent workflows** over your data.

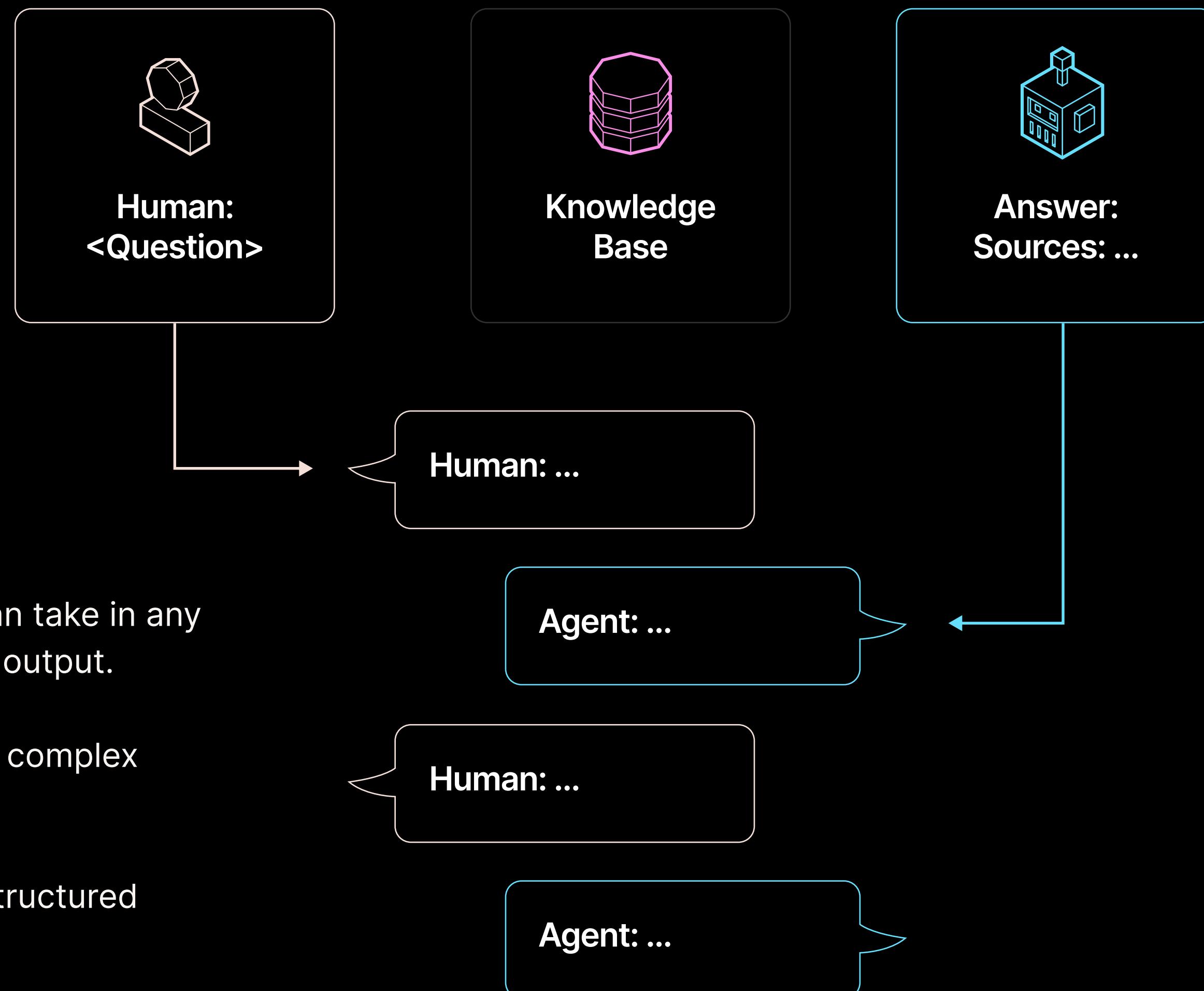


- Knowledge Assistants
- Report Generation
- Business Process Automation
- and more!

Prototyping is Easy.
Productionization is Hard.



Building a Knowledge Assistant



Goal: Build an interface that can take in any task as input and give back an output.

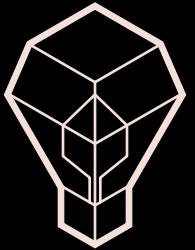
Input forms: simple questions, complex questions, research tasks

Output forms: short answer, structured output, research report

Can we do more?

There's many questions/tasks that naive RAG can't give an answer to

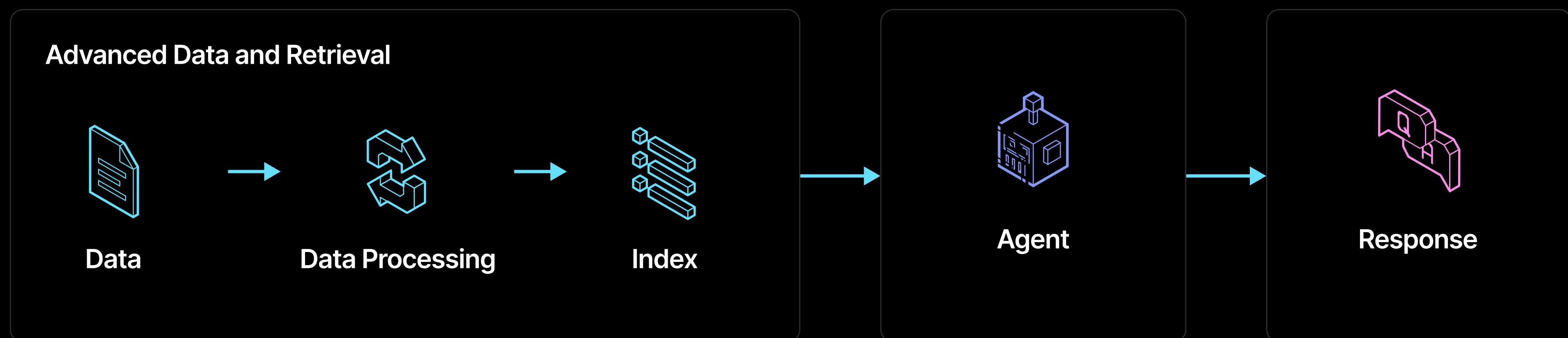
- 🚫 Hallucinations
- 🚫 Limited time savings
- 🚫 Limited decision-making enhancement



How do we aim to build
a production-ready
knowledge assistant?

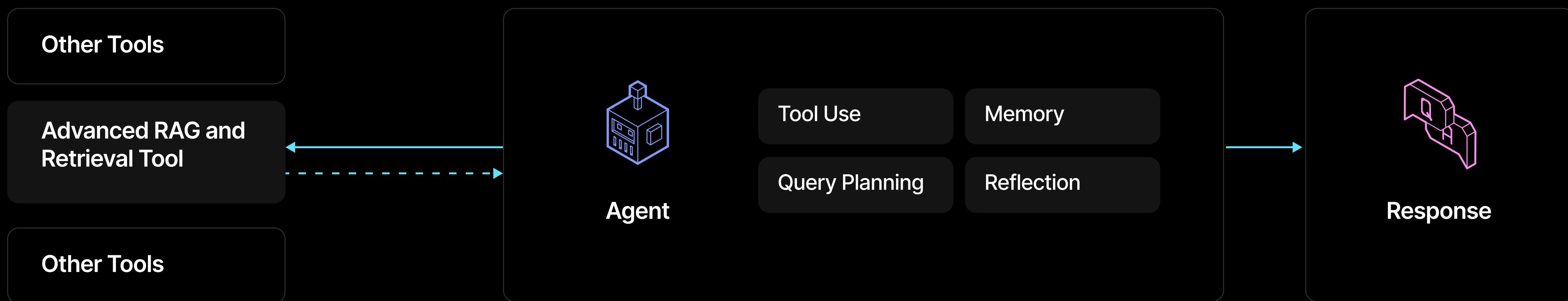
A Better Knowledge Assistant

1. High-quality data and retrieval
2. Agentic reasoning over complex inputs
3. Agentic decision-making and output generation
4. Towards a scalable, full-stack application



A Better Knowledge Assistant

1. High-quality data and retrieval
2. **Agentic reasoning over complex inputs**
3. Agentic decision-making and output generation
4. Towards a scalable, full-stack application



A Better Knowledge Assistant

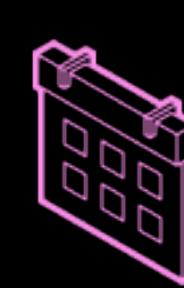
1. High-quality data and retrieval
2. Agentic reasoning over complex inputs
3. **Agentic decision-making and output generation**
4. Towards a scalable, full-stack application



Report
Generation



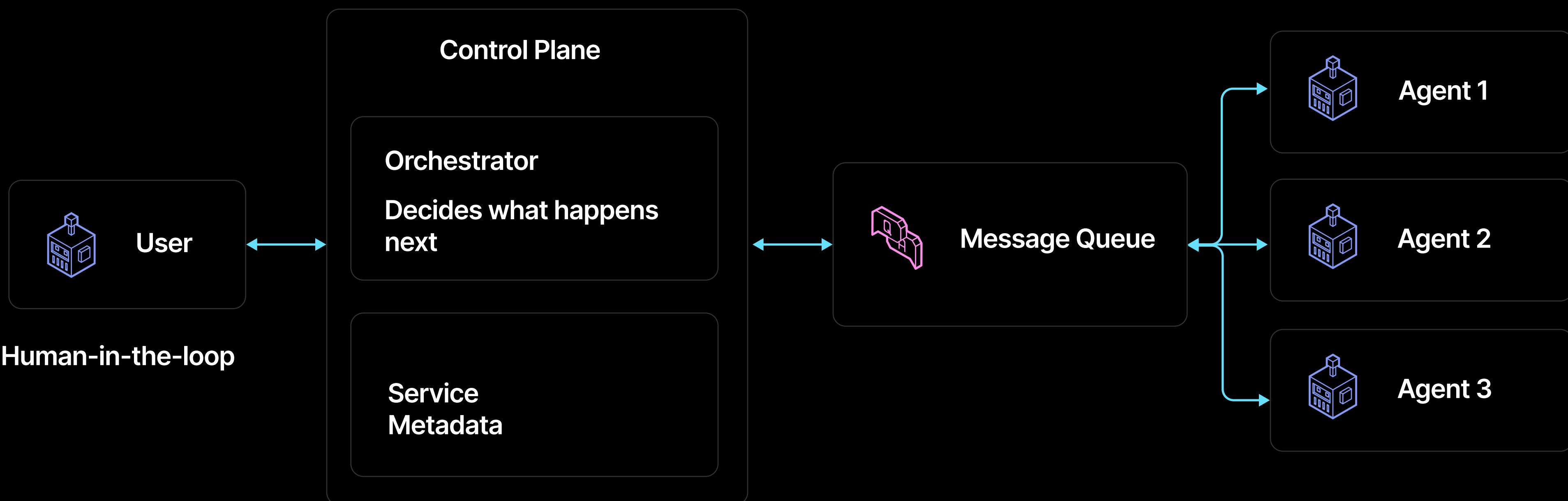
Data Analysis



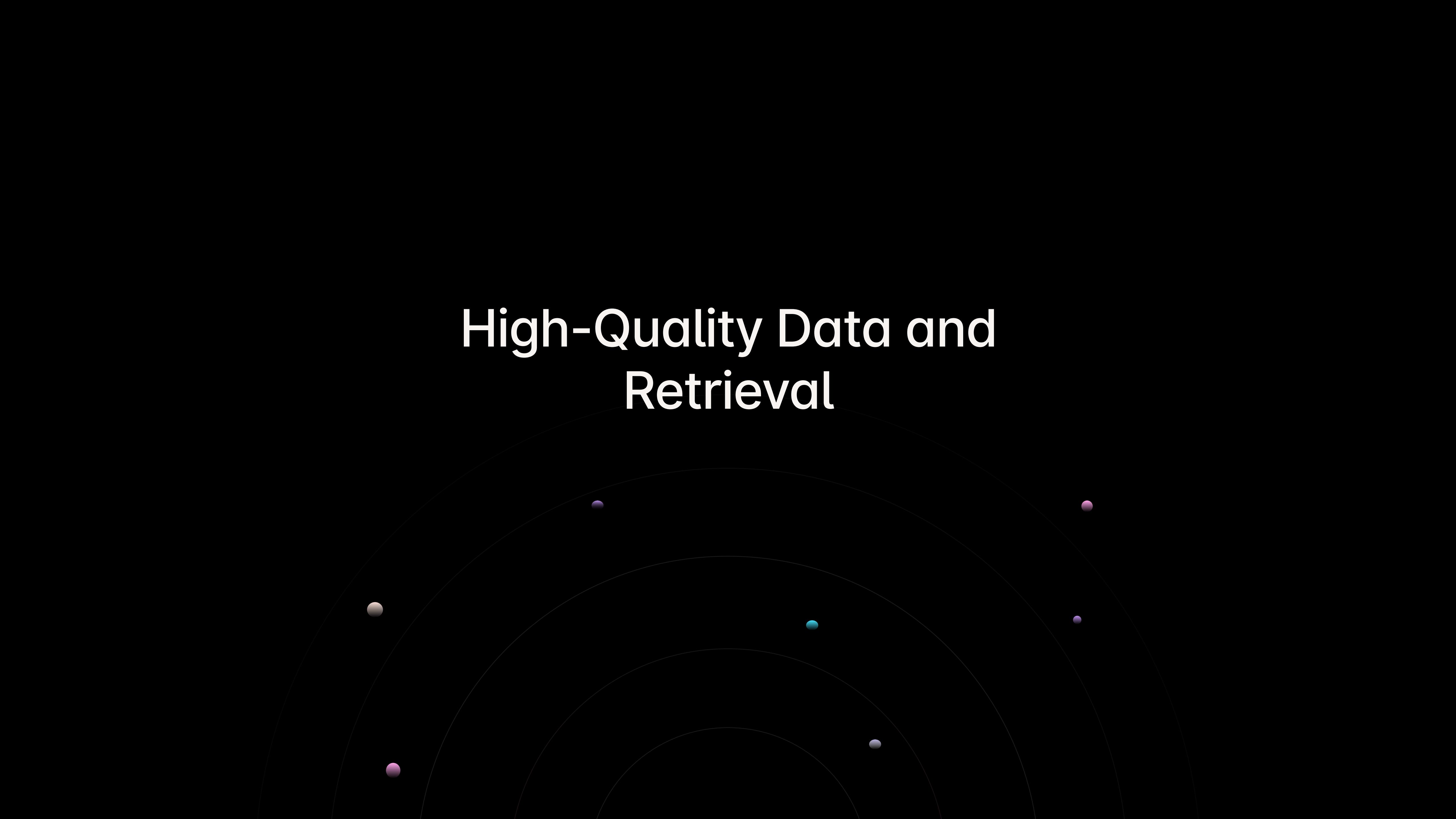
Action-Taking

A Better Knowledge Assistant

1. High-quality data and retrieval
2. Agentic reasoning over complex inputs
3. Agentic decision-making and output generation
4. **Towards a scalable, full-stack application**



High-Quality Data and Retrieval

The background features a dark gray gradient with several thin, light gray concentric circles. Scattered throughout are small, semi-transparent colored dots in shades of purple, pink, teal, and light gray.

Case Study: Complex Documents

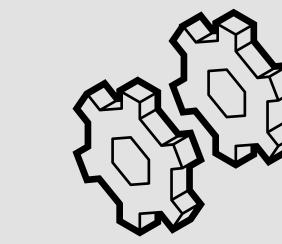
A lot of documents can be classified as **complex**:

- Embedded Tables, Charts, Images
- Irregular Layouts
- Headers/Footers

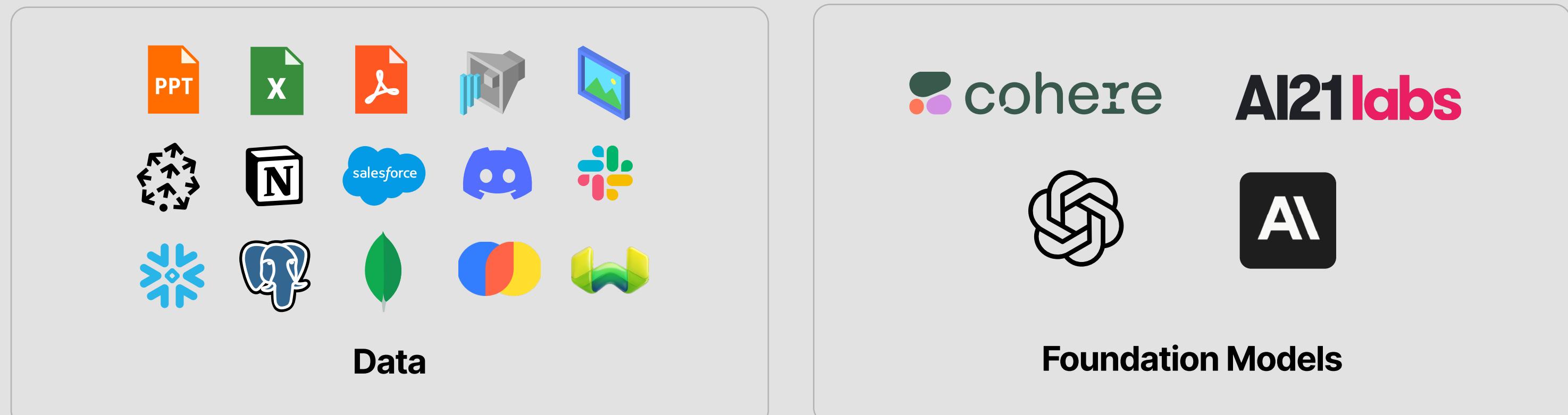
Users want to ask research questions over this data:

- Simple pointed questions
- Multi-document comparisons
- Research tasks

Building a production-ready knowledge assistants over this complex data is **challenging**.



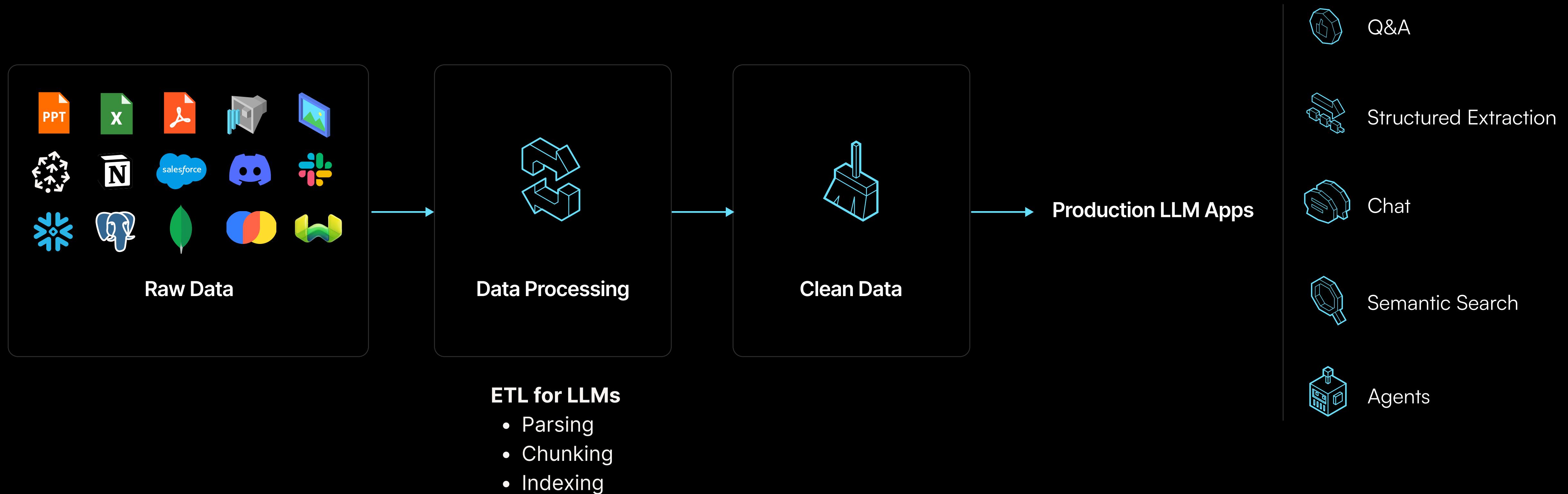
Developers



Any LLM App is only as Good as your Data

Garbage in = garbage out

Good data quality is a **necessary** component of any production LLM app.



LlamaParse

"As an AI Applied Data Scientist who was granted one of the first ML patents in the U.S., and who is building cutting-edge AI capabilities at one of the world's largest Private Equity Funds, I can confidently say that LlamaParse from LlamalIndex is currently the best technology I have seen for parsing complex document structures for Enterprise RAG pipelines. Its ability to preserve nested tables, extract challenging spatial layouts, and images is key to maintaining data integrity in advanced RAG and agentic model building."

Dean Barr, Applied AI Lead at Carlyle



LlamaParse

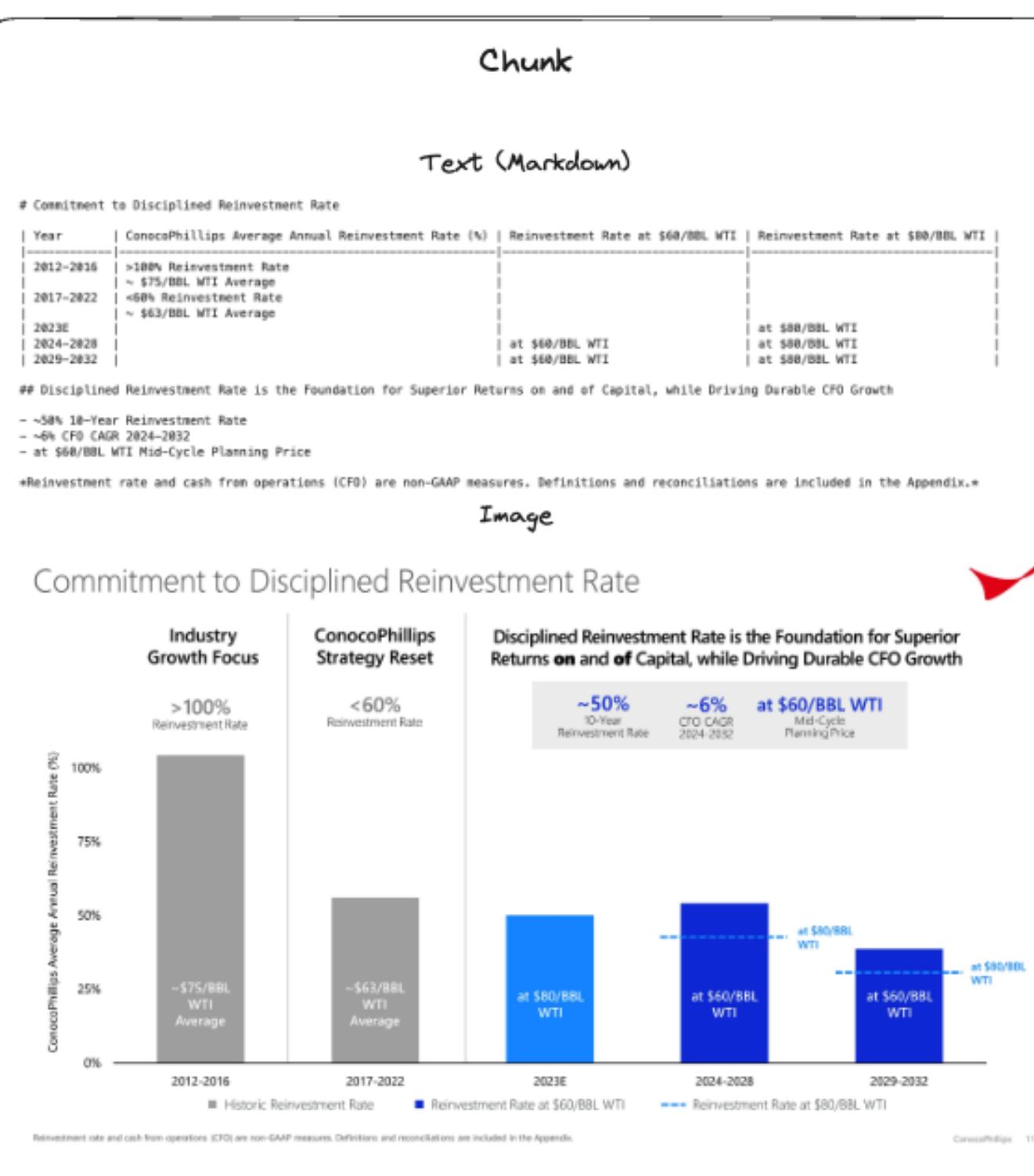
Advanced document parser specifically for reducing LLM hallucinations

20k+
unique users

25M+
pages processed

Use Cases

Earnings Decks (PPTX)



Vector Database

10K Reports (Text + Tables)

Cash flows from financing activities for Net

Query Engine*****")

Retriever Query Engine*****")

Retrieving

financing activities for Netflix is not provided in the 10-K report.

flows from financing activities were \$700,000 million on December 31, 2021.

Query Engine*****")

Retrieving

financing activities for the year ended December 31, 2021.

ds

Net cash provided by (used in) financing activities

Financials (Excel)

NVIDIA QUARTERLY REVENUE TREND REVENUE BY MARKET				
4	Q3 FY24	Q2 FY24	Q1 FY24	Q4 FY23
\$18,404	\$14,514	\$10,323	\$4,284	
2865	2856	2486	2240	
463	416	379	295	
281	261	253	296	
90	73	66	77	
\$22,103	\$18,120	\$13,507	\$7,192	
FY24	Q2 FY24	Q1 FY24	Q4 FY23	Q3 FY23
\$14,514	\$10,323	\$4,284	\$3,616	
2,856	2,486	2,240	1,831	
416	379	295	226	
261	253	296	294	
73	66	77	84	
\$18,120	\$13,507	\$7,192	\$6,051	
FY24	Q1 FY24	Q4 FY23	Q3 FY23	Q2 FY23
\$10,323	\$4,284	\$3,616	\$3,833	
2,486	2,240	1,831	1,574	
379	295	226	200	
253	296	294	251	
66	77	84	73	
\$13,507	\$7,192	\$6,051	\$5,931	
FY24	Q4 FY23	Q3 FY23	Q2 FY23	Q1 FY23
(407,729)	(524,585)			
—	(26,919)	\$4,284	\$3,616	\$3,833
(757,387)	(788,349)	\$2,240	\$1,831	\$1,574
(911,276)	—	\$295	\$226	\$200
(2,076,392)	(1,339,853)	\$296	\$294	\$251
		\$77	\$84	\$73
		\$7,192	\$6,051	\$5,931
FY23	Q3 FY23	Q2 FY23	Q1 FY23	Q4 FY23
(700,000)	(500,000)	\$3,616	\$3,833	\$3,806
35,746	174,414	\$1,831	\$1,574	\$2,042
—	(600,022)	\$226	\$200	\$496
		\$294	\$251	\$220
		\$77	\$84	\$73
		\$7,192	\$6,051	\$6,704
FY23	Q2 FY23	Q1 FY23	Q4 FY23	Q3 FY23
(664,254)	(1,149,776)	\$6,051	\$5,931	\$6,704
		\$84	\$73	\$140
		\$140	\$158	\$158
ds	ds	ds	ds	ds

Accident Claims (Forms)

AUTOMOBILE CLAIM

LOSS

Date 06/23/2023
 Location Intersection of Vine Street and Sunset Bl
 City Los Angeles State CA
 Police Dept. Involved LAPD Ticket Issued Traffic Violation

DESCRIPTION OF ACCIDENT

On October 15, 2023, at approximately 3:30 PM, I was driving my 2020 Honda Accord (License Plate: 7XYZ123) southbound on Vi approaching the intersection with Sunset Blvd. As I entered the intersection, a blue 2018 Ford Escape (License Plate: 8ABC456) traveling eastbound on Sunset Blvd ran a red light and collided with the front passenger side of my

INSURED VEHICLE

Year 2020 Make Honda Model Accord
 V.I.N. 1HGCV1F30LA123456 Plate 7XYZ123
 Extent of Damages The front passenger side of my Honda Accord sustained significant damage, including a dented fender and broken headlight. Estimated repair cost: \$3,500.
 Present Location Impound Lot
 Driver Michael Johnson (ASK IF OFFICE
 Date of Birth 01/15/1985 License No. 1111111111 State CA

OTHER VEHICLE

Year 2018 Make Ford Model Escape
 Extent of Damages The Ford Escape had damage to the front bumper and hood. Estimated repair cost: \$2,000.
 Owner Sarah Brown Phone 2139876543
 Address 405 Hilgard Av
 City Los Angeles State CA Zip 90095
 Address _____ State _____ Zip _____
 City _____ State _____ Zip _____

Insurance Information

Company Name Mors Mutual Insurance Policy No. 987654321
 Agent Name Emily Carter Phone 2131234567

INJURED

Name Michael Johnson Phone 3101234567
 Address 3470 Troutdale Pkwy
 City Los Angeles State CA Zip 90089
 Extent of Injury I sustained minor injuries, including neck pain and a bruise on my left arm. I sought medical attention at the local urgent care

WITNESSES

Name David Thompson Phone 3105678901
 Address 633 W 5th St
 City Los Angeles State CA Zip 90071

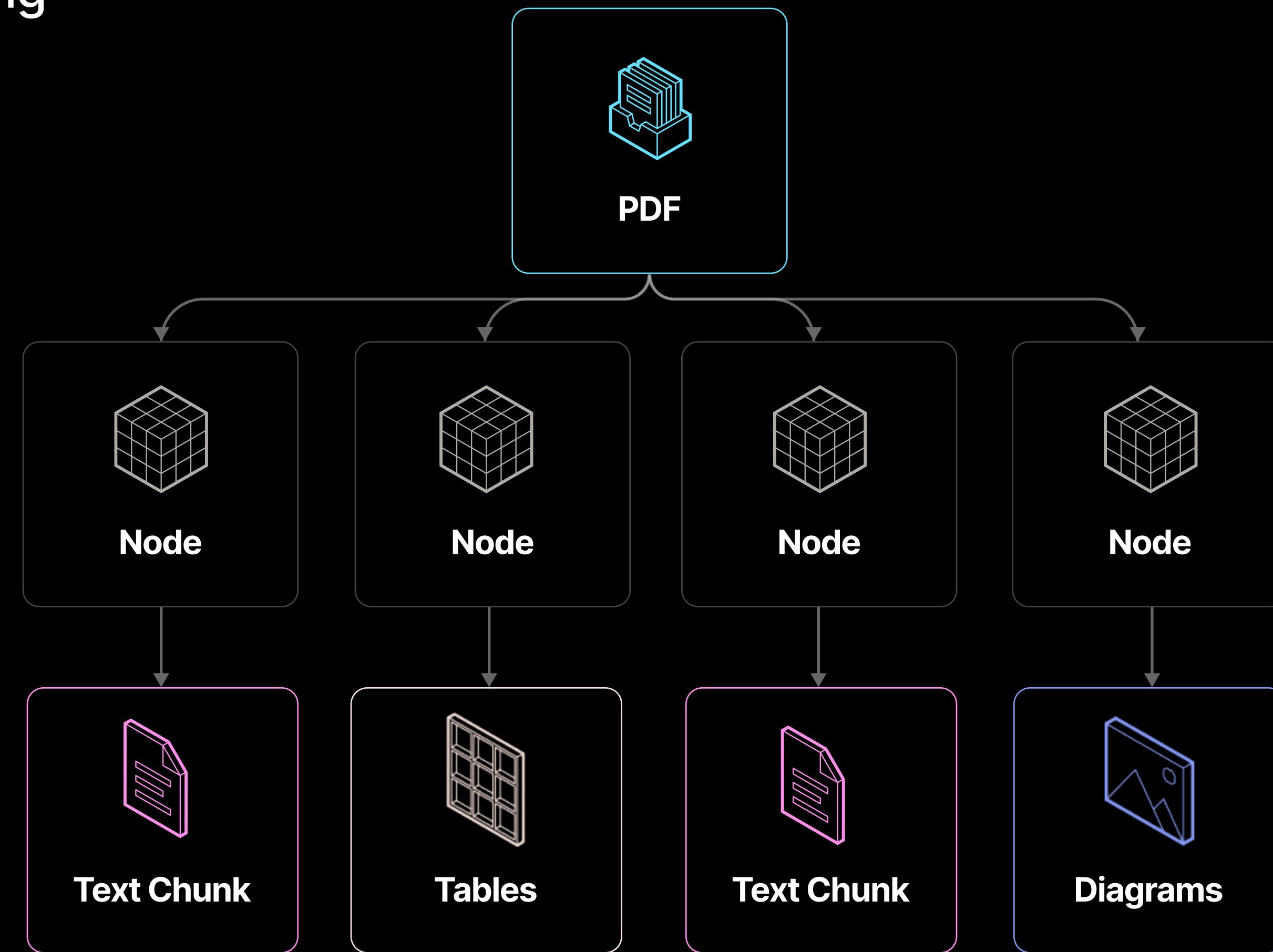
IMPACT

Is damaged auto essential to business? No
 How?

Advanced Parsing + Advanced Indexing

You can combine parsing with **hierarchical indexing and retrieval** to model heterogeneous unstructured/tabular/multimodal data within a document.

1. Parse documents into elements: text chunks, tables, images, and more.
2. For each element, extract **one or more** text representations that can be indexed.
3. Do **recursive retrieval**



LlamaCloud: An Enterprise RAG Platform

A production-ready RAG platform that allows developers to easily connect their unstructured data sources to LLM agent systems.

Instant Time-to-Value for building knowledge assistants

- Out-of-the-box advanced RAG capabilities
- Free up developer time to rapidly iterate on higher-level agent use cases

State-of-the-Performance leads to increased

satisfaction and reduced compliance risk

Reduced maintenance cost once application is deployed

Enterprise-ready security like access controls

Signup: <https://cloud.llamaindex.ai/>

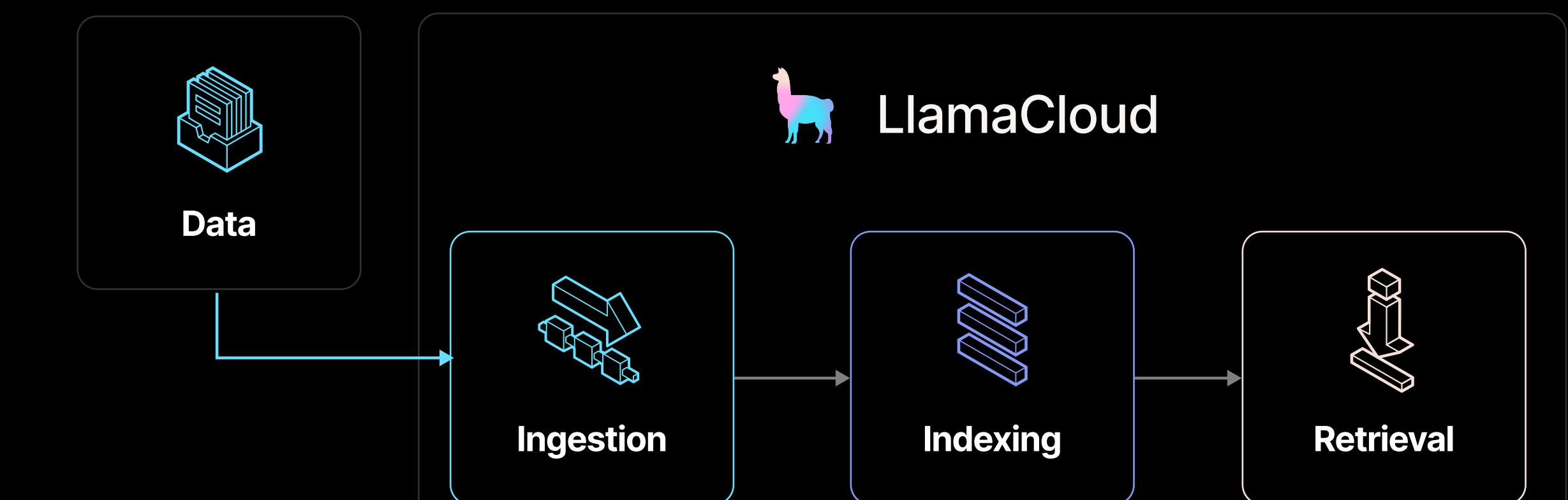
The screenshot shows the LlamaCloud web interface. At the top right, there's a user profile icon and links for 'Rerun' and 'Delete'. The main area displays a project titled 'uber_and_lyft' under 'test_projects_03_28_2024'. On the left sidebar, there are sections for 'PROJECT' (with 'test_projects_03_28_2024' selected), 'PROJECT NAVIGATION' (with 'Parse', 'Index', and 'Evals'), 'RESOURCES' (with 'API Key'), and 'YOUR PLAN' (showing 'Free Plan' and an 'Upgrade Plan' button). The central content area has tabs for 'Overview', 'Live Preview', 'Data Sources', and 'Activity'. Under 'Retrieval Endpoint', it shows an 'Endpoint URL' (a long API endpoint) with a 'Copy' button, and an 'Endpoint API' section with a code snippet for interacting with the index using Python:

```
from llama_index.indices.managed.llama_cloud import LlamaCloudIndex
# pip install llama-index-indices-llama-cloud

index = LlamaCloudIndex(
    name="uber_and_lyft",
    project_name="test_projects_03_28_2024",
    api_key="llx-..."
)

nodes = index.as_retriever().retrieve(query)
response = index.as_query_engine().query(query)
```

On the right side, there's an 'Index Information' panel with details: Status (Success), Index ID (edf2f149-b8b9-4b5b-9806-b5d4a182fe86), File Count (0), Date Updated (March 28, 2024, 10:43 AM), and Date Created (March 28, 2024, 1:54 PM).



E2E Multimodal RAG Capabilities

Setup multimodal indexing and retrieval in minutes

Signup here: <https://cloud.llamaindex.ai/>

Node 1 - Page Screenshot

Our Differentiated Portfolio: Deep, Durable and Diverse

~20 BBOE of Resource
Under \$40/BBL Cost of Supply

~\$32/BBL
Average Cost of Supply

WTI Cost of Supply (\$/BBL)

Resource (BBOE)

Legend: Lower 48, Canada, Alaska, EMENA, Asia Pacific

Costs assume a mid-cycle price environment of \$40/BBL WTI.

Diverse Production Base
10-Year Plan Cumulative Production (BBOE)

Lower 48

Alaska

GKA

GWA

GPA

WNS

EMENA

Norway

Qatar

Libya

Asia Pacific

Canada

Montney

APLNG

Malaysia

China

Surmont

Eagle Ford

Bakken

Other

Permian

ConocoPhillips 14

File: 2023-conocophillips-aim-presentation (1).pdf

File ID: d199c200-8ccb-4fa7-aab1-18c185afa389

Page: 14

Projects > llamacloud_demo > multimodal_test > playground

multimodal_test

Sync Edit Copy Delete

Overview Playground Data Sources

Chat Retrieval

Chat Options

tell me about the diverse geographies which represent the production bases

Show events >

ConocoPhillips has a diverse and extensive production base spread across various geographic regions. Here's a detailed breakdown of the regions where ConocoPhillips operates:

Lower 48 States (USA) - This includes significant operations in the Permian Basin, which is further divided into the Delaware Basin and the Midland Basin. The Delaware Basin, for instance, has shown improved performance from 2019 to 2022, and the company has a vast inventory in this area.

Alaska (USA) - ConocoPhillips has a strong presence in Alaska, where it engages in conventional oil production. The company's operations in this region are expected to leverage existing infrastructure for future growth.

Canada - While specific production figures are not detailed in the documents, Canada is listed as one of the regions in ConocoPhillips' diverse production base.

Norway - In Norway, ConocoPhillips has a production of 115 MBOED, with developments like the four subsea tie backs in the Greater Ekofisk Area, which are on track for onstream in 2024. The license for this area has been extended through 2048.

Libya - The company has increased its working interest to approximately 20% in the Waha Concession, with a production of 50 MBOED.

Malaysia - ConocoPhillips reports a production of 40 MBOED

Sources: 1 2 3

Page Screenshots: 1 2 3 4 5 6

Regenerate

Type a message

Send message

LlamaCloud Dramatically Accelerates Time-to-Value with High Performance

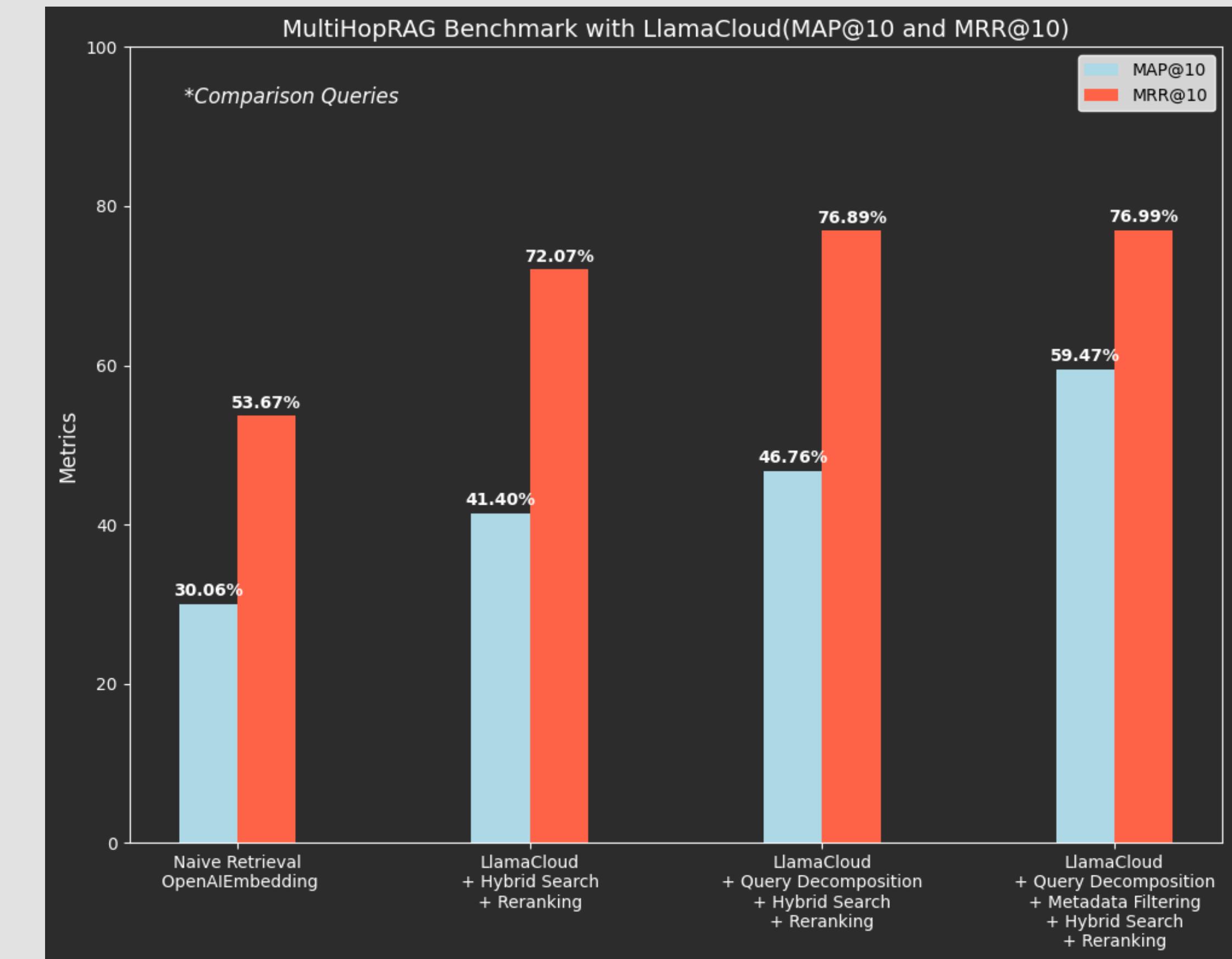
LlamaCloud helps users dramatically accelerate their RAG **time to production for multiple use cases**.

Vs. Pure DIY (open-source, roll-your-own)

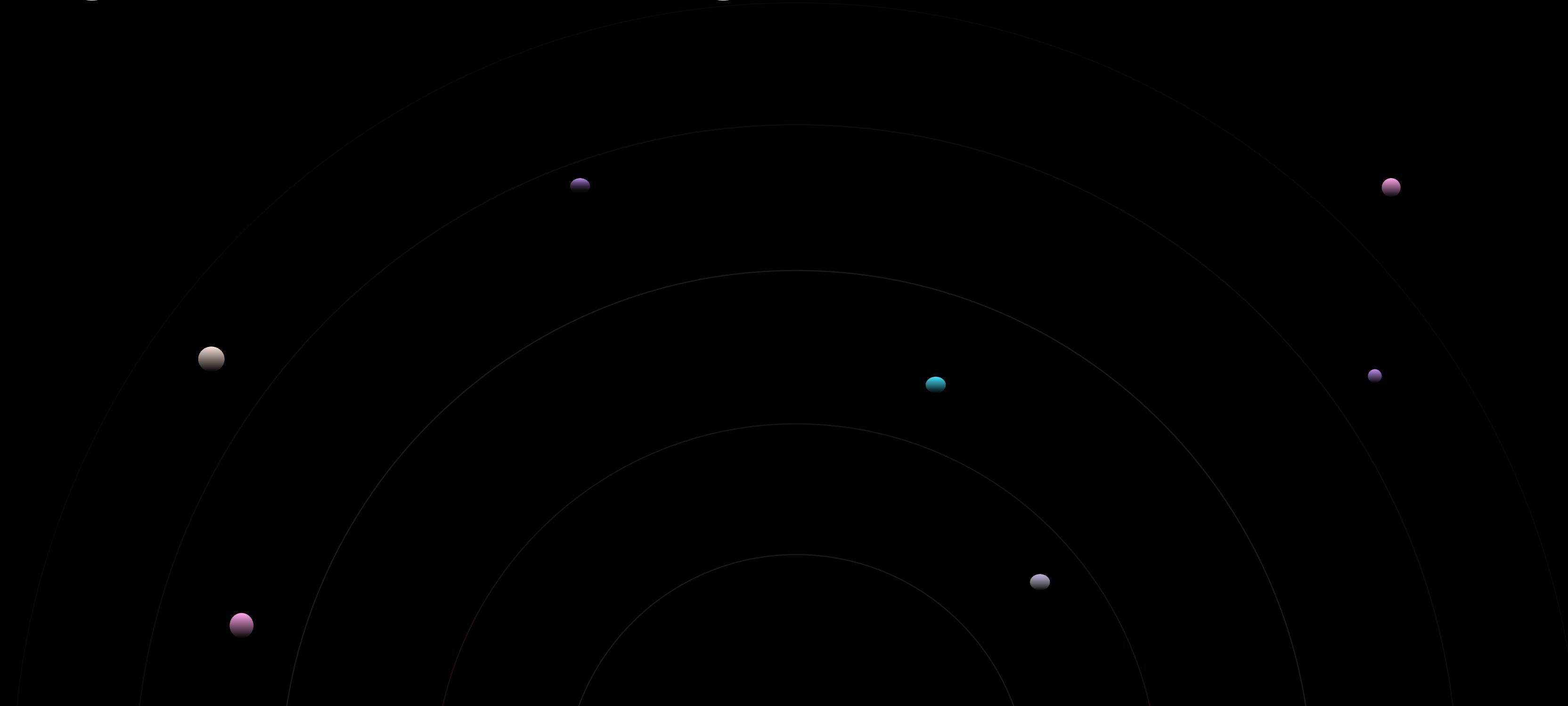
Vs. Pure Buy

Metric	In-House Solution	LlamaCloud
Dev Time for Data Sources	2 weeks	~5-30 mins
Dev Time for Parsing	2-3 weeks	~5-30 mins
Dev Time for Chunking	2-3 weeks	~5-30 mins

LlamaCloud enables users to get **state-of-the-art quality over their data**.



Agentic Reasoning over Complex Inputs



Complex Inputs

Naive RAG works well for pointed questions, but fails on more complex tasks.

Summarization Questions: “Give me a summary of the entire <company> 10K annual report”

Comparison Questions: “Compare the open-source contributions of candidate A and candidate B”

Multi-part Questions: “Tell me about the pro-X arguments in article A, and tell me about the pro-Y arguments in article B, make a table based on our internal style guide, then generate your own conclusion based on these facts.”

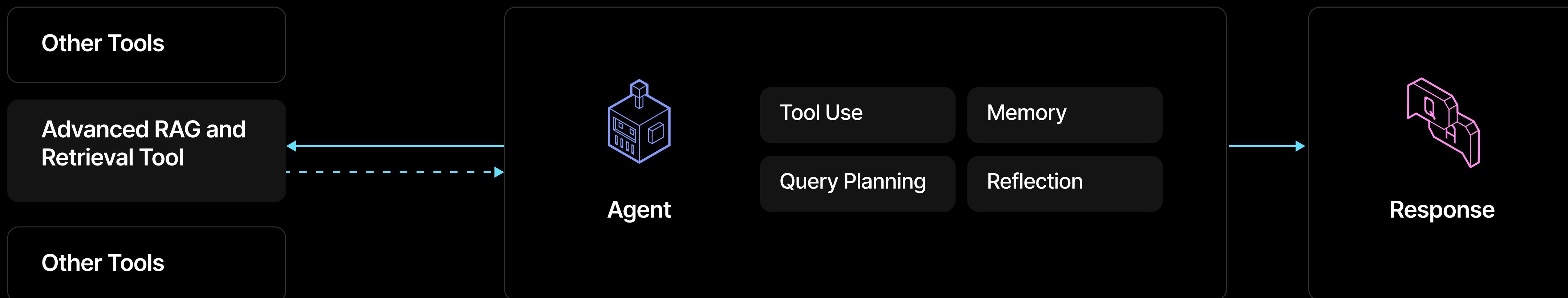
Research Tasks: “I want to create a research survey on current supervised fine-tuning techniques. Can you help?”

Agentic RAG

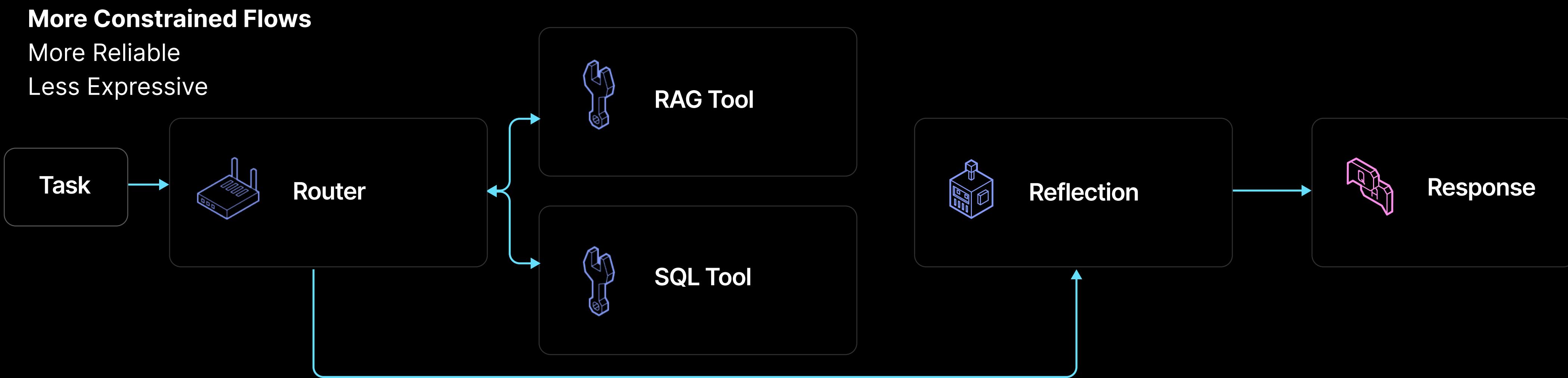
Every data interface is a tool

Use agent reasoning loops (sequential, DAG, tree) to tackle complex tasks

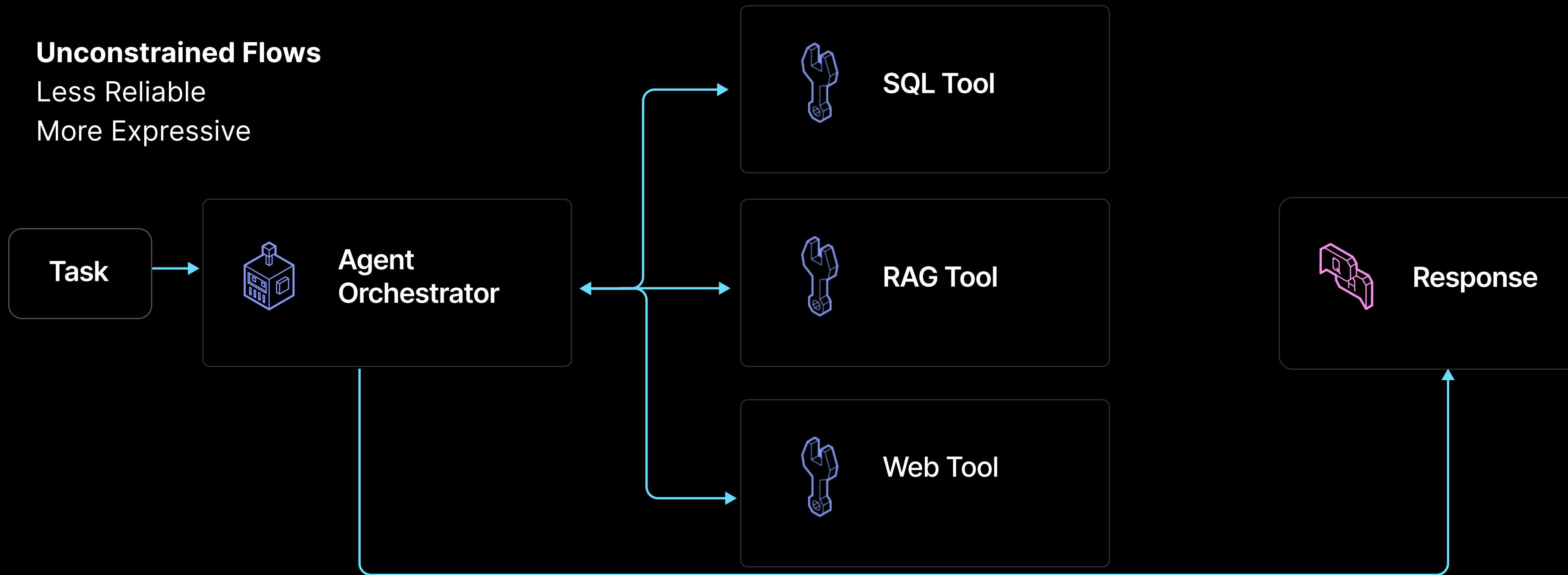
End Result: Build personalized QA systems capable of handling complex questions!



Unconstrained vs. Constrained Flows



Unconstrained vs. Constrained Flows

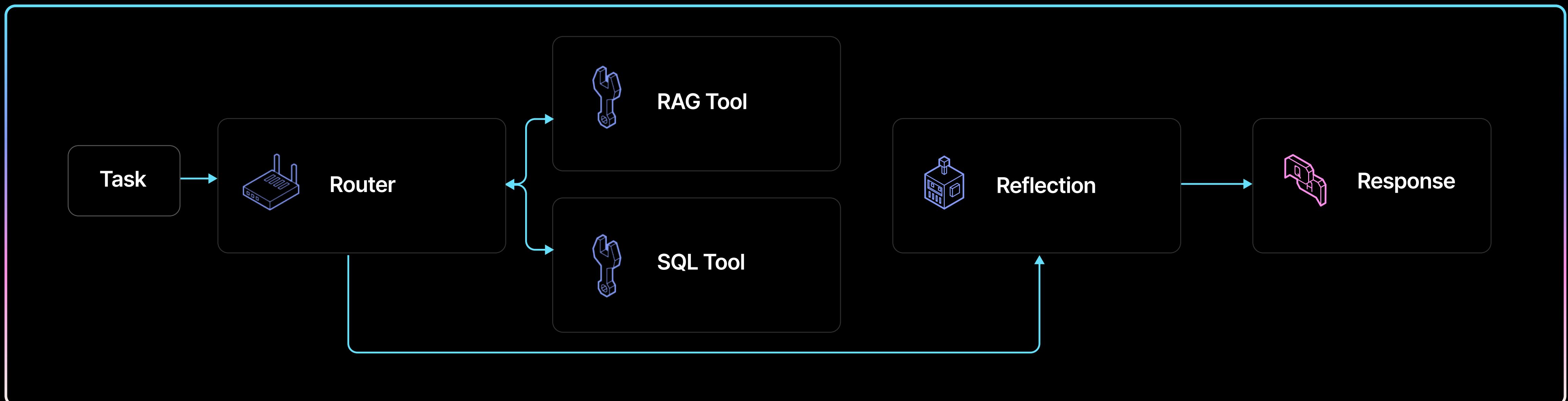


Agentic Orchestration Foundations

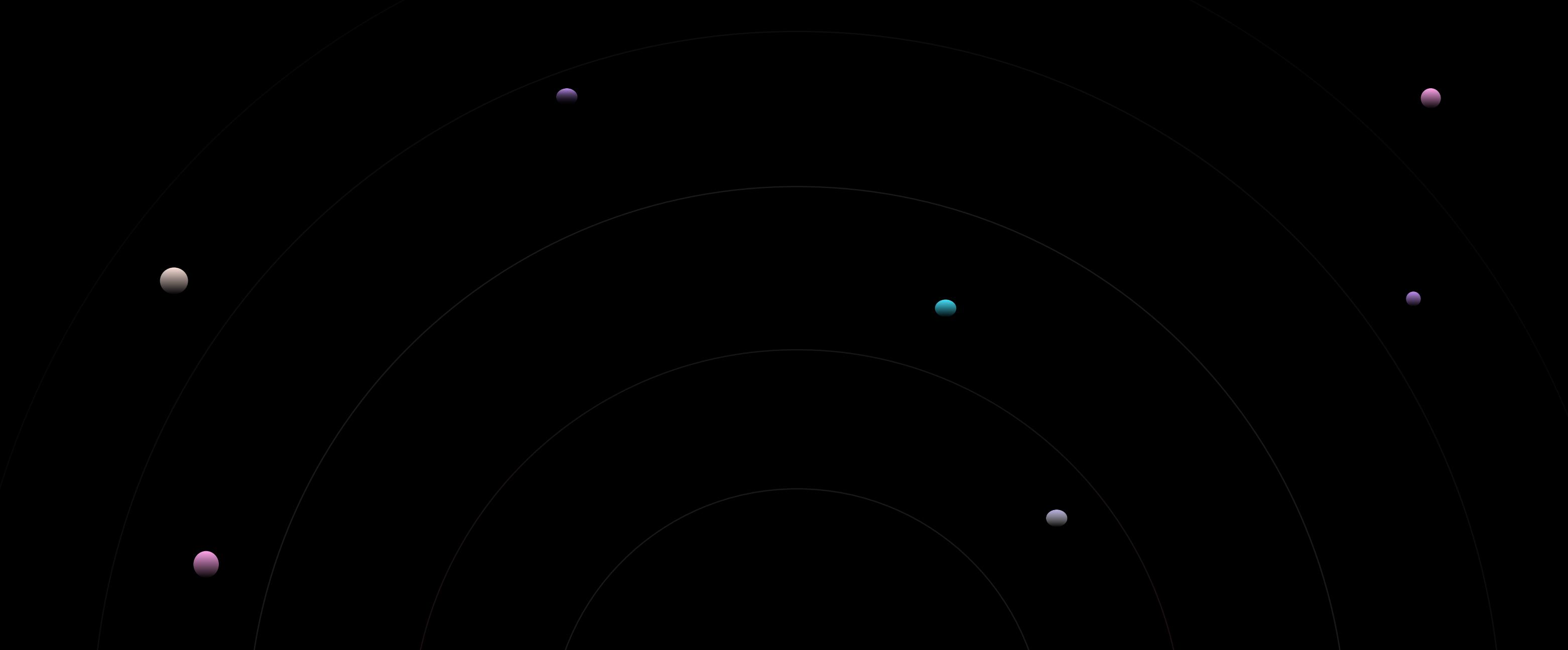
We believe an agent orchestration framework should have the following properties

- ✓ **Event-Driven:** Model each step as listening to input events and emitting output events
- ✓ **Composable:** Piece together granular workflows into higher-level workflows
- ✓ **Flexible:** Write logic through LLM calls or through plain Python
- ✓ **Code-first:** Express orchestration logic through code. Easy to read and easy to extend.
- ✓ **Debuggable and Observable:** Step through and observe states
- ✓ **Easily Deployable to Production:** Translate notebook code into services that run in production.

LlamaIndex Workflows



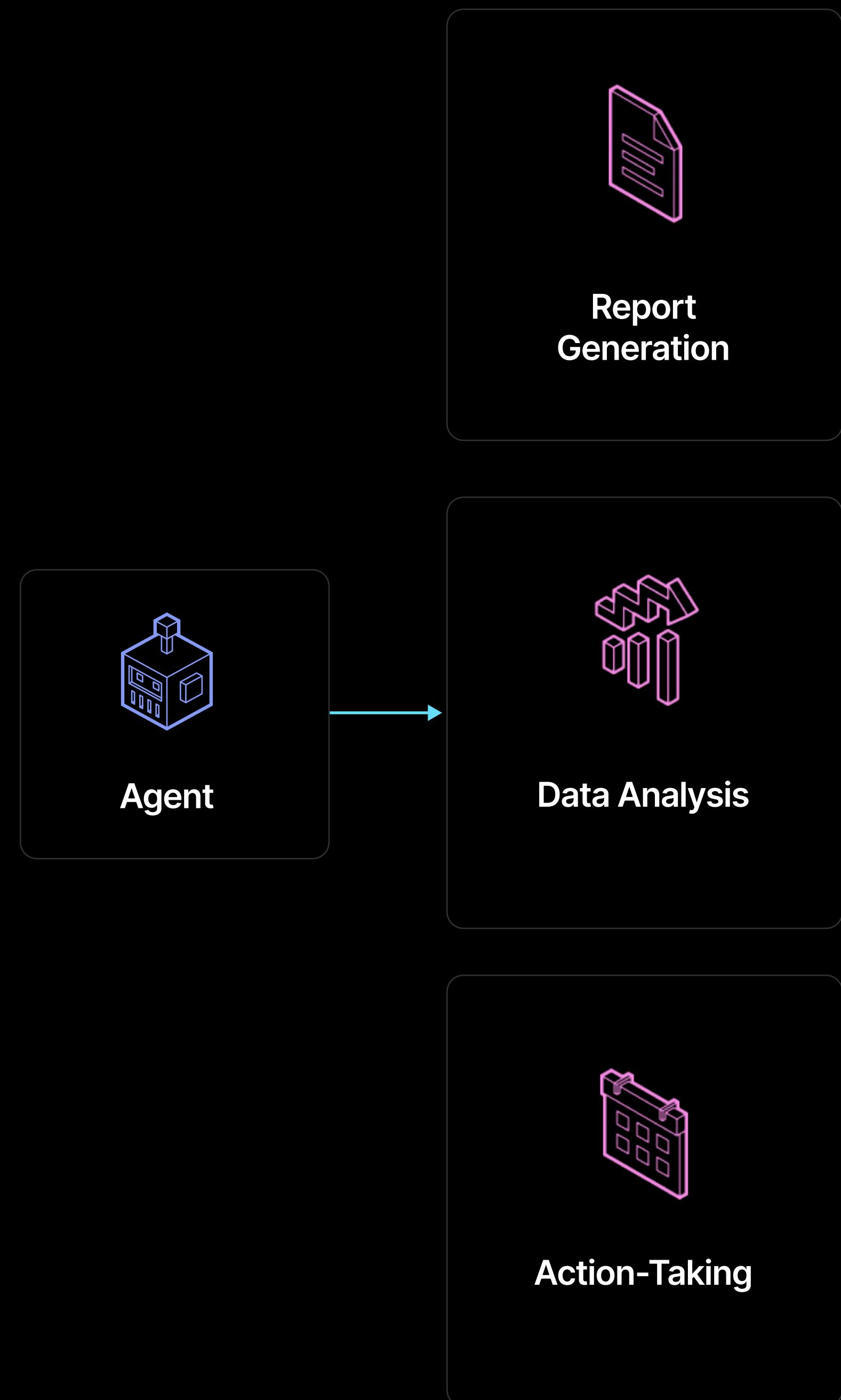
Agentic Decision Making and Output Generation



Automating Decision Making

Agents should have the capability to not only generate chatbot responses, but also

1. Produce knowledge work
2. Take actions



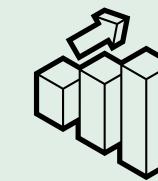
Report Generation

Report generation is one of the top enterprise agent use cases, and the natural next step for agentic RAG.

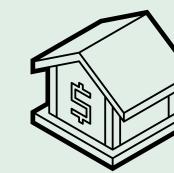
Use Cases:

1. Generate a full research report or presentation (with text, tables, images).
2. Fill out an example form or questionnaire.
3. Fill out an Excel sheet

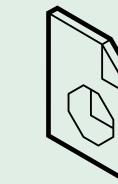
Top Use Case among our Customers



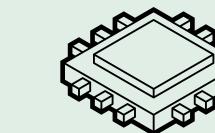
**Fortune 500
Investment Firm**



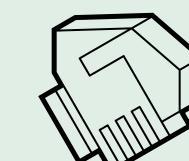
**Global
investment
bank**



**Big 4
Consulting
Firm**



**Leading Chinese Social
Media Company**



**Global 2000 Construction
Company**

Example: Multimodal Report Generation

Generate interleaving text-and-image responses with the help of **structured outputs**.

Example architecture: research and writer steps

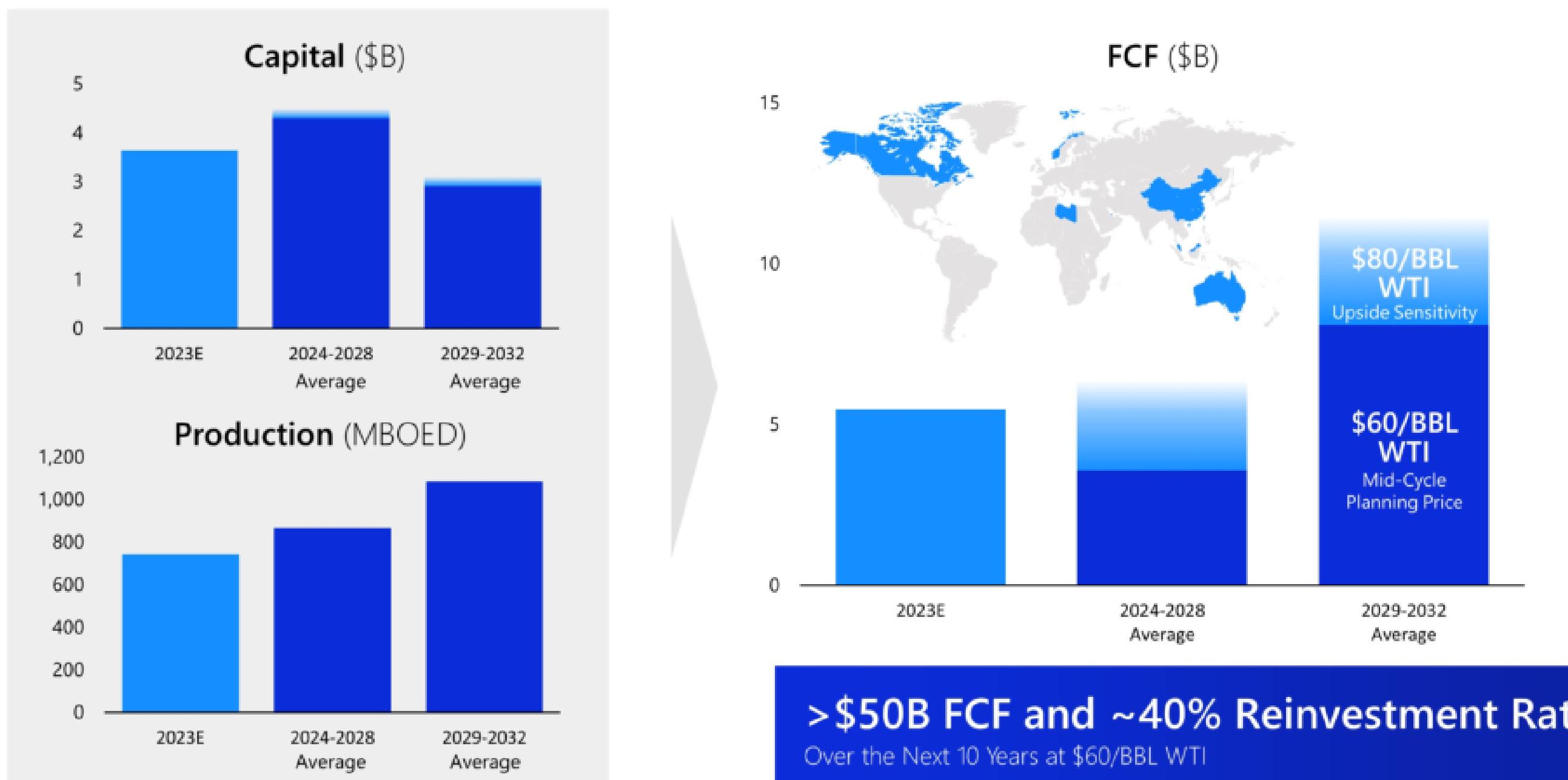
1. The **researcher** retrieves relevant chunks and documents, and puts them into a data cache.
2. The **writer** uses the data cache to generate a structured output of interleaving text and image blocks.

https://github.com/run-llama/llama_parse/blob/main/examples/multimodal/multimodal_report_generation_agent.ipynb

Lower 48 Segment

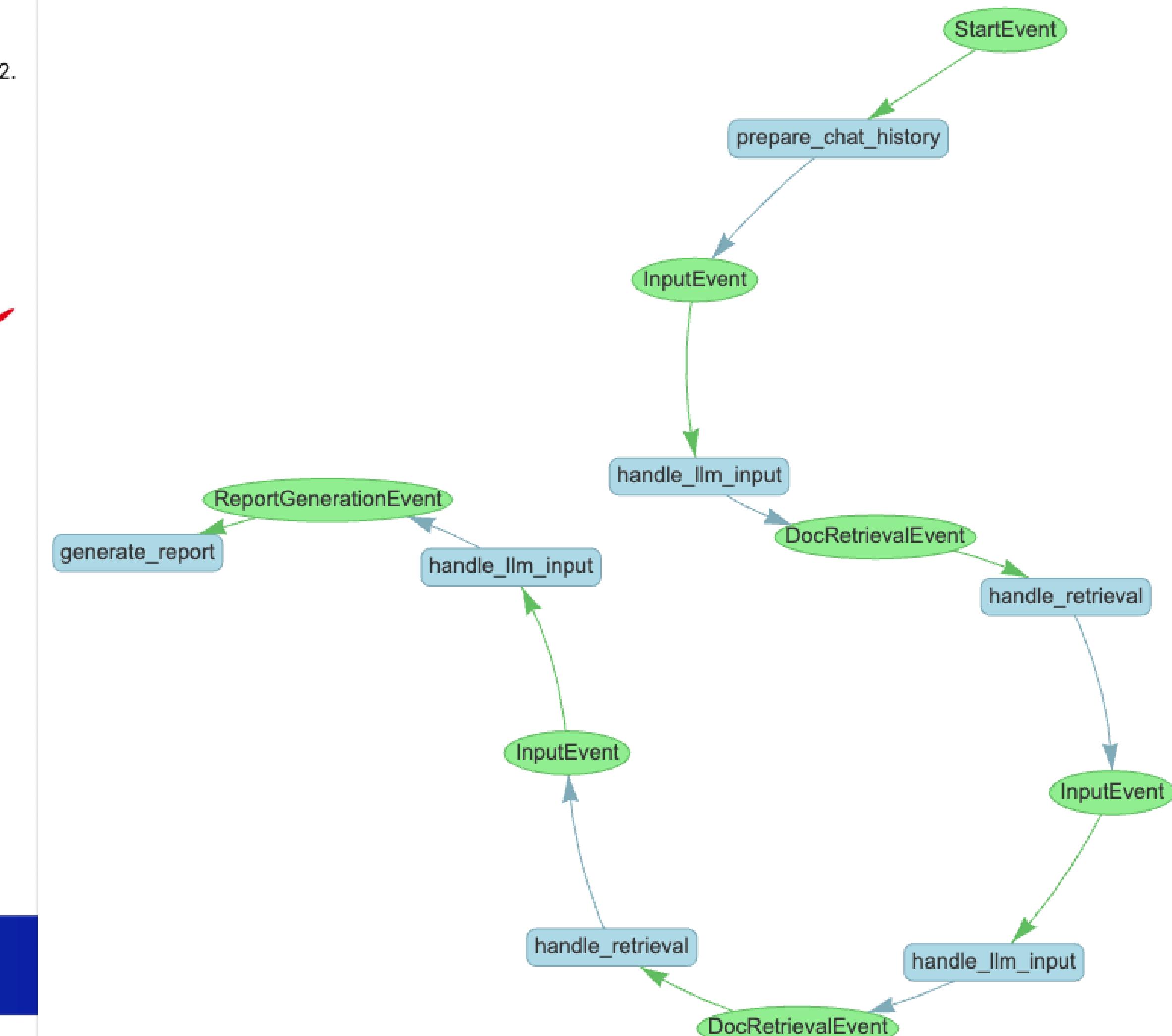
- **Capital Expenditures:** Expected to be *6.3 billion in 2023, averaging 6.5 billion from 2024-2028, and \$8.1 billion from 2029-2032.*
- **Production:** Projected to be around 1050 MBOED in 2023, increasing to 1220 MBOED on average from 2024-2028, and reaching 1530 MBOED on average from 2029-2032.
- **Free Cash Flow (FCF):** Estimated at *7 billion in 2023, averaging 5.5 billion from 2024-2028, and \$8 billion from 2029-2032.*
- **Key Projects:** Focused on the Permian Basin, Eagle Ford, and Bakken, with significant investments in technology and emissions reductions.

Alaska and International: Our Unique Diversification Advantage



Free cash flow (FCF) and reinvestment rate are non-GAAP measures defined in the Appendix.

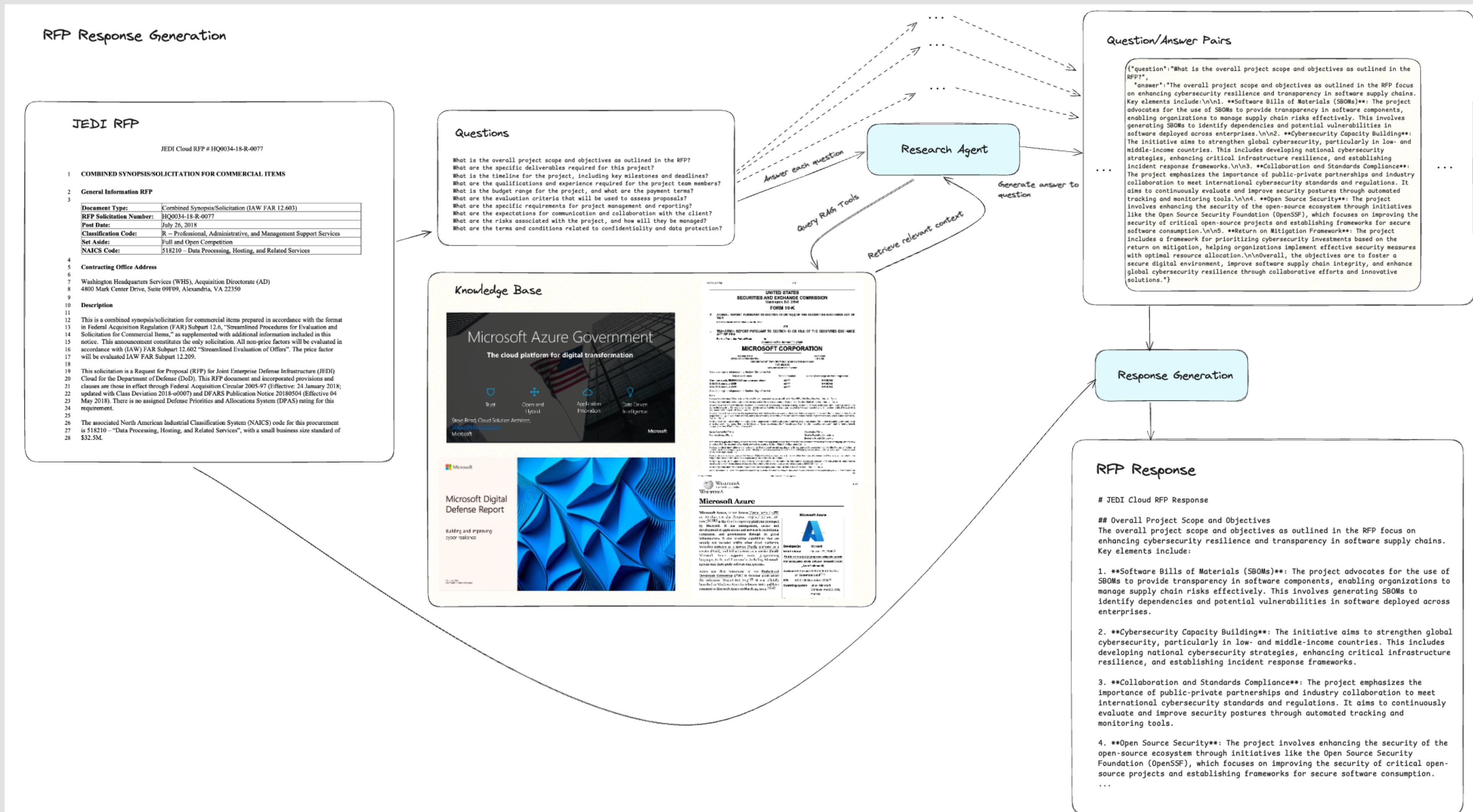
ConocoPhillips 29



Example: RFP Response Generation

As a vendor, generate a response that adheres to guidelines outlined in a Request for Proposal (RFP).

https://github.com/run-llama/llama_parse/blob/main/examples/report_generation/rfp_response/generate_rfp.ipynb



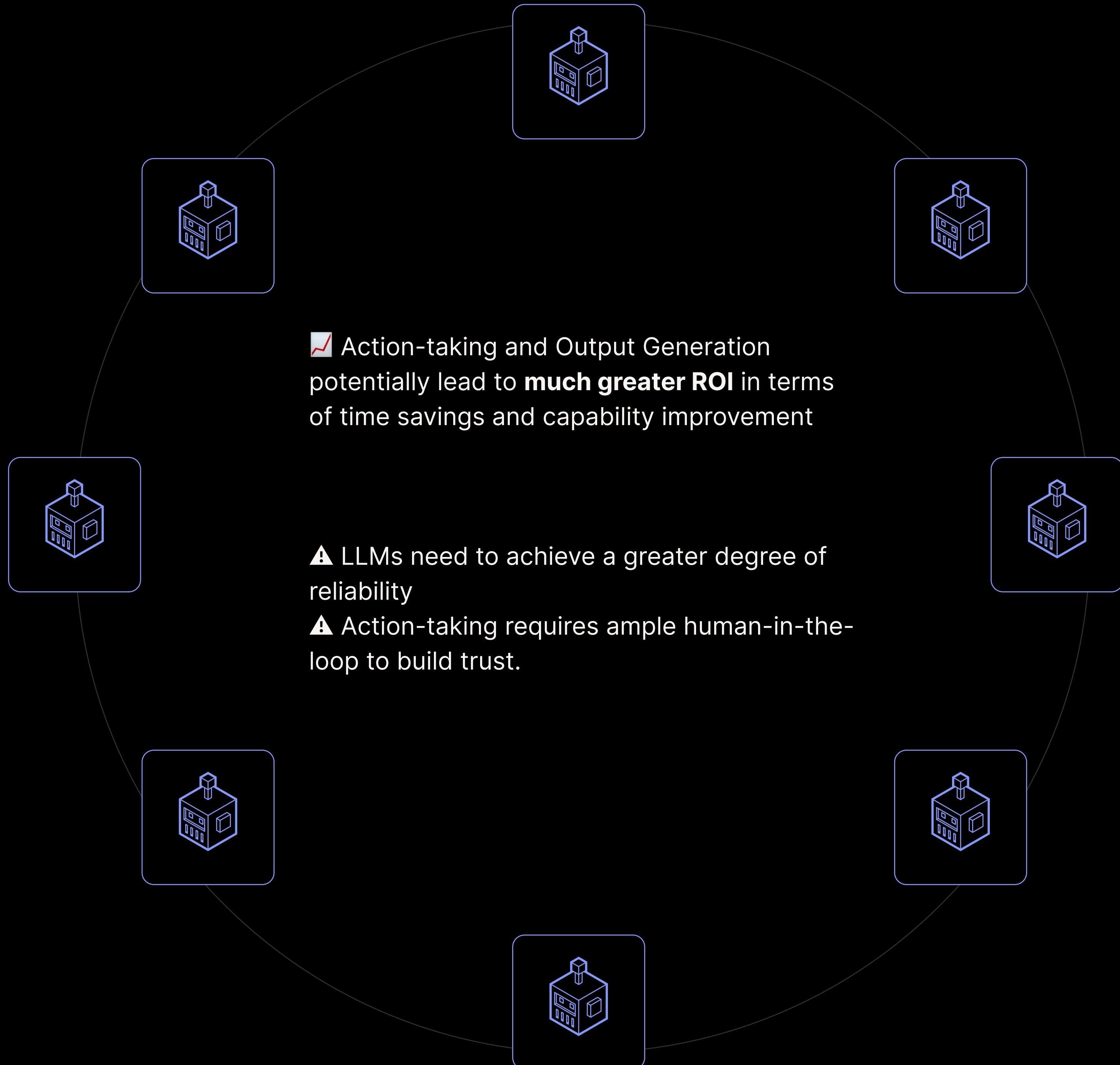
Example: Excel Form Filling

Parse and fill in a structured Excel template for financial comparisons

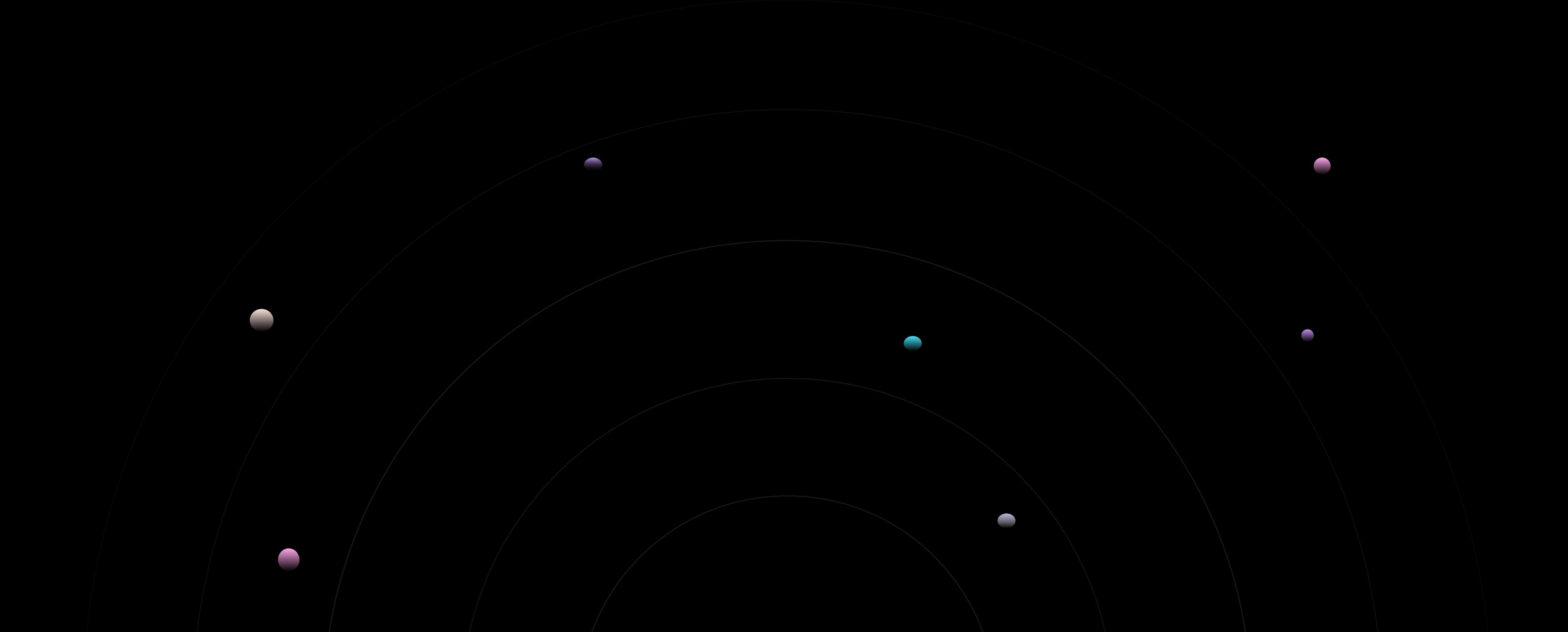
Example notebook: https://github.com/run-llama/llamacloud-demo/blob/main/examples/form_filling/Form_Filling_10K_SEC.ipynb

Parameter	2021			2022			Parameter	2021			2022								
	Amazon (AMZN)	Microsoft (MSFT)	Apple (AAPL)	Amazon (AMZN)	Microsoft (MSFT)	Apple (AAPL)		Amazon (AMZN)	Microsoft (MSFT)	Apple (AAPL)	Amazon (AMZN)	Microsoft (MSFT)	Apple (AAPL)						
1. Revenue							1. Revenue												
2. Net Income							2. Net Income												
3. Earnings Per Share (EPS)							3. Earnings Per Share (EPS)												
4. EBITDA							4. EBITDA												
5. Free Cash Flow							5. Free Cash Flow												
6. Return on Equity (ROE)							6. Return on Equity (ROE)												
7. Return on Assets (ROA)							7. Return on Assets (ROA)												
8. Debt-to-Equity Ratio							8. Debt-to-Equity Ratio												
9. Current Ratio							9. Current Ratio												
10. Gross Margin							10. Gross Margin												
11. Operating Margin							11. Operating Margin												
12. Net Profit Margin							12. Net Profit Margin												
13. Inventory Turnover							13. Inventory Turnover												
14. Accounts Receivable Turnover							14. Accounts Receivable Turnover												
15. Capital Expenditures							15. Capital Expenditures												
16. Research and Development Expenses							16. Research and Development Expenses												
17. Market Cap							17. Market Cap												
18. Price-to-Earnings (P/E) Ratio							18. Price-to-Earnings (P/E) Ratio												
19. Dividend Yield							19. Dividend Yield												
20. Year-over-Year Growth Rate							20. Year-over-Year Growth Rate												
Year	Company	Accounts Receivable Turnover	Capital Expenditures	Current Ratio	Debt-to-Equity Ratio	Dividend Yield	EBITDA	Earnings Per Share (EPS)	Free Cash Flow	Gross Margin	Inventory Turnover	Market Cap	Net Income	Net Profit Margin	Operating Margin	Price-to-Earnings (P/E) Ratio	Research and Development Expenses	Return on Assets (ROA)	
0	2021	Amazon (AMZN)	5.6 times	\$58.3 billion	1.02	0.63	NaN	24,879	64.81	-\$11,569	41.0%	4.8 times	\$1.66 trillion	33,364	6.8%	13.7%	73.60	Not significant	5.9%
1	2021	Microsoft (MSFT)	7.5 times	\$9.5 billion	1.79	0.30	1. 0.8%	\$76,632 million	8.05	28.7 billion	\$115.9 billion	2.5 times	\$2.5 trillion	\$61,271 million	19.7%	32%	34.50	20,716 million	10.9%
2	2021	Apple (AAPL)	6.2 times	\$9,000 million	1.12	0.93	1. 0.0065	19,863	5.61	\$73,000	\$152,836 million	Not available	\$2.46 trillion	94,680 million	21.7%	44.7%	28.11	\$21,914 million	5.1%
3	2022	Amazon (AMZN)	6.4 times	\$58.3 billion	1.10	1.39	NaN	\$15,432 million	9.70	-\$11,569 million	NaN	NaN	\$1.47 trillion	(2,722)	NaN	16.9%	NaN	\$73,213 million	6.9%
4	2022	Microsoft (MSFT)	6.1 times	\$8.5 billion	1.78	0.59	2.0%	\$107,895 billion	9.65	\$58.7 billion	\$135,620 billion	3.4 times	\$1.87 trillion	72,738	19%	19%	38.60	\$24,512 million	9.99%
5	2022	Apple (AAPL)	6.0 times	\$42,117 million	2.78	0.68	0.0054	\$145,787 million	6.15	\$88,531	\$170,782 million	1. 6.87	\$2.9 trillion	99,803	21.9%	37.9%	24.60	\$26,251 million	6.7%

Benefits and Risks



Towards a Scalable, Full-Stack Application

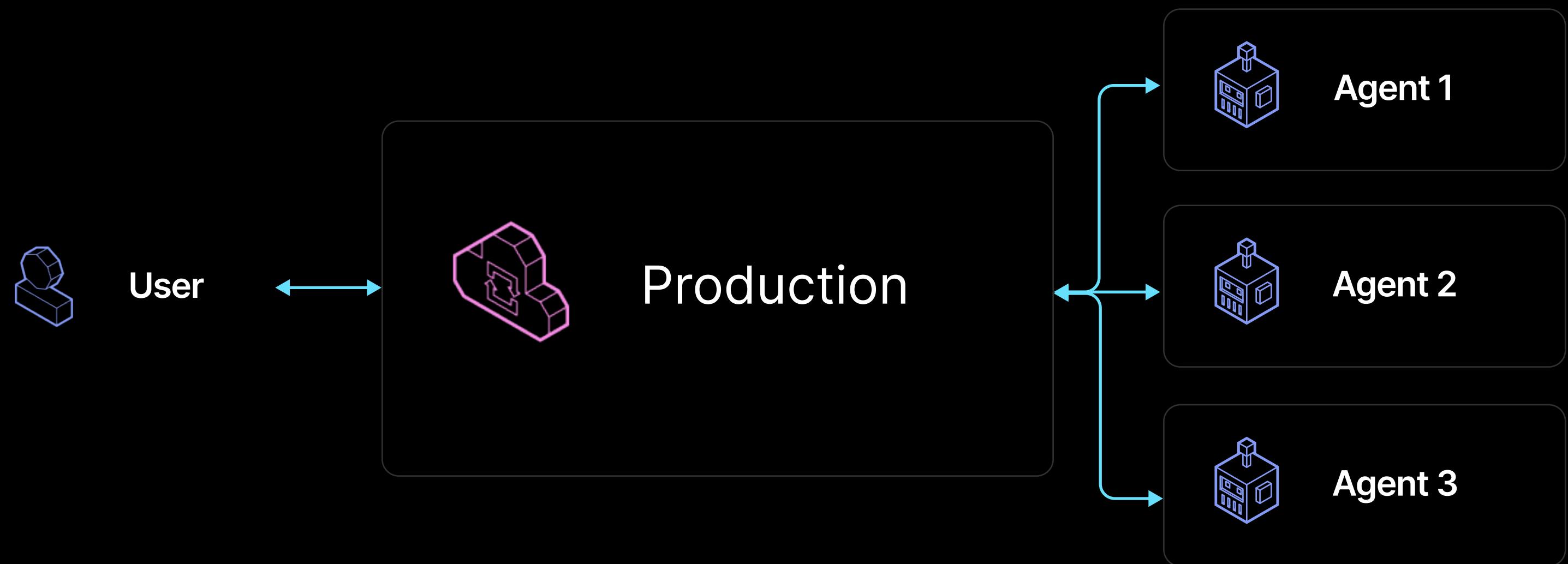
The background features a dark gray gradient with several thin, light gray concentric circles. Scattered across these circles are small, semi-transparent colored dots in shades of purple, pink, teal, and light gray.

Running Agents in Production

You need the right architecture and infra components to serve complex, agentic workflows to end-users as a production application.

Requirements:

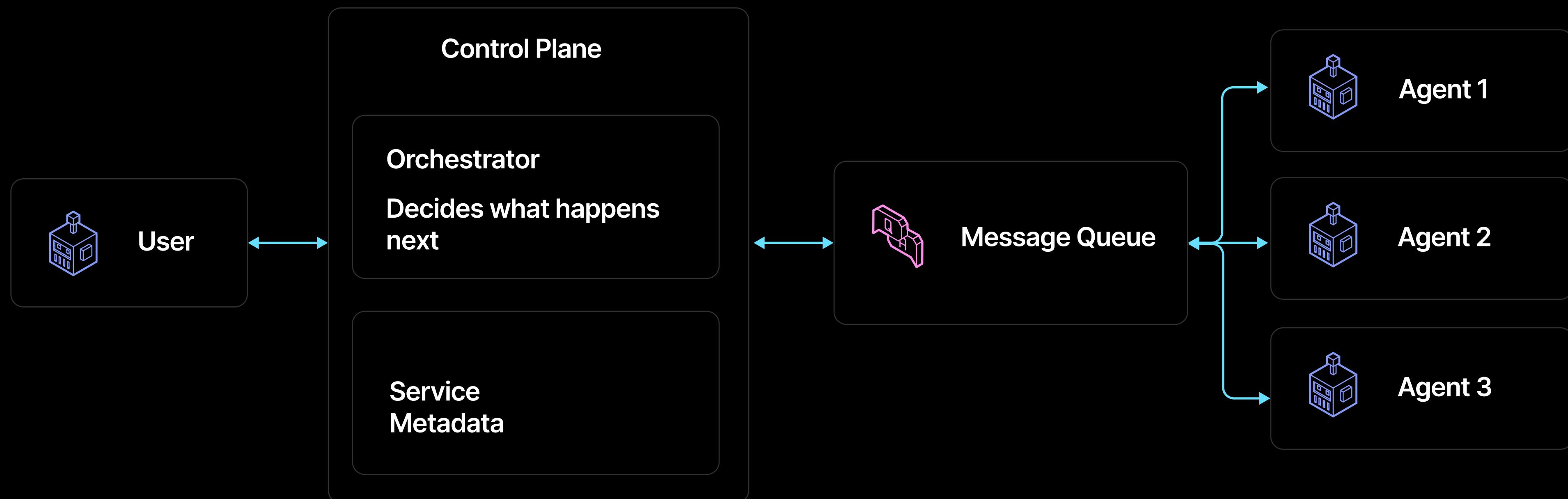
1. Encapsulation and re-use
2. Standardized communication interfaces between agents and with the client.
3. Scalability in number of users and number of agents
4. Human-in-the-loop for the end-user
5. Debugging and observability tools for the developer



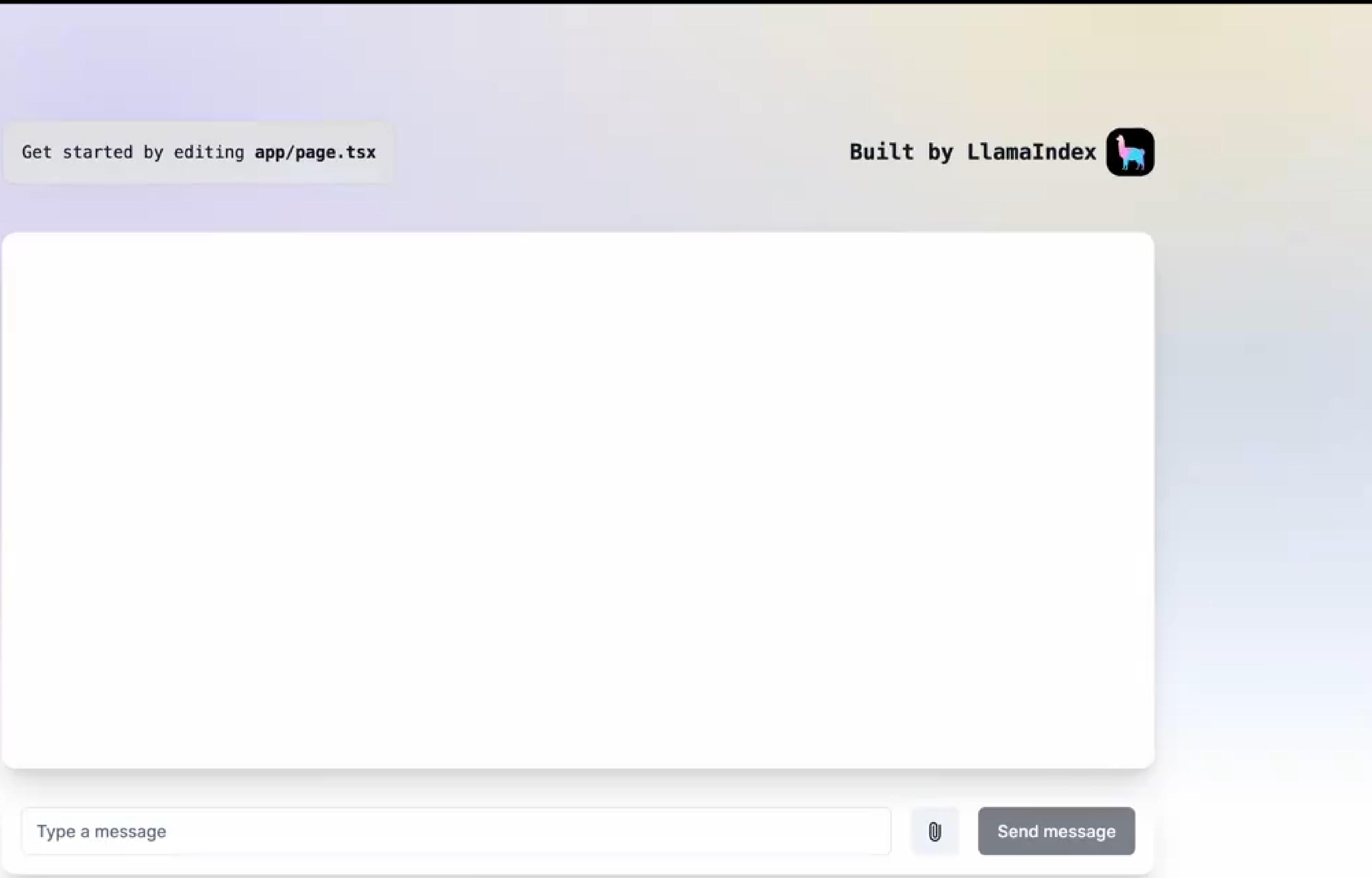
llama-deploy

Deploy agentic workflows as **microservices**.

- Model every agent workflow as a service API
- All agent communication occurs via a central message queue
- Distributed tool-execution
- Human-in-the-loop as a service
- Easy deployment with docker-compose and Kubernetes



- https://github.com/run-llama/llama_deploy



Add knowledge or test the chat below. Once you're satisfied, [start the app or use the API](#).

Agents

Configure tools and agents

Researcher Analyst Reporter [+](#)

Agent Name
Researcher
Unique name to identify the agent.

Agent Role
Expert in researching news material
Helps RAGapp to assign the right agent for a task.

System Prompt
You are an expert in researching material for news-related questions. Don't make up any information; instead, use your tools to gather the information.

Tools

Use Knowledge
Query information from your knowledge base

Web Search
Find information on the internet by searching DuckDuckGo

Wikipedia
Search for information on Wikipedia

OpenAPI
Make requests to external APIs using the information from the OpenAPI spec

Code Interpreter
Execute python code in a sandboxed environment using E2B code interpreter

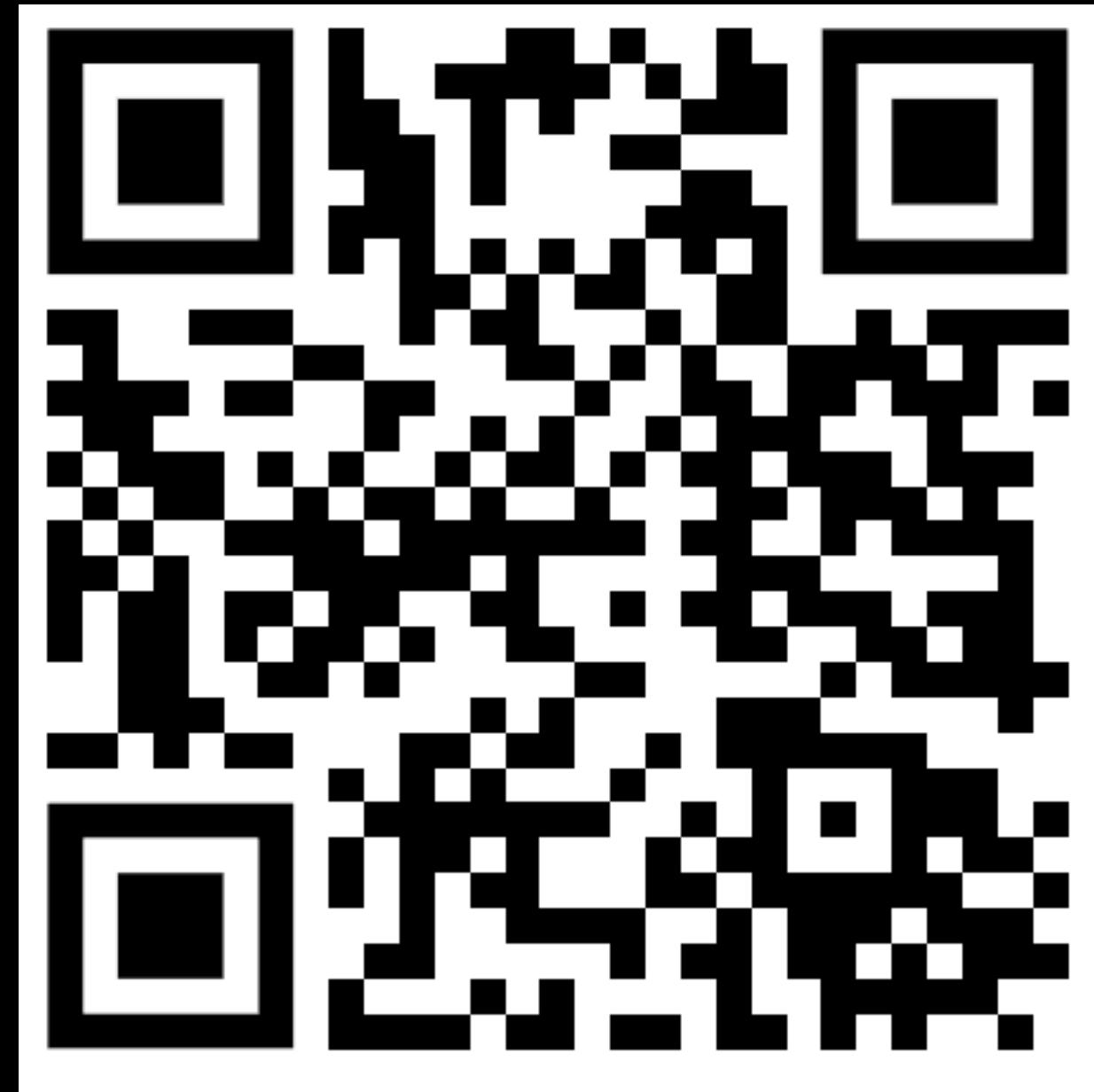
Image Generator
Generate images from the provided text using the Stability AI API

Knowledge

Type a message [0](#) [Send message](#)

Questions, feature requests or found a bug? Open an issue on [GitHub](#). © 2024 by Schiesser IT, LLC.

Thank you!



Get in Touch



LlamaCloud
Signup

