

## Building Production RAG Over Complex Documents

 Databricks  
142K subscribers

[Subscribe](#)

 1.7K   Share  Ask  Download ...

61,492 views Jul 23, 2024

Large Language Models (LLMs) are revolutionizing how users search for, interact with, and generate new content. Some recent stacks and toolkits around Retrieval-Augmented Generation (RAG) have emerged, enabling users to build applications such as chatbots using LLMs on their private data. However, while setting up naive RAG is straightforward, building production RAG is very challenging, especially as users scale to larger and more complex data sources. A classic example is a large number of PDFs with embedded tables. RAG is only as good as your data, and developers must carefully consider how to parse, ingest, and retrieve their data to successfully build RAG over complex documents. This session provides an in-depth exploration of this entire process; you will get an overview of the process around building a RAG pipeline that can handle messy, complicated PDF documents. This includes implementing a parsing strategy for parsing a complex document with embedded objects. This consists of an indexing strategy to process these documents beyond simple chunking techniques. We will then explore various advanced retrieval algorithms to handle questions about the tabular and unstructured data and discuss their use cases and tradeoffs.

# DATA+AI SUMMIT

DATA INTELLIGENCE FOR ALL

BY  databricks

DATA+AI  
SUMMIT  
BY  databricks

## Building Advanced RAG Over Complex Documents

Jerry Liu  
June 11, 2024

DATA+AI SUMMIT

©2024 Databricks Inc. — All rights reserved



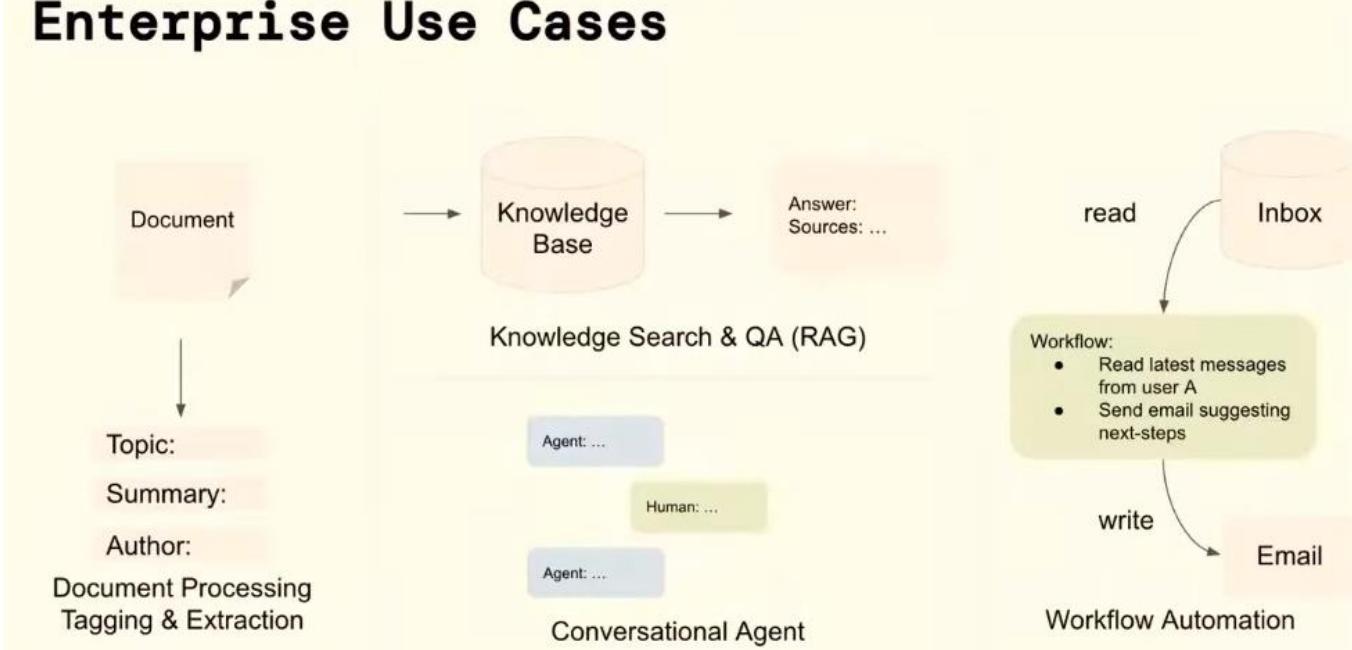
1

## Agenda

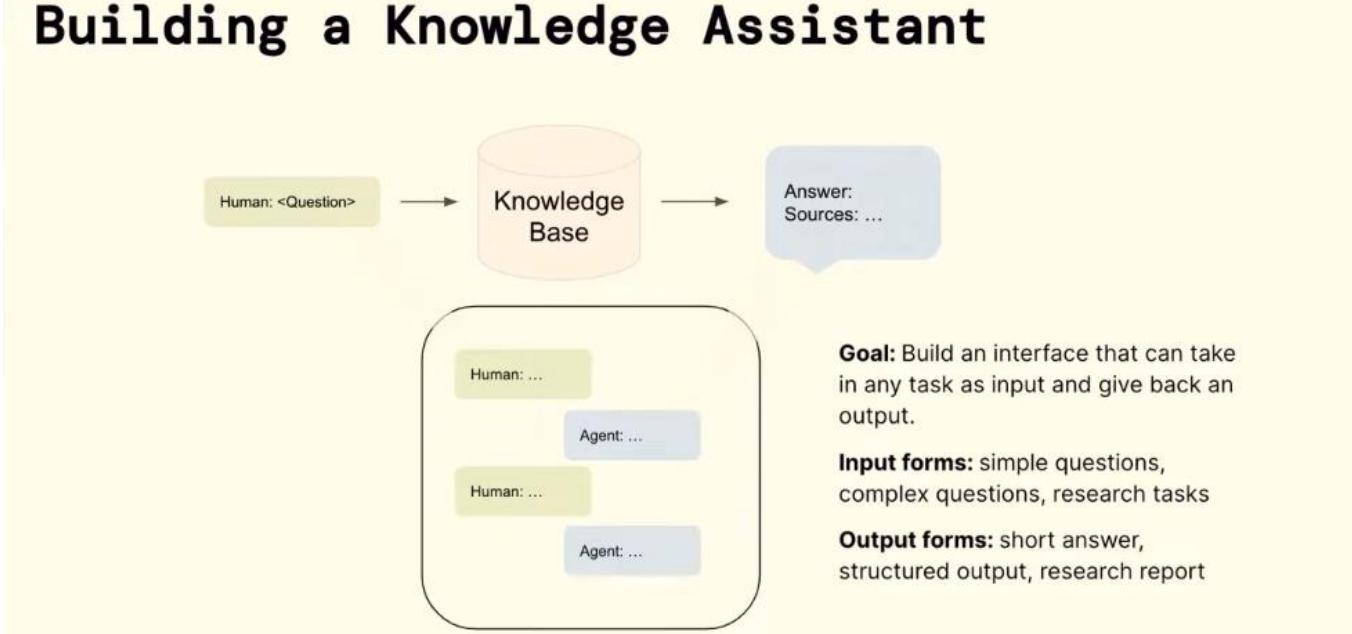
- 1. Building a Knowledge Assistant**
- 2. RAG Overview: Basic RAG and where it goes wrong**
- 3. Improving Data Quality:**
  - Improve LLM reasoning over complex data
  - **Workshop:** LlamaParse over Complex Documents
- 4. Improving Query Complexity:** from RAG to agents
  - **Workshop:** LlamaParse-powered document agent
- 5. What's next?**

We will be seeing how to build a full end to end RAG pipeline from the data side to the query side.

# Enterprise Use Cases

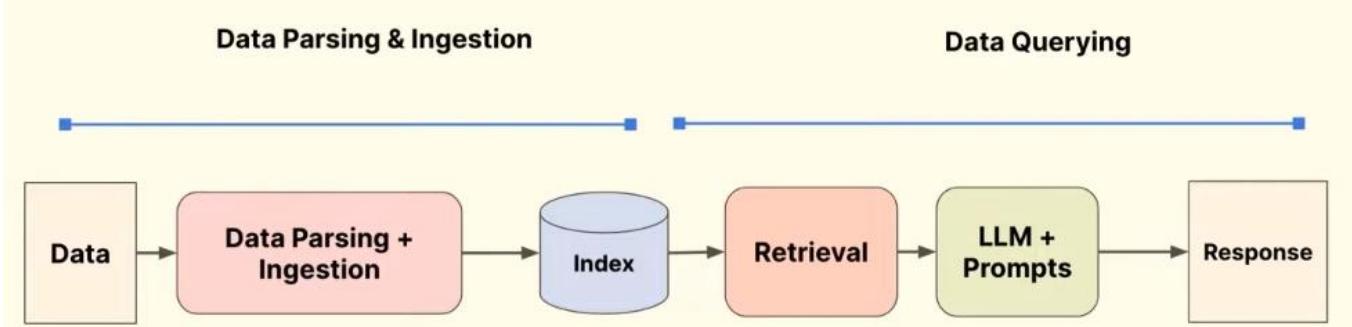


## Building a Knowledge Assistant

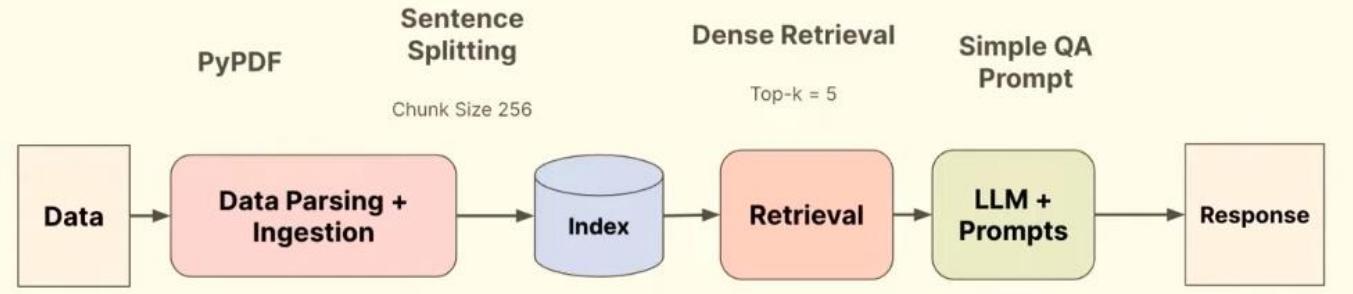


## Retrieval Augmented Generation (RAG)

An overview of a RAG Pipeline



# Naive RAG



## Challenges with Naive RAG

### Easy to Prototype, Hard to Productionize

Naive RAG approaches tend to work well for **simple** questions over a **simple, small** set of documents.

- "What are the main risk factors for Tesla?" (over Tesla 2021 10K)
- "What did the author do during his time at YC?" (Paul Graham essay)

But productionizing RAG over more questions and a larger set of data **is hard!**

### Easy to Prototype, Hard to Productionize

#### Failure Modes:

- Simple Questions over Complex Data
- Simple Questions over Multiple Documents
- Complex Questions

The top priority goal should be figuring out how to **get high-response quality** from the set of representative questions you want to ask.

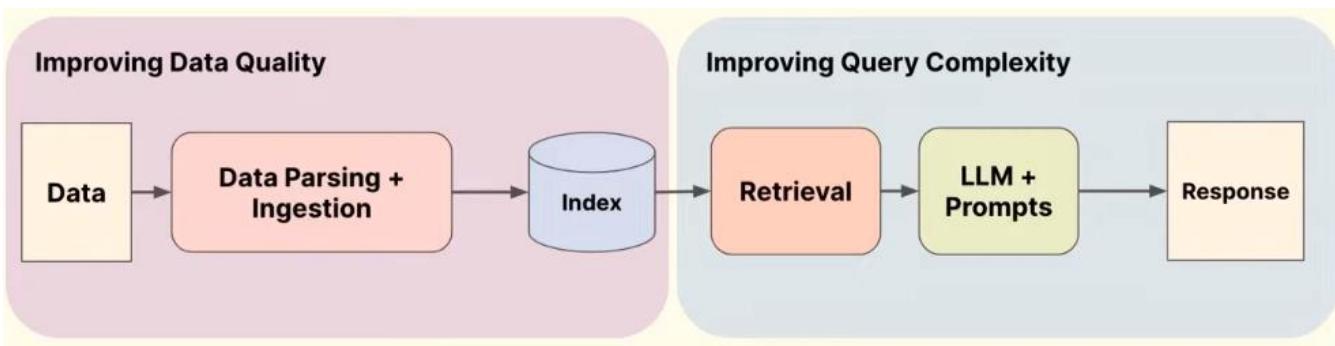
# Can we do more?

In the naive setting, RAG is boring.

- 🚫 It's just a glorified search system
- 🚫 There's many questions/tasks that naive RAG can't give an answer to.

💡 Can we go beyond simple search/QA to building a general context-augmented research assistant?

## Main Focus Areas

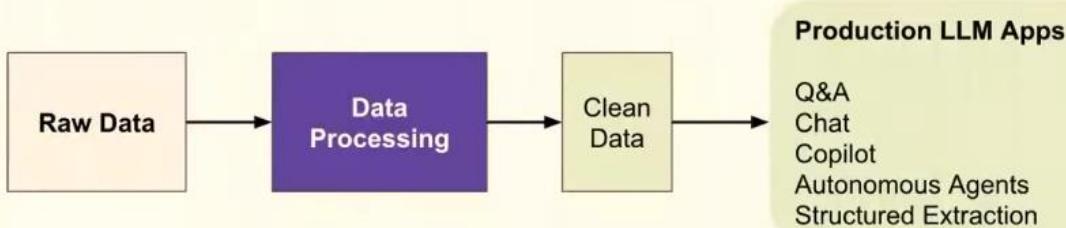


## Improving Data Quality

### RAG is only as Good as your Data

Garbage in = garbage out

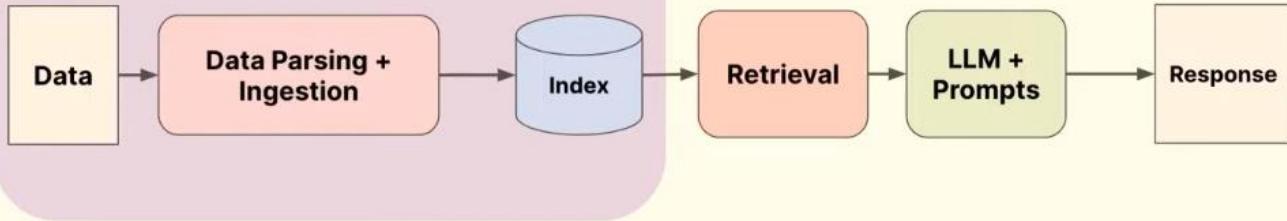
Good data quality is a **necessary** component of any production LLM app.



# RAG is only as Good as your Data

## Main Components of Data Processing:

- Parsing
- Chunking
- Indexing



You can add metadata to your chunks to provide more details about the parent document ID.

## General Principles

### Parsing:

- Bad parsers are a key cause of garbage in == garbage out.
- Badly formatted text/tables confuse even the best LLMs

### Chunking:

- Try to preserve semantically similar content.
  - [5 Levels of Text Splitting](#)
- **Strong baseline:** page-level chunking.

### Indexing:

- Raw text oftentimes confuse the embedding model.
- Don't just embed the raw text, embed **references**.
- Having multiple embeddings point to the same chunk is a **good practice!**

## Case Study: Complex Documents

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

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

Naive RAG indexing pipelines fail over these documents.

Let's build an **advanced RAG indexing pipeline**.

Item	Liabilities (in 000's of CHF)		
	31 Dec 2022	31 Dec 2021	Change
Payables and accruals	4,685	4,066	619
Employee benefits	(127,215)	84,676	42,539
Contributions received in advance	6,975	10,192	(3,217)
Unearned revenue from exchange transactions	20	651	(631)
Deferred Revenue	(71,301)	55,737	15,564
Borrowings	28,229	29,002	(773)
Funds held in trust	30,373	29,014	1,359
Provisions	1,706	1,910	(204)
<b>Total Liabilities</b>	<b>270,504</b>	<b>215,248</b>	<b>55,256</b>

# Most PDF Parsing is Inadequate

Extracts into a messy format that is impossible to pass down into more advanced ingestion/retrieval algorithms.

Please find below AXA's rankings and market shares in the main countries where it operates:

Property & Casualty Ranking	Market share Ranking	Life & Savings	
		Market share	Life Sources
France	2	11.9	8.4 "Tous Assureurs" as of December 31, 2022.
Sweden	1	11.1	Market share based on statutory premiums and market approach by STA (Swiss Insurance Association) Figures as of December 31, 2022.
Germany	6	9.8	GIV (German association of insurance companies) as of December 31, 2021.
Belgium	1	11.7	AXA Belgium Professional Union of Insurance Companies based on gross written premium as of September 30, 2022.
United Kingdom	4	8.2	AXA UK Insurance Company Competitor Analytics 2021, AXA Data, as of December 31, 2021.
Ireland	1	11.9	AXA Insurance Ireland Ltd. Statistics 2021 as of December 31, 2021.
Spain	5	4.0	1.1 as of December 31, 2022.
Italy	3	3.8	Associazione Nazionale Imprese Assicurative (ANIA) Insurance statistics as of December 31, 2022.
Japan	11	0.6	Discounted financial reports including Kango Life
Hong Kong	1	7.0	Insurance Authority statistics based on gross written premiums as of December 31, 2021.
US Insurance in the United States	16	1.8	AM Best 2021 as of December 31, 2021, in the United States in Commercial Lines.
Switzerland	24	2.3	AM Best 2021 as of December 31, 2021.
Thailand	18	1.0	7.2 TGA (The General Insurance Association) as of December 31, 2022 and TLA (The Life Association) as of December 31, 2022.
Indonesia	n/a	n/a	8.7 AXA Statistics measured on Weighted New Business Premium as of September 30, 2022.
Philippines	n/a	n/a	9.8 Insurance statistics measured on total premium income as of June 30, 2022.
China	n/a	0.4	CBIRC (China Banking and Insurance Regulatory Commission) as of December 31, 2021.
Mexico	5	8.0	12.3 RMS (Reservas Mexicanas de Instituciones de Seguros) as of September 30, 2022.
Brazil	13	1.4	10.97 (Superintendencia de Seguros Privados) as of September 2020.

Please find below AXA's rankings and market shares in the main countries where it operates:

Property & Casualty Ranking	Market share Ranking	Life & Savings Ranking	Market share Ranking	Life Sources Ranking
France	2	11.9	8.4 "Tous Assureurs" as of December 31, 2022.	8.4 "Tous Assureurs" as of December 31, 2022.
Sweden	1	11.1	Market share based on statutory premiums and market approach by STA (Swiss Insurance Association) Figures as of December 31, 2022.	8.2 "Tous Assureurs" as of December 31, 2022.
Germany	6	9.8	GIV (German association of insurance companies) as of December 31, 2021.	8.3 "Tous Assureurs" as of December 31, 2021.
Belgium	1	11.7	AXA Belgium Professional Union of Insurance Companies based on gross written premium as of September 30, 2022.	8.7 "Tous Assureurs" as of December 31, 2022.
United Kingdom	4	8.2	AXA UK Insurance Company Competitor Analytics 2021, AXA Data, as of December 31, 2021.	8.0 "Tous Assureurs" as of December 31, 2021.
Ireland	1	11.9	AXA Insurance Ireland Ltd. Statistics 2021 as of December 31, 2021.	8.1 "Tous Assureurs" as of December 31, 2021.
Spain	5	4.0	1.1 as of December 31, 2022.	8.2 "Tous Assureurs" as of December 31, 2021.
Italy	3	3.8	Associazione Nazionale Imprese Assicurative (ANIA) Insurance statistics as of December 31, 2022.	8.3 "Tous Assureurs" as of December 31, 2021.
Japan	11	0.6	Discounted financial reports including Kango Life	8.4 "Tous Assureurs" as of December 31, 2021.
Hong Kong	1	7.0	Insurance Authority statistics based on gross written premiums as of December 31, 2021.	8.5 "Tous Assureurs" as of December 31, 2021.
US Insurance in the United States	16	1.8	AM Best 2021 as of December 31, 2021, in the United States in Commercial Lines.	8.6 "Tous Assureurs" as of December 31, 2021.
Switzerland	24	2.3	AM Best 2021 as of December 31, 2021.	8.7 "Tous Assureurs" as of December 31, 2021.
Thailand	18	1.0	7.2 TGA (The General Insurance Association) as of December 31, 2022 and TLA (The Life Association) as of December 31, 2022.	8.8 "Tous Assureurs" as of December 31, 2021.
Indonesia	n/a	n/a	8.7 AXA Statistics measured on Weighted New Business Premium as of September 30, 2022.	8.9 "Tous Assureurs" as of December 31, 2021.
Philippines	n/a	n/a	9.8 Insurance statistics measured on total premium income as of June 30, 2022.	9.0 "Tous Assureurs" as of December 31, 2021.
China	n/a	0.4	CBIRC (China Banking and Insurance Regulatory Commission) as of December 31, 2021.	0.5 "Tous Assureurs" as of December 31, 2021.
Mexico	5	8.0	12.3 RMS (Reservas Mexicanas de Instituciones de Seguros) as of September 30, 2022.	12.4 "Tous Assureurs" as of December 31, 2021.
Brazil	13	1.4	10.97 (Superintendencia de Seguros Privados) as of September 2020.	11.0 "Tous Assureurs" as of December 31, 2021.

Please find below AXA's rankings and market shares in the main countries where it operates:

Property & Casualty Ranking	Market share Ranking	Life & Savings Ranking	Market share Ranking	Life Sources Ranking
France	2	11.9	8.4 "Tous Assureurs" as of December 31, 2022.	8.4 "Tous Assureurs" as of December 31, 2022.
Sweden	1	11.1	Market share based on statutory premiums and market approach by STA (Swiss Insurance Association) Figures as of December 31, 2022.	8.2 "Tous Assureurs" as of December 31, 2022.
Germany	6	9.8	GIV (German association of insurance companies) as of December 31, 2021.	8.3 "Tous Assureurs" as of December 31, 2021.
Belgium	1	11.7	AXA Belgium Professional Union of Insurance Companies based on gross written premium as of September 30, 2022.	8.7 "Tous Assureurs" as of December 31, 2022.
United Kingdom	4	8.2	AXA UK Insurance Company Competitor Analytics 2021, AXA Data, as of December 31, 2021.	8.0 "Tous Assureurs" as of December 31, 2021.
Ireland	1	11.9	AXA Insurance Ireland Ltd. Statistics 2021 as of December 31, 2021.	8.1 "Tous Assureurs" as of December 31, 2021.
Spain	5	4.0	1.1 as of December 31, 2022.	8.2 "Tous Assureurs" as of December 31, 2021.
Italy	3	3.8	Associazione Nazionale Imprese Assicurative (ANIA) Insurance statistics as of December 31, 2022.	8.3 "Tous Assureurs" as of December 31, 2021.
Japan	11	0.6	Discounted financial reports including Kango Life	8.4 "Tous Assureurs" as of December 31, 2021.
Hong Kong	1	7.0	Insurance Authority statistics based on gross written premiums as of December 31, 2021.	8.5 "Tous Assureurs" as of December 31, 2021.
US Insurance in the United States	16	1.8	AM Best 2021 as of December 31, 2021, in the United States in Commercial Lines.	8.6 "Tous Assureurs" as of December 31, 2021.
Switzerland	24	2.3	AM Best 2021 as of December 31, 2021.	8.7 "Tous Assureurs" as of December 31, 2021.
Thailand	18	1.0	7.2 TGA (The General Insurance Association) as of December 31, 2022 and TLA (The Life Association) as of December 31, 2022.	8.8 "Tous Assureurs" as of December 31, 2021.
Indonesia	n/a	n/a	8.7 AXA Statistics measured on Weighted New Business Premium as of September 30, 2022.	8.9 "Tous Assureurs" as of December 31, 2021.
Philippines	n/a	n/a	9.8 Insurance statistics measured on total premium income as of June 30, 2022.	9.0 "Tous Assureurs" as of December 31, 2021.
China	n/a	0.4	CBIRC (China Banking and Insurance Regulatory Commission) as of December 31, 2021.	0.5 "Tous Assureurs" as of December 31, 2021.
Mexico	5	8.0	12.3 RMS (Reservas Mexicanas de Instituciones de Seguros) as of September 30, 2022.	12.4 "Tous Assureurs" as of December 31, 2021.
Brazil	13	1.4	10.97 (Superintendencia de Seguros Privados) as of September 2020.	11.0 "Tous Assureurs" as of December 31, 2021.

# LlamaParse

A special Document Parser designed to let you build RAG over Complex docs

[https://github.com/run-llama/llama\\_parse](https://github.com/run-llama/llama_parse)

The screenshot shows the GitHub repository for LlamaParse. The README file is open, displaying the following content:

```

# llama_parse
[...]
## LlamaParse (Preview)

LlamaParse is an API created by Llamaindex to efficiently parse and represent files for efficient retrieval and context augmentation using Llamaindex frameworks.

LlamaParse directly integrates with Llamaindex.

```

The repository has 5 branches and 5 tags. It includes examples, tests, ignore, LICENSE, README.md, poetry.lock, and pyproject.toml files. The repository is public and has 59 commits. It is associated with the user logan-markwich and has 4 contributors. It is licensed under MIT and published on PyPI. The repository is part of the Llamaindex ecosystem.

# LlamaParse

## Capabilities

- ✓ Extracts tables / charts
- ✓ Input natural language parsing instructions
- ✓ JSON mode
- ✓ Image Extraction
- ✓ Support for ~10+ document types (.pdf, .pptx, .docx, .xml)

The screenshot shows the GitHub repository for LlamaParse. The README file is open, displaying the following content:

```

# llama_parse
[...]
## LlamaParse (Preview)

LlamaParse is an API created by Llamaindex to efficiently parse and represent files for efficient retrieval and context augmentation using Llamaindex frameworks.

LlamaParse directly integrates with Llamaindex.

```

The repository has 5 branches and 5 tags. It includes examples, tests, ignore, LICENSE, README.md, poetry.lock, and pyproject.toml files. The repository is public and has 59 commits. It is associated with the user logan-markwich and has 4 contributors. It is licensed under MIT and published on PyPI. The repository is part of the Llamaindex ecosystem.

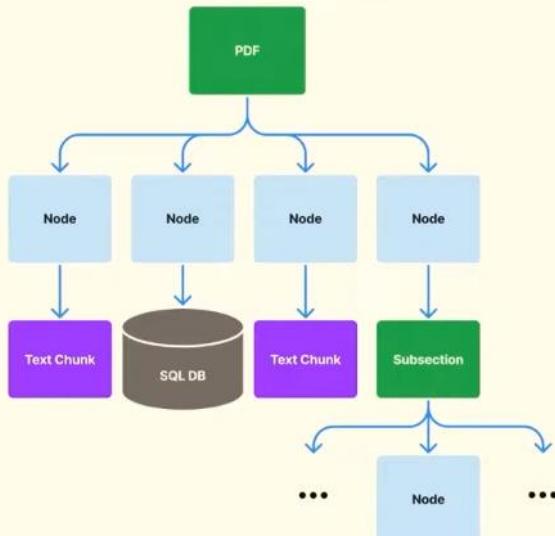
# LlamaParse Results

**Expanded:**

LlamaParse	PyPDF	PyMuPDF	Textract	PdfMiner
<img alt="Screenshot of LlamaParse output showing a detailed financial report with				

# LlamaParse + Advanced Indexing

1. Use LlamaParse to parse a document into a semi-structured markdown representation (text + tables)
  2. Use a markdown parser to extract out text and table chunks
  3. Use an LLM to extract a summary from each table → link to underlying table chunk.
  4. Index a **graph** of text and table chunks.



# Advanced Table Understanding

<pre>[54]: query = "what is the Repayments of debt in the Cash flows from financing activities for Netflix?"</pre>	response_0 = baseline_pdf_query_engine.query(query)	print("*****Baseline PDF Query Engine*****")	print(response_0)																																																												
	response_1 = raw_query_engine.query(query)	print("*****New LlamaParse+ Basic Query Engine*****")	print(response_1)																																																												
	response_2 = llmqa_query_engine.query(query)	print("*****New LlamaParse+ Recursive Retriever Query Engine*****")	print(response_2)																																																												
	*****Baseline PDF Query Engine*****	The Repayments of debt in the Cash flows from financing activities for Netflix is not provided in the given information.																																																													
	*****New LlamaParse+ Basic Query Engine*****	The repayments of debt for Netflix in the Cash flows from financing activities were \$700,000 for the year ended as of December 31, 2022, and \$500,000 for the year ended as of December 31, 2021.																																																													
	*****New LlamaParse+ Recursive Retriever Query Engine*****	The repayments of debt in the cash flows from financing activities for the year ended December 31, 2021 was \$500,000.																																																													
	<table border="1"><thead><tr><th></th><th>2022</th><th>2021</th><th>2020</th></tr></thead><tbody><tr><td><b>Cash flows from investing activities:</b></td><td></td><td></td><td></td></tr><tr><td>Purchases of property and equipment</td><td>( 407,729 )</td><td>( 524,585 )</td><td>( 497,923 )</td></tr><tr><td>Change in other assets</td><td>—</td><td>( 26,919 )</td><td>( 7,431 )</td></tr><tr><td>Acquisitions</td><td>( 757,387 )</td><td>( 788,349 )</td><td>—</td></tr><tr><td>Purchases of short-term investments</td><td>( 911,276 )</td><td>—</td><td>—</td></tr><tr><td>Net cash used in investing activities</td><td>( 2,076,392 )</td><td>( 1,339,853 )</td><td>( 505,354 )</td></tr><tr><td><b>Cash flows from financing activities:</b></td><td></td><td></td><td></td></tr><tr><td>Proceeds from issuance of debt</td><td>—</td><td>—</td><td>1,009,464</td></tr><tr><td>Debt issuance costs</td><td>—</td><td>—</td><td>( 7,559 )</td></tr><tr><td>Repayments of debt</td><td>( 700,000 )</td><td>( 500,000 )</td><td>—</td></tr><tr><td>Proceeds from issuance of common stock</td><td>35,746</td><td>174,414</td><td>235,406</td></tr><tr><td>Repurchases of common stock</td><td>—</td><td>( 600,022 )</td><td>—</td></tr><tr><td>Taxes paid related to net share settlement of equity awards</td><td>—</td><td>( 224,168 )</td><td>—</td></tr><tr><td>Net cash provided by (used in) financing activities</td><td>( 664,254 )</td><td>( 1,149,776 )</td><td>1,237,311</td></tr></tbody></table>		2022	2021	2020	<b>Cash flows from investing activities:</b>				Purchases of property and equipment	( 407,729 )	( 524,585 )	( 497,923 )	Change in other assets	—	( 26,919 )	( 7,431 )	Acquisitions	( 757,387 )	( 788,349 )	—	Purchases of short-term investments	( 911,276 )	—	—	Net cash used in investing activities	( 2,076,392 )	( 1,339,853 )	( 505,354 )	<b>Cash flows from financing activities:</b>				Proceeds from issuance of debt	—	—	1,009,464	Debt issuance costs	—	—	( 7,559 )	Repayments of debt	( 700,000 )	( 500,000 )	—	Proceeds from issuance of common stock	35,746	174,414	235,406	Repurchases of common stock	—	( 600,022 )	—	Taxes paid related to net share settlement of equity awards	—	( 224,168 )	—	Net cash provided by (used in) financing activities	( 664,254 )	( 1,149,776 )	1,237,311		
	2022	2021	2020																																																												
<b>Cash flows from investing activities:</b>																																																															
Purchases of property and equipment	( 407,729 )	( 524,585 )	( 497,923 )																																																												
Change in other assets	—	( 26,919 )	( 7,431 )																																																												
Acquisitions	( 757,387 )	( 788,349 )	—																																																												
Purchases of short-term investments	( 911,276 )	—	—																																																												
Net cash used in investing activities	( 2,076,392 )	( 1,339,853 )	( 505,354 )																																																												
<b>Cash flows from financing activities:</b>																																																															
Proceeds from issuance of debt	—	—	1,009,464																																																												
Debt issuance costs	—	—	( 7,559 )																																																												
Repayments of debt	( 700,000 )	( 500,000 )	—																																																												
Proceeds from issuance of common stock	35,746	174,414	235,406																																																												
Repurchases of common stock	—	( 600,022 )	—																																																												
Taxes paid related to net share settlement of equity awards	—	( 224,168 )	—																																																												
Net cash provided by (used in) financing activities	( 664,254 )	( 1,149,776 )	1,237,311																																																												

## Parsing Instructions

### CALCULATING THE DERIVATIVE OF A CONSTANT, LINEAR, OR QUADRATIC FUNCTION

- Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = \alpha$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(\alpha + \varepsilon) - f(\alpha)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\alpha - \alpha}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} 0 = 0$$

Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.

To parse this, we take the same instructions as before and add one sentence: Output any math equation in LATEX markdown (between \$\$). The result of parsing is clear LaTeX instructions, which render the equations perfectly:

### Calculating the Derivative of a Constant, Linear, or Quadratic Function

- Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = \alpha$  is

$$\lim_{\varepsilon \rightarrow 0} \left( \frac{f(\alpha + \varepsilon) - f(\alpha)}{\varepsilon} \right) = \lim_{\varepsilon \rightarrow 0} \left( \frac{\alpha - \alpha}{\varepsilon} \right) = \lim_{\varepsilon \rightarrow 0} 0 = 0 \text{ Thus, the derivative of } f(x) \text{ is } f'(x) = 0.$$

This makes sense, since our function is constant—the rate of change is 0.

LlamaIndex Talk (Data + AI) | Introducing LlamaParse Parsing Instructions | h4r.drawio.png - Google Drive | financial\_gpt\_databricks | research\_agent\_databricks

colab.research.google.com/drive/1dO2cwDCXj9pS9yQDZ2Vjg-0b5sRXQY/o?usp=sharing

## Introducing LlamaParse Parsing Instructions

Parsing instructions allow you to instruct our parsing model the same way you would instruct an LLM!

They can be useful to help the parser get better results on complex document layouts, to extract data in a specific format, or to transform the document in other ways.

Using Parsing Instruction you will get better results out of LlamaParse on complicated documents, and also be able to simplify your application code.

### Installation

Parsing instructions are part of the llamaParse API. They can be accessed by directly specifying the parsing\_instruction parameter in the API or by using the LlamaParse python module (which we will use for this tutorial).

To install llama-parse, just get it from PIP:

```
[ ] !pip install llama-parse
```

```
Collecting llama-parse
  Downloading llama_parse-0.3.8-py3-none-any.whl (6.7 kB)
Collecting llama-index-core>=0.10.7 (from llama-parse)
  Downloading llama_index_core-0.18.19-py3-none-any.whl (15.3 MB)
      15.3/15.3 MB 31.9 MB/s eta 0:00:00
Requirement already satisfied: PyYAML>=6.0.1 in /usr/local/lib/python3.10/dist-packages (from llama-index-core>=0.10.7->llama-parse) (6.0.1)
Requirement already satisfied: SQLAlchemy[asyncio]>=1.4.49 in /usr/local/lib/python3.10/dist-packages (from llama-index-core>=0.10.7->llama-parse) (2.0.28)
Requirement already satisfied: aiohttp<4.0.0,>=3.8.6 in /usr/local/lib/python3.10/dist-packages (from llama-index-core>=0.10.7->llama-parse) (3.9.3)
Collecting dataclasses-json (from llama-index-core>=0.10.7->llama-parse)
```

colab.research.google.com/drive/1dO2cwDCXj9pS9yQDZ2Vjg-0b5sRXQY/o?usp=sharing

## Introducing LlamaParse Parsing Instructions

### API key

The use of LlamaParse requires an API key which you can get here: <https://cloud.llamaindex.ai/pars>

```
[ ] import os
os.environ["LLAMA_CLOUD_API_KEY"] = "llx-..."
```

### Async (Notebook only)

llama-parse is async-first, so running the code in a notebook requires the use of nest\_asyncio

```
[ ] import nest_asyncio
nest_asyncio.apply()
```

### Import the package

```
[ ] from llama_parse import LlamaParse
```

### Using llamacore for getting better results (on Manga!)

Sometimes the layout of a page is unusual and you will get sub-optimal reading order results with LlamaParse. For example, when parsing manga you expect the reading order to be right to left even if the content is in English!

Introducing LlamaParse Parsing Instructions

Sometimes the layout of a page is unusual and you will get sub-optimal reading order results with LlamaParse. For example, when parsing manga you expect the reading order to be right to left even if the content is in English!

Let's download an extract of a great manga "The manga guide to calculus", by Hiroyuki Kojima (<https://www.amazon.com/Manga-Guide-Calculus-Hiroyuki-Kojima/dp/1593271948>)

```
! wget "https://drive.usercontent.google.com/uc?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download" -O ./manga.pdf
```

2024-03-13 13:57:19 -- https://drive.usercontent.google.com/uc?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download  
Resolving drive.usercontent.google.com (drive.usercontent.google.com)... 173.194.211.132, 2607:f8b0:400:c10::84  
Connecting to drive.usercontent.google.com (drive.usercontent.google.com)|173.194.211.132|:443... connected.  
HTTP request sent, awaiting response... 303 See Other  
Location: https://drive.usercontent.google.com/download?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download [following]  
--2024-03-13 13:57:19-- https://drive.usercontent.google.com/download?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download  
Reusing existing connection to drive.usercontent.google.com:443.  
HTTP request sent, awaiting response... 200 OK  
Length: 3841634 (2.9M) [application/octet-stream]  
Saving to: './manga.pdf'  
  
. ./manga.pdf 100%[=====] 2.90M --.-KB/s in 0.04s  
2024-03-13 13:57:20 (78.6 MB/s) - './manga.pdf' saved [3841634/3841634]

Without parsing instructions

Introducing LlamaParse Parsing Instructions

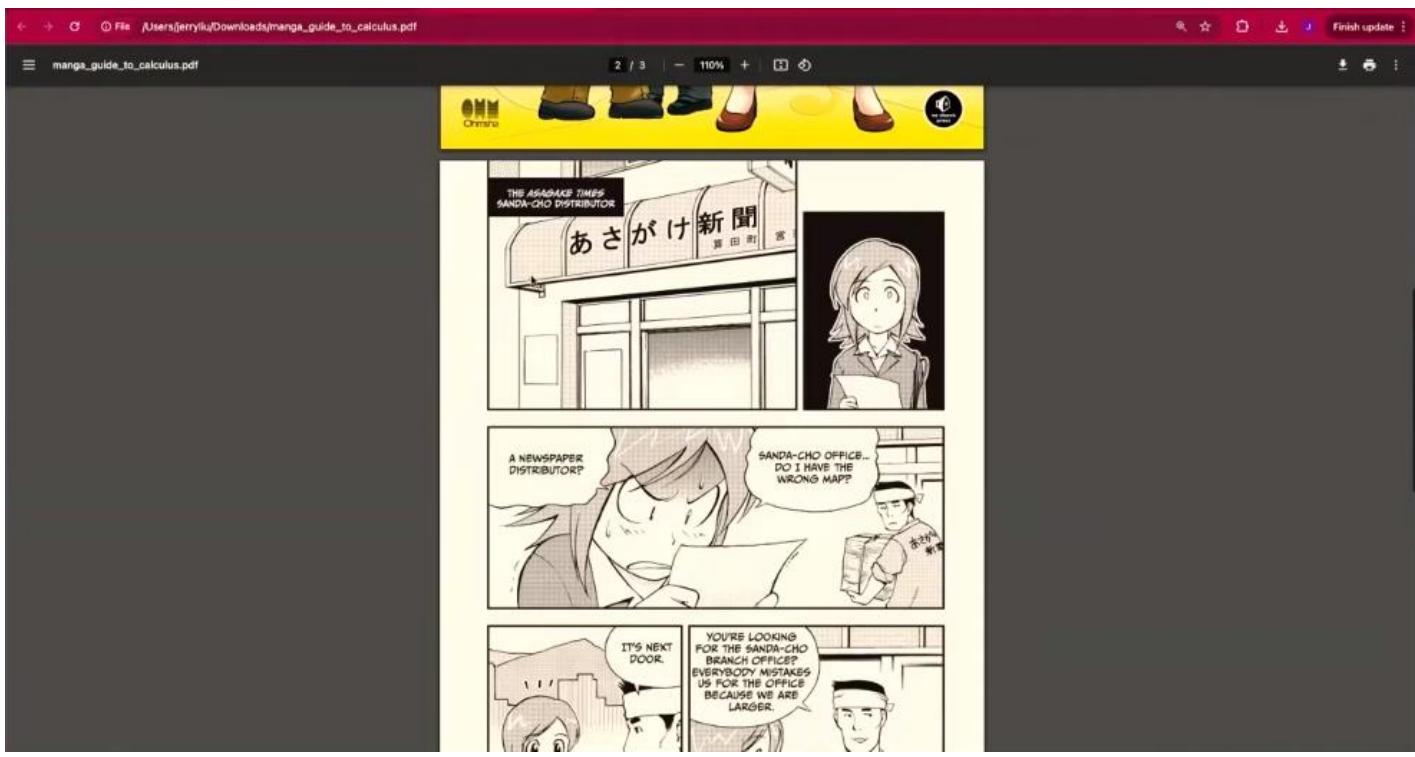
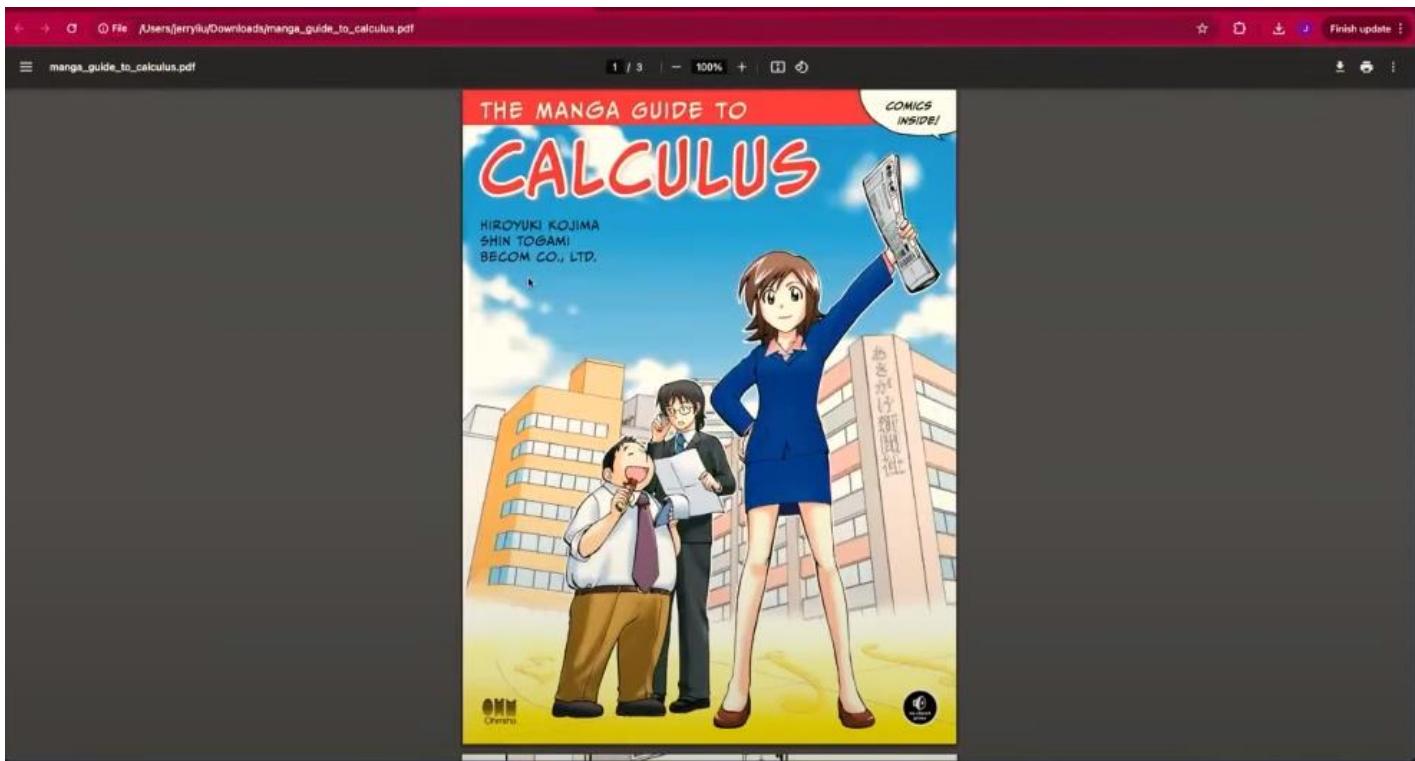
Sometimes the layout of a page is unusual and you will get sub-optimal reading order results with LlamaParse. For example, when parsing manga you expect the reading order to be right to left even if the content is in English!

Let's download an extract of a great manga "The manga guide to calculus", by Hiroyuki Kojima (<https://www.amazon.com/Manga-Guide-Calculus-Hiroyuki-Kojima/dp/1593271948>)

```
! wget "https://drive.usercontent.google.com/uc?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download" -O ./manga.pdf
```

2024-03-13 13:57:19 -- https://drive.usercontent.google.com/uc?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download  
Resolving drive.usercontent.google.com (drive.usercontent.google.com)... 173.194.211.132, 2607:f8b0:400:c10::84  
Connecting to drive.usercontent.google.com (drive.usercontent.google.com)|173.194.211.132|:443... connected.  
HTTP request sent, awaiting response... 303 See Other  
Location: https://drive.usercontent.google.com/download?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download [following]  
--2024-03-13 13:57:19-- https://drive.usercontent.google.com/download?id=1tZJhcpeplRdQFJFCFX500IqLyLgqzZsY&export=download  
Reusing existing connection to drive.usercontent.google.com:443.  
HTTP request sent, awaiting response... 200 OK  
Length: 3841634 (2.9M) [application/octet-stream]  
Saving to: './manga.pdf'  
  
. ./manga.pdf 100%[=====] 2.90M --.-KB/s in 0.04s  
2024-03-13 13:57:20 (78.6 MB/s) - './manga.pdf' saved [3841634/3841634]

Without parsing instructions



File /Users/jerryliu/Downloads/manga\_guide\_to\_calculus.pdf

3 / 3 | - 110% + [undo] [redo]

### CALCULATING THE DERIVATIVE OF A CONSTANT, LINEAR, OR QUADRATIC FUNCTION

1. Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = \alpha$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(\alpha + \varepsilon) - f(\alpha)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\alpha - \alpha}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} 0 = 0$$

Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.

NOTE: The differential coefficient of  $f(x)$  at  $x = \alpha$  is often simply called the derivative of  $f(x)$  at  $x = \alpha$ , or just  $f'(\alpha)$ .

2. Let's calculate the derivative of linear function  $f(x) = \alpha x + \beta$ . The derivative of  $f(x)$  at  $x = \alpha$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(\alpha + \varepsilon) - f(\alpha)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\alpha(\alpha + \varepsilon) + \beta - (\alpha\alpha + \beta)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \alpha = \alpha$$

Thus, the derivative of  $f(x)$  is  $f'(x) = \alpha$ , a constant value. This result should also be intuitive—linear functions have a constant rate of change by definition.

3. Let's find the derivative of  $f(x) = x^2$ , which appeared in the story. The differential coefficient of  $f(x)$  at  $x = \alpha$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(\alpha + \varepsilon) - f(\alpha)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{(\alpha + \varepsilon)^2 - \alpha^2}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{2\alpha\varepsilon + \varepsilon^2}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} (2\alpha + \varepsilon) = 2\alpha$$

Thus, the differential coefficient of  $f(x)$  at  $x = \alpha$  is  $2\alpha$ , or  $f'(\alpha) = 2\alpha$ . Therefore, the derivative of  $f(x)$  is  $f'(x) = 2x$ .

### SUMMARY

- The calculation of a limit that appears in calculus is simply a formula calculating an error.
- A limit is used to obtain a derivative.
- The derivative is the slope of the tangent line at a given point.
- The derivative is nothing but the rate of change.

CALCULATING THE DERIVATIVE OF A CONSTANT, LINEAR, OR QUADRATIC FUNCTION

i. Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = a$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(a + \varepsilon) - f(a)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\alpha - \alpha}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} 0 = 0$$

Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.

NOTE The differential coefficient of  $f(x)$  at  $x = a$  is often simply called the derivative of  $f(x)$  at  $x = a$ , or just  $f'(a)$ .

ii. Let's calculate the derivative of linear function  $f(x) = \alpha x + \beta$ . The derivative of  $f(x)$  at  $x = a$  is

$$\lim_{\varepsilon \rightarrow 0} \frac{f(a + \varepsilon) - f(a)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\alpha(a + \varepsilon) + \beta - (\alpha a + \beta)}{\varepsilon} = \lim_{\varepsilon \rightarrow 0} \alpha = \alpha$$

Thus, the derivative of  $f(x)$  is  $f'(x) = \alpha$ , a constant value. This result should also be intuitive—linear functions have a constant rate of change

Introducing LLamaParse Parsing Instructions

File Edit View Insert Runtime Tools Help Changes will not be saved

+ Code + Text Copy to Drive Connect Gemini

As you can see below, LlamaParse is not doing a great job here. It is interpreting the grid of comic panels as a table, and trying to fit the dialogue into a table. It's very hard to follow.

```
print(vanillaParsing[0].text[100:1000])
```

The Asagake Times Sanda-Chō Distributor

A newspaper distributor? do I have the wrong map?

You're looking It's next for the Sanda-chō door. branch office? Everybody mistakes us for the office because we are larger. What Is a Function? 3

## Calculating the Derivative of a Constant, Linear, or Quadratic Function

| 1. | Let's find the derivative of constant function  $f(x) = a$ . The differential coefficient of  $f(x)$  at  $x = a$  is|  
| ---|  
| |  $\lim_{\epsilon \rightarrow 0} (f(a + \epsilon) - f(a)) / \epsilon = \lim_{\epsilon \rightarrow 0} (\alpha - a) = \lim_{\epsilon \rightarrow 0} 0 = 0$ |  
| | Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.|

Note: The differential coefficient of  $f(x)$  at  $x = a$  is often simply called the derivative of  $f(x)$  at  $x = a$ , or just  $f'(a)$ .

| 2. | Let's calculate the derivative of linear function  $f(x) = ax + \beta$ . The derivative of  $f(x)$  at  $x = a$  is|  
| ---|  
| |  $\lim_{\epsilon \rightarrow 0} (f(a + \epsilon) - f(a)) =$

Using parsing instructions

Let's try to parse the manga with custom instructions:

The provided document is a manga comic book. Most pages do NOT have title. It does not contain tables. Try to reconstruct the dialogue happening in a cohesive way.

To do so just pass the parsing instruction as a parameter to LlamaParse:

Introducing LLamaParse Parsing Instructions

File Edit View Insert Runtime Tools Help Changes will not be saved

+ Code + Text Copy to Drive Connect Gemini

Note: The differential coefficient of  $f(x)$  at  $x = a$  is often simply called the derivative of  $f(x)$  at  $x = a$ , or just  $f'(a)$ .

| 2. | Let's calculate the derivative of linear function  $f(x) = ax + \beta$ . The derivative of  $f(x)$  at  $x = a$  is|  
| ---|  
| |  $\lim_{\epsilon \rightarrow 0} (f(a + \epsilon) - f(a)) =$

Using parsing instructions

Let's try to parse the manga with custom instructions:

The provided document is a manga comic book. Most pages do NOT have title. It does not contain tables. Try to reconstruct the dialogue happening in a cohesive way.

To do so just pass the parsing instruction as a parameter to LlamaParse:

```
parsingInstructionManga = """The provided document is a manga comic book, most page do NOT have title.  
It does not contain table.  
Try to reconstruct the dialog happening in a cohesive way."""  
withInstructionParsing = LlamaParse(result_type="markdown", parsing_instruction=parsingInstructionManga).load_data("./manga.pdf")
```

Started parsing the file under job\_id 88ab273e-b2a7-4f84-8e72-e9367cf6b114

Let's see how it compare with page 3! We encourage you to play with the target name and explore other names. As you will see, the parsing

Introducing LLamaParse Parsing Instructions

The derivative is the slope of the tangent line at a given point.  
The derivative is nothing but the rate of change.

40 CHAPTER 1 LET'S DIFFERENTIATE A FUNCTION!

```
parsingInstructionMangaLatex = """The provided document is a manga comic book, most page do NOT have title.  
It does not contain table. Do not output table.  
Try to reconstruct the dialog happening in a cohesive way.  
Output any math equation in LATEX markdown (between $$)"""  
withLatex = LlamaParse(result_type="markdown", parsing_instruction=parsingInstructionMangaLatex).load_data("./manga.pdf")  
  
Started parsing the file under job_id 3a055e64-d91e-484e-b9b0-99a2e637c08d  
  
target_page=2  
print("\n\n[Without instruction]-----\n")  
print(vanillaParsing[0].text.split('\n---\n')[target_page])  
print("\n\n[With instruction to output math in LATEX!]-----\n")  
print(withLatex[0].text.split('\n---\n')[target_page])
```

Introducing LLamaParse Parsing Instructions

```
target_page=2  
print("\n\n[Without instruction]-----\n")  
print(vanillaParsing[0].text.split('\n---\n')[target_page])  
print("\n\n[With instruction to output math in LATEX!]-----\n")  
print(withLatex[0].text.split('\n---\n')[target_page])  
  
1.|Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = a$  is|  
|---|---|  
| | $\lim_{\epsilon \rightarrow 0} (f(a + \epsilon) - f(a)) / \epsilon = \lim_{\epsilon \rightarrow 0} (\alpha - \alpha) = \lim_{\epsilon \rightarrow 0} 0 = 0$ |  
| |Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.|  
Note: The differential coefficient of  $f(x)$  at  $x = a$  is often simply called the derivative of  $f(x)$  at  $x = a$ , or just  $f'(a)$ .  
2.|Let's calculate the derivative of linear function  $f(x) = \alpha x + \beta$ . The derivative of  $f(x)$  at  $x = a$  is|  
|---|---|  
| | $\lim_{\epsilon \rightarrow 0} (f(a + \epsilon) - f(a)) / \epsilon = \lim_{\epsilon \rightarrow 0} (\alpha(a + \epsilon) + \beta - (\alpha a + \beta)) = \lim_{\epsilon \rightarrow 0} \alpha = \alpha$ |  
| |Thus, the derivative of  $f(x)$  is  $f'(x) = \alpha$ , a constant value. This result should also be intuitive—linear functions have a constant rate of change.|  
3.|Let's find the derivative of  $f(x) = x^2$ , which appeared in the story. The differential coefficient of  $f(x)$  at  $x = a$  is|  
|---|---|  
| | $\lim_{\epsilon \rightarrow 0} ((a + \epsilon)^2 - a^2) / \epsilon = \lim_{\epsilon \rightarrow 0} (a^2 + 2a\epsilon + \epsilon^2 - a^2) / \epsilon = \lim_{\epsilon \rightarrow 0} (2a\epsilon + \epsilon^2) = \lim_{\epsilon \rightarrow 0} 2a = 2a$ |  
| |Thus, the differential coefficient of  $f(x)$  at  $x = a$  is  $2a$ , or  $f'(a) = 2a$ . Therefore, the derivative of  $f(x)$  is  $f'(x) = 2x$ .|  
## Summary
```

Introducing LLamaParse Parsing Instructions

+ Code + Text | Copy to Drive

[With instruction to output math in LATEX!]

```
# Derivative of Constant, Linear, or Quadratic Function
## Calculating the Derivative of a Constant, Linear, or Quadratic Function

1. Let's find the derivative of constant function  $f(x) = \alpha$ . The differential coefficient of  $f(x)$  at  $x = a$  is

$$
\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} \left( \frac{f(a + \varepsilon) - f(a)}{\varepsilon} \right) = \lim_{\varepsilon \rightarrow 0} \frac{\alpha - \alpha}{\varepsilon} = 0
\end{aligned}

Thus, the derivative of  $f(x)$  is  $f'(x) = 0$ . This makes sense, since our function is constant—the rate of change is 0.

Note: The differential coefficient of  $f(x)$  at  $x = a$  is often simply called the derivative of  $f(x)$  at  $x = a$ , or just  $f'(a)$ .
```

2. Let's calculate the derivative of linear function  $f(x) = ax + \beta$ . The derivative of  $f(x)$  at  $x = a$  is

```
$$
\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} \left( \frac{f(a + \varepsilon) - f(a)}{\varepsilon} \right) = \lim_{\varepsilon \rightarrow 0} \frac{a(a + \varepsilon) + \beta - (a\alpha + \beta)}{\varepsilon} = a
\end{aligned}

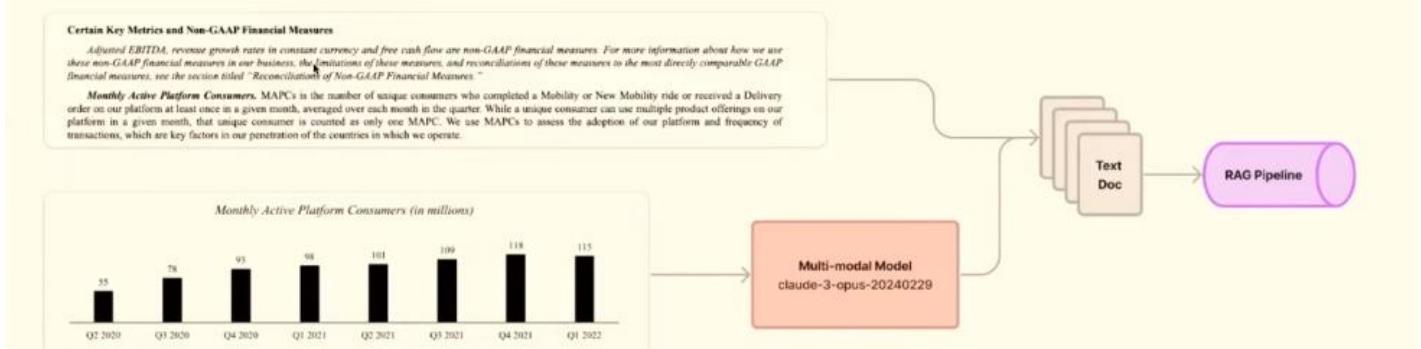
Thus, the derivative of  $f(x)$  is  $f'(x) = a$ , a constant value. This result should also be intuitive—linear functions have a constant rate of change.
```

3. Let's find the derivative of  $f(x) = x^2$ . The differential coefficient of  $f(x)$  at  $x = a$  is

Using a prompt, we now get back the equations in Latex as above

## JSON Mode

[https://github.com/run-llama/llama\\_parse/blob/main/examples/demo\\_json.ipynb](https://github.com/run-llama/llama_parse/blob/main/examples/demo_json.ipynb)



# RAG over Powerpoints

[https://github.com/run-llama/llama\\_parse/blob/main/examples/other\\_files/demo\\_ppt\\_financial.ipynb](https://github.com/run-llama/llama_parse/blob/main/examples/other_files/demo_ppt_financial.ipynb)

DATA4AI SUMMIT

```
[27]: print(llama_parse_documents[0].get_content()[-2800:-2300])
```

```
ation and mitigation
_____
|Item|31 Dec 2022|31 Dec 2021|Change|
|---|---|---|
|Payables and accruals|4,685|4,066|619|
|Employee benefits|127,215|84,676|42,539|
|Contributions received in advance|6,975|10,192|(3,217)|
|Unearned revenue from exchange transactions|20|651|(631)|
|Deferred Revenue|71,301|55,737|15,564|
|Borrowings|28,229|29,082|(773)|
|Funds held in trust|30,373|29,014|1,359|
|Provisions|1,706|1,910|(204)|
|Total Liabilities|270,504|215,248|55,256|
_____
## Liabilities
```

Employee Ben

Compared against the original slide image.

## Liabilities

(in '000's of CHF)

Item	31 Dec 2022	31 Dec 2021	Change
Payables and accruals	4,685	4,066	619
Employee benefits	127,215	84,676	42,539
Contributions received in advance	6,975	10,192	(3,217)
Unearned revenue from exchange transactions	20	651	(631)
Deferred Revenue	71,301	55,737	15,564
Borrowings	28,229	29,002	(773)
Funds held in trust	30,373	29,014	1,359
Provisions	1,706	1,910	(204)
<b>Total Liabilities</b>	<b>270,504</b>	<b>215,248</b>	<b>55,256</b>

©2024 Databricks Inc. — All rights reserved



31

# Workshop

Let's build a RAG pipeline with Databricks LLMs + local embeddings

<https://colab.research.google.com/drive/1iB3HPOZPpY4qMFpw1Dtr1rJDb44k1TL2?usp=sharing>



# Improving Query Complexity

# Complex Questions

There's certain questions we want to ask where naive RAG will fail.

Examples:

- **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"
- **Structured Analytics + Semantic Search:** "Tell me about the risk factors of the highest-performing rideshare company in the US"
- **General 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."

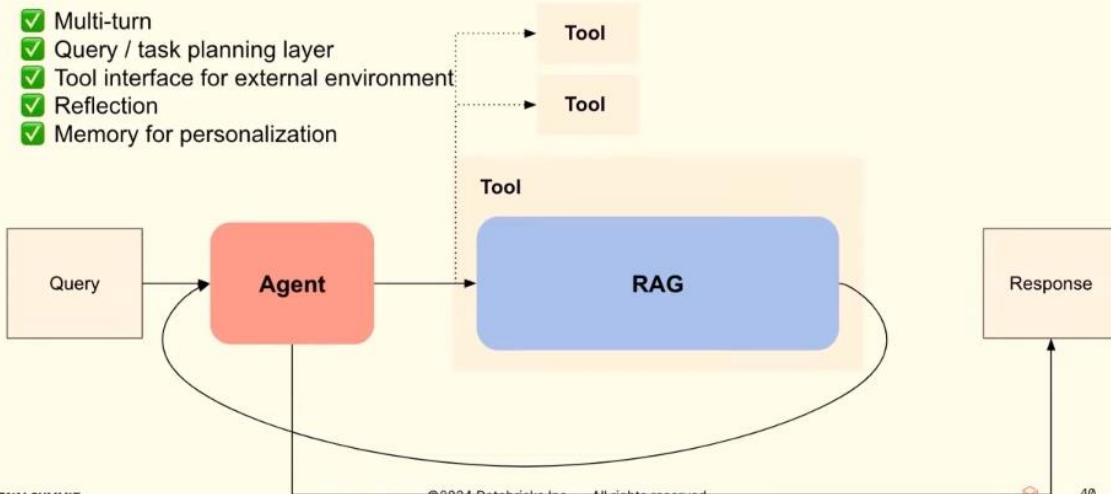
## From RAG to Agents



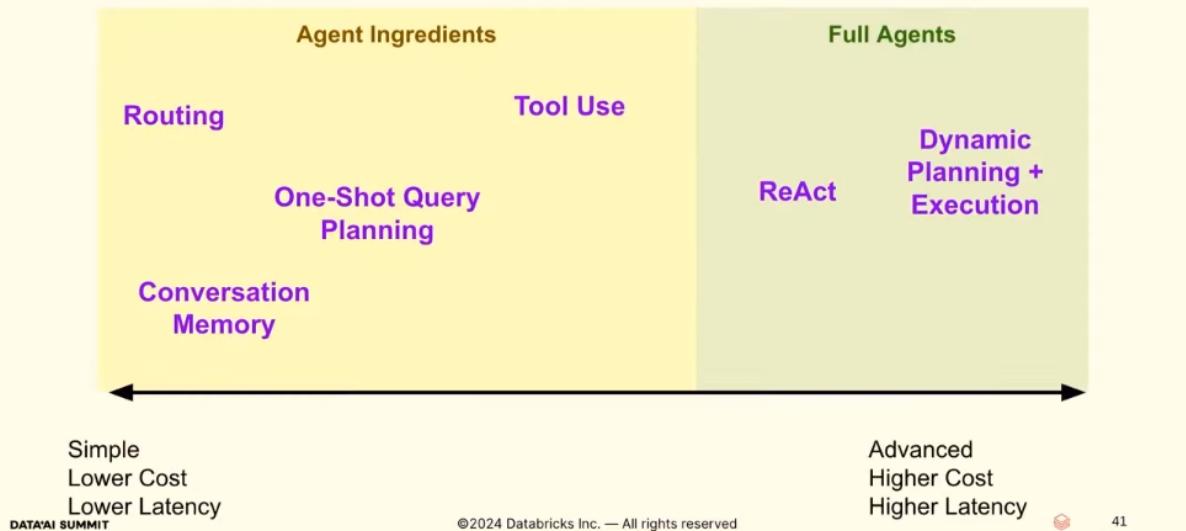
- ⚠ Single-shot
- ⚠ No query understanding/planning
- ⚠ No tool use
- ⚠ No reflection, error correction
- ⚠ No memory (stateless)

## From RAG to Agents

- ✓ Multi-turn
- ✓ Query / task planning layer
- ✓ Tool interface for external environment
- ✓ Reflection
- ✓ Memory for personalization



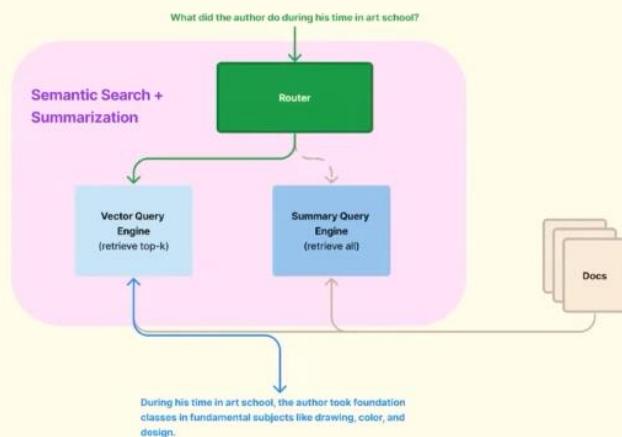
# From Simple to Advanced Agents



## Routing

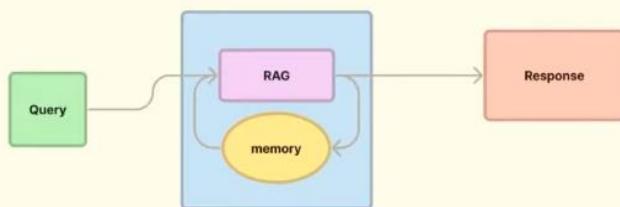
Simplest form of agentic reasoning.

Given user query and set of choices, output subset of choices to route query to.



## Conversation Memory

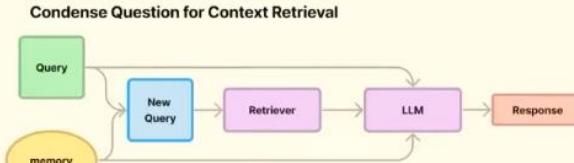
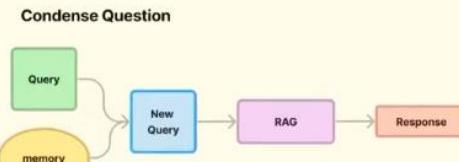
In addition to current query, take into account **conversation history** as input to your RAG pipeline.



# Conversation Memory

How to account for conversation history in a RAG pipeline?

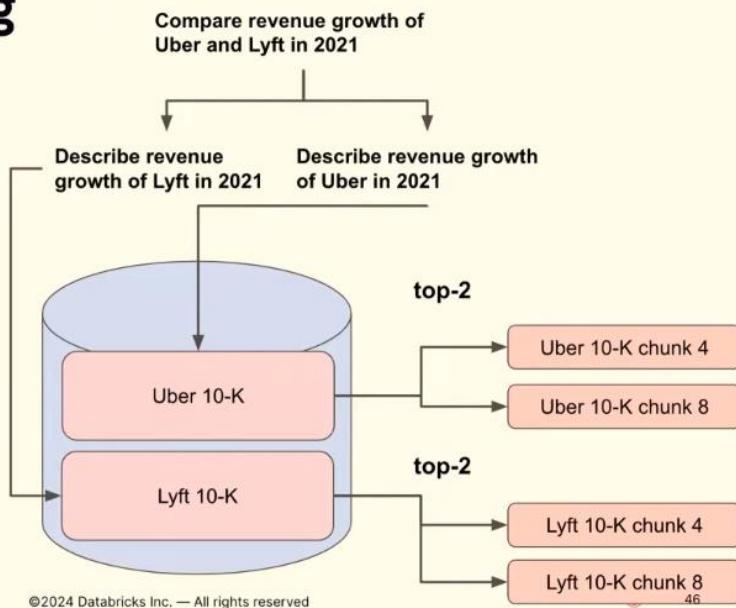
- Condense question
- Condense question + context



# Query Planning

Break down query into parallelizable sub-queries.

Each sub-query can be executed against any set of RAG pipelines



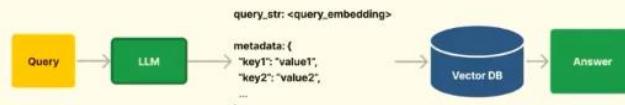
DATAVI SUMMIT

# Tool Use

Use an LLM to call an API

Infer the parameters of that API

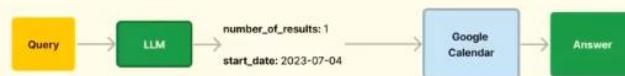
## Auto-Retrieval



## Text-to-SQL



## Calendar

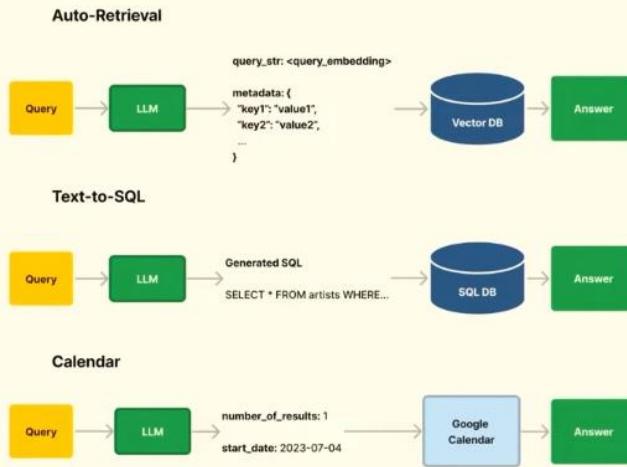


# Tool Use

In normal RAG you just pass through the query.

But what if you used the LLM to infer all the parameters for the API interface?

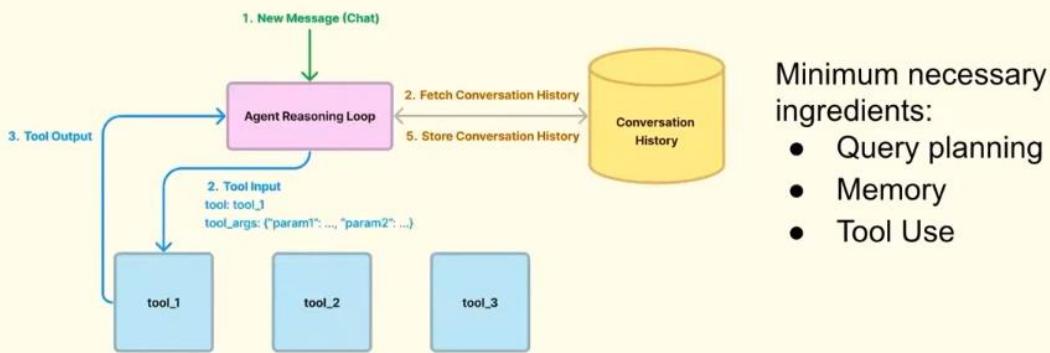
A key capability in many QA use cases (auto-retrieval, text-to-SQL, and more)



## Let's put them together

- All of these are **agent ingredients**
- Let's put them together for a full agent system
  - Query planning
  - Memory
  - Tool Use
- Let's add additional components
  - Reflection
  - Controllability
  - Observability

## Core Components of a Full Agent



# Agent Reasoning Loops

**Sequential:** Generate next step given previous steps (chain-of-thought prompt)

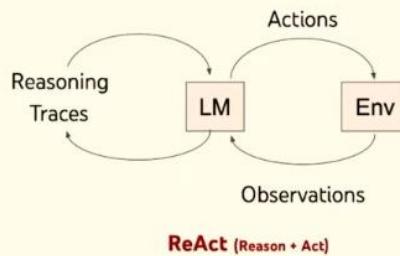
**DAG-based planning (deterministic):** Generate a deterministic DAG of steps. Replan if steps don't reach desired state.

**Tree-based planning (stochastic):** Sample multiple future states at each step. Run Monte-Carlo Tree Search (MCTS) to balance exploration vs. exploitation.

## Agent Reasoning: Sequential

**ReAct:** Chain-of-thought and tool use through prompting.

**Function Calling Loop:** Call LLM Function Calling APIs in a loop until done.



[ReAct + RAG Guide](#)

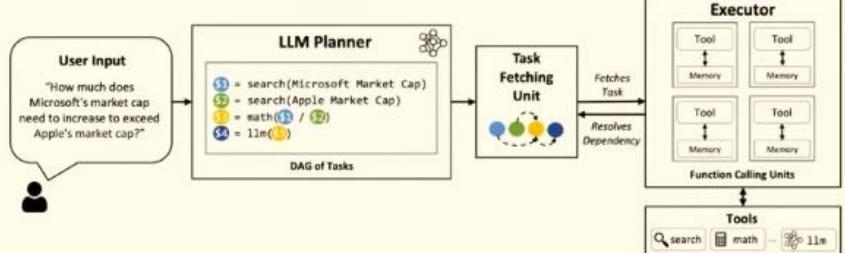
[Function Calling Anthropic Agent](#)

## Agent Reasoning: DAG-based Planning

**LLM Compiler (Kim et al. 2023):** An agent compiler for parallel multi-function planning + execution.

[LLMCompiler Agent](#)

[Structured Planner Agent](#)



# Agent Reasoning: Tree-based Planning

Tree of Thoughts (Yao et al. 2023)

Reasoning via Planning (Hao et al. 2023)

Language Agent Tree Search (Zhou et al. 2023)

LATS Guide

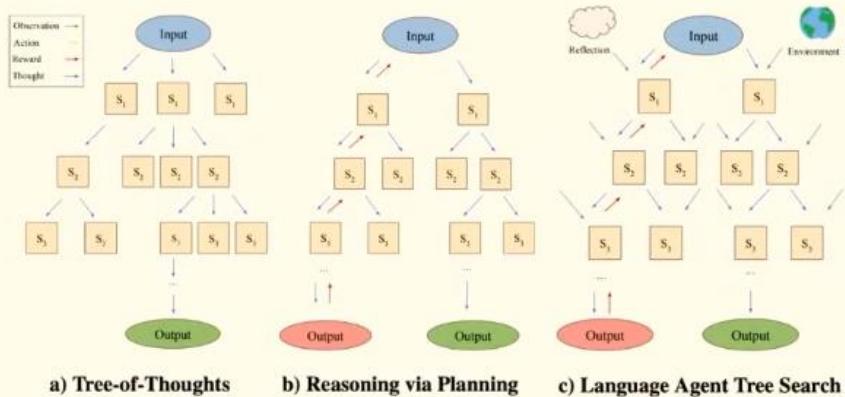


Figure 2: An overview of the differences between LATS and recently proposed LM search algorithms ToT (Yao et al., 2023a) and RAP (Hao et al., 2023). LATS leverages environmental feedback and self-reflection to further adapt search and improve performance.

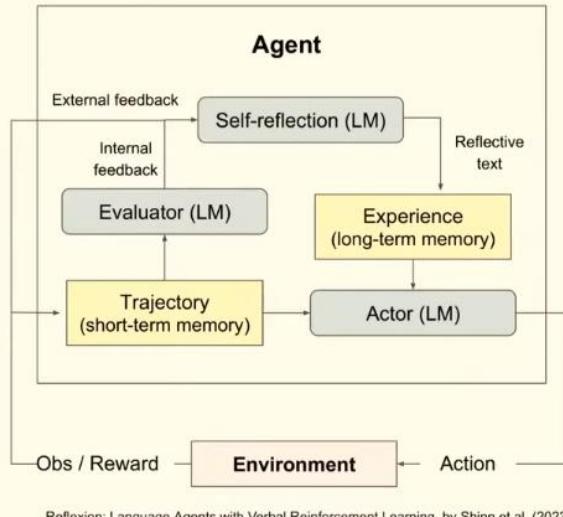
## Self-Reflection

Use feedback to improve agent execution and reduce errors

💡 Human feedback

🌐 LLM feedback

**Use few-shot examples instead of retraining the model!**



## Additional Requirements

- **Observability:** see the full trace of the agent
  - [Observability Guide](#)
- **Control:** Be able to guide the intermediate steps of an agent *step-by-step*
  - [Lower-Level Agent API](#)
- **Customizability:** Define your own agentic logic around any set of tools.
  - [Custom Agent Guide](#)
  - [Custom Agent with Query Pipeline Guide](#)
- **Multi-agents:** Define multi-agent interactions!
  - Synchronously: Define an explicit flow between agents
  - Asynchronously: Treat each agent as a microservice that can communicate with each other.
    - Upcoming in LlamaIndex!
  - Current Frameworks: Autogen, CrewAI

# LlamaCloud for Advanced RAG

Building RAG/agents in an enterprise setting? We'd love to chat!

## Workshop

Let's extend our RAG pipeline into an agent!

<https://colab.research.google.com/drive/18RUkf8IpHVSJF-rDh8cOj0QJ6UwQonfh?usp=sharing>

The screenshot shows a Jupyter Notebook interface with the title "research\_agent\_databricks". The notebook contains a section titled "Building a Research Agent with Databricks" which includes a note about reasoning over tools in multiple steps and using a FunctionCallingAgent implementation. Below this, there is a "Setup" section with a code cell containing pip installation commands:

```
[1] !pip install llama-index==0.10.28
!pip install llama-index-lms-databricks
!pip install llama-index-embeddings-huggingface
!pip install llama-parse
```

After running the command, the terminal output shows the download and collection of various packages, including "llama-index-agent-openai", "llama\_index\_agent\_openai", "llama-index-embeddings-openai", and "llama\_index\_indices\_managed\_llama\_cloud". The total download time is 2m 4s, completed at 3:54PM.

colab.research.google.com/drive/18RUkfBipIVSJF-rDh8c0j0QJ6UwQorfh?usp=sharing#scrollTo=dC2zbATkgM

Comment Share Gemini J

research\_agent\_databricks

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Requirement already satisfied: typing<3.0.0,>=4.00.1 in /usr/local/lib/python3.10/dist-packages (from llama-index-core<0.11.0,>=0.10.28)

[1] Requirement already satisfied: typing-extentions>=4.5.0 in /usr/local/lib/python3.10/dist-packages (from llama-index-core<0.11.0,>=0.10.28)

Collecting typing-inspect>=0.8.0 (from llama-index-core<0.11.0,>=0.10.28->llama-index==0.10.28)

  ↳ Downloading typing\_inspect-0.9.0-py3-none-any.whl (8.8 kB)

Requirement already satisfied: wrapt in /usr/local/lib/python3.10/dist-packages (from llama-index-core<0.11.0,>=0.10.28)

[2m] import os  
os.environ["LLAMA\_CLOUD\_API\_KEY"] = "llx..."

[2m] import nest\_asyncio  
nest\_asyncio.apply()

[ ] # databricks api key  
api\_key = "<api\_key>"

[ ] from llama\_index.llms.databricks import Databricks  
from llama\_index.embeddings.huggingface import HuggingFaceEmbedding  
from llama\_index.core import Settings

embed\_model = HuggingFaceEmbedding(model\_name="BAAI/bge-small-en-v1.5")  
llm = Databricks(  
 model="databricks-meta-llama-3-70b-instruct",  
 api\_key=api\_key,  
 api\_base="https://<cluster\_id>.cloud.databricks.com/serving-endpoints",  
)

2m 4s completed at 3:54PM

colab.research.google.com/drive/18RUkfBipIVSJF-rDh8c0j0QJ6UwQorfh?usp=sharing#scrollTo=dC2zbATkgN

Comment Share Gemini J

research\_agent\_databricks

File Edit View Insert Runtime Tools Help

+ Code + Text

COLLECTING typing-inspect>=0.8.0 (from llama-index-core<0.11.0,>=0.10.28->llama-index==0.10.28)

  ↳ Downloading typing\_inspect-0.9.0-py3-none-any.whl (8.8 kB)

Requirement already satisfied: wrapt in /usr/local/lib/python3.10/dist-packages (from llama-index-core<0.11.0,>=0.10.28->llama-index==0.10.28)

[1] import os  
os.environ["LLAMA\_CLOUD\_API\_KEY"] = "llx..."

[2] import nest\_asyncio  
nest\_asyncio.apply()

[ ] # databricks api key  
api\_key = "<api58d99aab09a284792217973389528f86"

[ ] from llama\_index.llms.databricks import Databricks  
from llama\_index.embeddings.huggingface import HuggingFaceEmbedding  
from llama\_index.core import Settings

embed\_model = HuggingFaceEmbedding(model\_name="BAAI/bge-small-en-v1.5")  
llm = Databricks(  
 model="databricks-meta-llama-3-70b-instruct",  
 api\_key=api\_key,  
 api\_base="https://<cluster\_id>.cloud.databricks.com/serving-endpoints",  
)

Settings.llm = llm  
Settings.embed\_model = embed\_model

0s completed at 5:13PM

colab.research.google.com/drive/1B8Ukf8ipIV5JF-Drh8cOjOQJ6UwQonfh?usp=sharing#scrollTo=tM57oGkxTKgN

research\_agent\_databricks ☆

File Edit View Insert Runtime Tools Help

+ Code + Text

[ ] import os  
os.environ["LLAMA\_CLOUD\_API\_KEY"] = "llx-..."

{x} 0s [2] import nest\_asyncio  
nest\_asyncio.apply()

0s [3] # databricks api key  
api\_key = "dap158d99aab09a284792217973389528f86"

from llama\_index.llms.databricks import Databricks  
from llama\_index.embeddings.huggingface import HuggingFaceEmbedding  
from llama\_index.core import Settings

embed\_model = HuggingFaceEmbedding(model\_name="BAII/bge-small-en-v1.5")  
llm = Databricks(  
 model="databricks-meta-llama-3-70b-instruct",  
 api\_key=api\_key,  
 api\_base="https://<cluster\_id>.cloud.databricks.com/serving-endpoints",  
)

Settings.llm = llm  
Settings.embed\_model = embed\_model

0s completed at 5:13 PM

colab.research.google.com/drive/1B8Ukf8ipIV5JF-Drh8cOjOQJ6UwQonfh?usp=sharing#scrollTo=tM57oGkxTKgN

research\_agent\_databricks ☆

File Edit View Insert Runtime Tools Help

+ Code + Text

0s [3] # databricks api key  
api\_key = "dap158d99aab09a284792217973389528f86"

0s [ ] from llama\_index.llms.databricks import Databricks  
from llama\_index.embeddings.huggingface import HuggingFaceEmbedding  
from llama\_index.core import Settings

embed\_model = HuggingFaceEmbedding(model\_name="BAII/bge-small-en-v1.5")  
llm = Databricks(  
 model="databricks-meta-llama-3-70b-instruct",  
 api\_key=api\_key,  
 api\_base="https://<cluster\_id>.cloud.databricks.com/serving-endpoints",  
)

Settings.llm = llm  
Settings.embed\_model = embed\_model

...

colab.research.google.com/drive/1B8Ukf8ipIV5JF-Drh8cOjOQJ6UwQonfh?usp=sharing#scrollTo=tM57oGkxTKgN

research\_agent\_databricks ☆

File Edit View Insert Runtime Tools Help

+ Code + Text

Please note that authentication is recommended but still optional to access public models or datasets.  
warnings.warn(...)

modules.json: 100% 349/349 [00:00<00:00, 6.89kB/s]

config\_sentence\_transformers.json: 100% 124/124 [00:00<00:00, 2.43kB/s]

README.md: 100% 94.8k/94.8k [00:00<00:00, 1.28MB/s]

sentence\_bert\_config.json: 100% 52.0/52.0 [00:00<00:00, 554B/s]

/usr/local/lib/python3.10/dist-packages/huggingface\_hub/file\_download.py:1132: FutureWarning: `resume\_download` is deprecated and will be  
warnings.warn(...)

config.json: 100% 743/743 [00:00<00:00, 12.1kB/s]

research\_agent\_databricks

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk Gemini

Download ~3 ICLR 2024 papers, use LlamaParse

{x} Let's parse 3 ICLR 2024 research papers using LlamaParse.

```
urls = [
    "https://openreview.net/pdf?id=VtmBAGCN7o",
    "https://openreview.net/pdf?id=6PmJoRfdak",
    "https://openreview.net/pdf?id=hSyW5go0v8",
]

papers = [
    "metappt.pdf",
    "longlora.pdf",
    "selfrag.pdf",
]

[] for url, paper in zip(urls, papers):
    !wget "{url}" -O "{paper}"

--2024-06-11 06:04:43-- https://openreview.net/pdf?id=VtmBAGCN7o
Resolving openreview.net (openreview.net)... 35.184.86.251
Connecting to openreview.net (openreview.net)|35.184.86.251|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 16911937 (16M) [application/pdf]
```

0s completed at 5:14 PM

research\_agent\_databricks

File Edit View Insert Runtime Tools Help

+ Code + Text

RAM Disk Gemini

```
2024-06-12 00:14:23 (4.48 MB/s) - 'selfrag.pdf' saved [1244749/1244749]
```

{x} [8] from llama\_parse import LlamaParse

```
from llama_index.core.schema import Document
def _load_data(file_path: str) -> Document:
    parser = LlamaParse(result_types="text")
    json_objs = parser.get_json_result(file_path)
    json_list = json_objs[0]["pages"]
    docs = []
    for item in json_list:
        doc = Document(
            text=item["text"], metadata={"page_label": item["page"]})
        docs.append(doc)
    return docs
```

Convert papers to Tools

```
[ ] # TODO: abstract all of this into a function that takes in a PDF file name
from llama_index.core import VectorStoreIndex, SummaryIndex
from llama_parse import LlamaParse
```

0s completed at 5:14 PM

colab.research.google.com/drive/18RUkFBiphIVSJF-0h8c0j0QJ6UwQnfh?usp=sharing#scrollTo=u3w2YtPcTKgP

Comment Share ⚙️ J

+ Code + Text

Convert papers to Tools

```
▶ # TODO: abstract all of this into a function that takes in a PDF file name

from llama_index.core import VectorStoreIndex, SummaryIndex
from llama_parse import LlamaParse
from llama_index.core.node_parser import SentenceSplitter
from llama_index.core.tools import FunctionTool, QueryEngineTool
from llama_index.core.vector_stores import MetadataFilters, FilterCondition
from typing import List, Optional

def get_doc_tools(
    file_path: str,
    name: str,
) -> str:
    """Get vector query and summary query tools from a document."""

    # load documents
    # documents = LlamaParse(result_type="text").load_data(file_path)
    documents = _load_data(file_path)
    splitter = SentenceSplitter(chunk_size=1024)
    nodes = splitter.get_nodes_from_documents(documents)
    vector_index = VectorStoreIndex(nodes)

    def vector_query(
        query: str,
        page_numbers: Optional[List[int]] = None
    ) -> str:
        """Use to answer questions over a given paper.

        Useful if you have specific questions over the paper.
        Always leave page_numbers as None UNLESS there is a specific page you want to search for.

        Args:
            query (str): the string query to be embedded.
            page_numbers (Optional[List[int]]): Filter by set of pages. Leave as NONE
                if we want to perform a vector search
                over all pages. Otherwise, filter by the set of specified pages.

```

✓ 0s completed at 5:14 PM

colab.research.google.com/drive/18RUkFBiphIVSJF-0h8c0j0QJ6UwQnfh?usp=sharing#scrollTo=m6XU4gBtTKgP

Comment Share ⚙️ J

+ Code + Text

All changes saved

```
▶ def get_doc_tools(
    file_path: str,
    name: str,
) -> str:
    """Get vector query and summary query tools from a document."""

    # load documents
    # documents = LlamaParse(result_type="text").load_data(file_path)
    documents = _load_data(file_path)
    splitter = SentenceSplitter(chunk_size=1024)
    nodes = splitter.get_nodes_from_documents(documents)
    vector_index = VectorStoreIndex(nodes)

    def vector_query(
        query: str,
        page_numbers: Optional[List[int]] = None
    ) -> str:
        """Use to answer questions over a given paper.

        Useful if you have specific questions over the paper.
        Always leave page_numbers as None UNLESS there is a specific page you want to search for.

        Args:
            query (str): the string query to be embedded.
            page_numbers (Optional[List[int]]): Filter by set of pages. Leave as NONE
                if we want to perform a vector search
                over all pages. Otherwise, filter by the set of specified pages.

```

✓ 0s completed at 5:14 PM

research\_agent\_databricks

```
name: str,
) -> str:
    """Get vector query and summary query tools from a document."""

    # load documents
    # documents = LlamaParse(result_type="text").load_data(file_path)
    documents = _load_data(file_path)
    splitter = SentenceSplitter(chunk_size=1024)
    nodes = splitter.get_nodes_from_documents(documents)
    vector_index = VectorStoreIndex(nodes)

    def vector_query(
        query: str,
        page_numbers: Optional[List[int]] = None
    ) -> str:
        """Use to answer questions over a given paper.

        Useful if you have specific questions over the paper.
        Always leave page_numbers as None UNLESS there is a specific page you want to search for.

        Args:
            query (str): the string query to be embedded.
            page_numbers (Optional[List[int]]): Filter by set of pages. Leave as NONE
                if we want to perform a vector search
                over all pages. Otherwise, filter by the set of specified pages.

        ....
    
```

✓ 0s completed at 5:14PM

research\_agent\_databricks

```
VECTOR_INDEX = VectorStoreIndex(nodes)

def vector_query(
    query: str,
    page_numbers: Optional[List[int]] = None
) -> str:
    """Use to answer questions over a given paper.

    Useful if you have specific questions over the paper.
    Always leave page_numbers as None UNLESS there is a specific page you want to search for.

    Args:
        query (str): the string query to be embedded.
        page_numbers (Optional[List[int]]): Filter by set of pages. Leave as NONE
            if we want to perform a vector search
            over all pages. Otherwise, filter by the set of specified pages.

    ....
    page_numbers = page_numbers or []
    metadata_dicts = [
        {"key": "page_label", "value": p} for p in page_numbers
    ]

    query_engine = vector_index.as_query_engine(
        similarity_top_k=2,
        filters=MetadataFilters.from_dicts(
            metadata_dicts,
            condition=FilterCondition.OR
        )
    )
    response = query_engine.query(query)
    return response

```

✓ 0s completed at 5:14PM

The screenshot shows a Jupyter Notebook interface with the following code:

```
query_engine = vector_index.as_query_engine(
    similarity_top_k=2,
    filters=MetadataFilters.from_dicts(
        metadata_dicts,
        condition=FilterCondition.OR
    )
)
response = query_engine.query(query)
return response

vector_query_tool = FunctionTool.from_defaults(
    name=f"vector_tool_{name}",
    fn=vector_query
)

summary_index = SummaryIndex(nodes)
summary_query_engine = summary_index.as_query_engine(
    response_mode="tree_summarize",
    use_async=True,
)
summary_tool = QueryEngineTool.from_defaults(
    name=f"summary_tool_{name}",
    query_engine=summary_query_engine,
    description=(
        f"Useful for summarization questions related to {name}"
    ),
)
return vector_query_tool, summary_tool
```



The screenshot shows a Jupyter Notebook interface with the following details:

- Header:** The URL is `colab.research.google.com/drive/18RUkfBphVSJF-rDh8cOj0QJ6UwQonfh?usp=sharing#scrollTo=aCwgBPcITKgQ`. The top right corner has a "Finish update" button.
- Toolbar:** Includes icons for back, forward, search, and file operations like File, Edit, View, Insert, Runtime, Tools, Help, Saving... A "Comment" and "Share" button are also present.
- Sidebar:** Shows "research\_agent\_databricks PRO" and a "Gemini" section with RAM and Disk options.
- Code Cell:** Contains Python code to import `pathlib` and define a function `paper_to_tools_dict` that iterates over papers, prints their names, and returns a dictionary mapping each paper to its vector tool and summary tool. The code cell is preceded by "+ Code" and "+ Text".
- Output Cell:** Displays the output of the code execution, showing the message "... Getting tools for paper: metapgt.pdf" and "Started parsing the file under job\_id cac11eca-0617-449c-98ef-c43594d8b8ff".

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

Setup an agent over 3 papers

We now setup our function calling agent over 3 papers. We do this by combining the vector/summary tools for each document into a list and passing it to the agent.

```
[ ] (variable) initial_tools: list "selfrag.pdf", "longlora.pdf"
initial_tools= [t for paper in initial_papers for t in paper_to_tools_dict[paper]]
```

```
[ ] # tmp = paper_to_tools_dict["selfrag.pdf"][1]("summary")
# print(str(tmp))
```

```
[ ] from llama_index.core.agent import ReActAgentWorker
from llama_index.core.agent import AgentRunner

agent_worker = ReActAgentWorker.from_tools(
    initial_tools,
    # llm=llm,
    verbose=True
)
agent = AgentRunner(agent_worker)
```

We first query the agent.

```
[ ] response = agent.query(
    "Tell me about the evaluation dataset used in Self-RAG, and then tell me about the evaluation results"
)
```

Thought: The current language of the user is: English. I need to use a tool to help me answer the question.  
Action: summary\_tool.selfrag

colab.research.google.com/drive/18RUKfBipHVSJF-rDh8cOj0QJ6UwQonfh?usp=sharing#scrollTo=aCwgBPcfTKgQ

**research\_agent\_databricks**

File Edit View Insert Runtime Tools Help

Comment Share Gemini J

RAM Disk Gemini

+ Code + Text

Setup an agent over 3 papers

We now setup our function calling agent over 3 papers. We do this by combining the vector/summary tools for each document into a list and passing it to the agent.

```
[ ] initial_papers = ["metagpt.pdf", "selfrag.pdf", "longlora.pdf"]
initial_tools = [t for paper in initial_papers for t in paper_to_tools_dict[paper]]
```

```
[ ] # tmp = paper_to_tools_dict["selfrag.pdf"][1]("summary")
# print(str(tmp))
```

```
[ ] from llama_index.core.agent import ReActAgentWorker
from llama_index.core.agent import AgentRunner
```

```
agent_worker = ReActAgentWorker.from_tools(
    initial_tools,
    # llm_llm,
    verbose=True
)
agent = AgentRunner(agent_worker)
```

We first query the agent.

```
[ ] response = agent.query(
    "Tell me about the evaluation dataset used in Self-RAG, and then tell me about the evaluation results"
)
```

Executing (1m 9s) <--> ge... > ... > build... > \_build... > \_add... > \_get\_no... > e... > w... > get\_tex... > \_get\_t... > ... > e... > f... > \_wra... > ... > fo... >

colab.research.google.com/drive/18RUKfBipHVSJF-rDh8cOj0QJ6UwQonfh?usp=sharing#scrollTo=aCwgBPcfTKgQ

**research\_agent\_databricks**

File Edit View Insert Runtime Tools Help

Comment Share Gemini J

RAM Disk Gemini

+ Code + Text

```
[10]     query_organic_summary_query_engage,
        description=(
            f"Useful for summarization questions related to {name}"
        ),
    )

    return vector_query_tool, summary_tool
```

```
from pathlib import Path
```

```
paper_to_tools_dict = {}
for paper in papers:
    print(f"Getting tools for paper: {paper}")
    vector_tool, summary_tool = get_doc_tools(paper, Path(paper).stem)
    paper_to_tools_dict[paper] = [vector_tool, summary_tool]
```

... Getting tools for paper: metagpt.pdf  
Started parsing the file under job\_id cac11eca-0617-449c-98ef-c43594d8b8ff  
Getting tools for paper: longlora.pdf  
Started parsing the file under job\_id cac11eca-06e5-45c6-9470-84d6903584a0  
Getting tools for paper: selfrag.pdf  
Started parsing the file under job\_id cac11eca-4420-43e8-8b78-1bcab890bd15

Executing (1m 14s) <--> ge... > ... > build... > \_build... > \_add... > \_get\_no... > e... > w... > get\_tex... > \_get\_t... > ... > e... > f... > \_wra... > ... > fo... >

research\_agent\_databricks

```
[ ] initial_papers = ["metagpt.pdf", "selffrag.pdf", "longlora.pdf"]
initial_tools = [t for paper in initial_papers for t in paper_to_tools_dict[paper]]  
  
[ ] # tmp = paper_to_tools_dict["selffrag.pdf"][1]("summary")
# print(str(tmp))  
  
from llama_index.core.agent import ReActAgentWorker
from llama_index.core.agent import AgentRunner
agent_worker = ReActAgentWorker.from_tools(
    initial_tools,
    # llm=llm,
    verbose=True
)
agent = AgentRunner(agent_worker)
```

We first query the agent.

```
[ ] response = agent.query(
    "Tell me about the evaluation dataset used in Self-RAG, and then tell me about the evaluation results"
)  
  
Thought: The current language of the user is: English. I need to use a tool to help me answer the question.
Action: summary_tool_selffrag
Action Input: {'input': 'evaluation dataset used in Self-RAG'}  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.9971250610078541 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.857028781313381 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.1347305779126 seconds as it raised RateLimitError: Error code: 429 - {'  
Executing (1m 17s) <--> ge... > _... > ... > build_... > _build_... > _add_... > _get_no... > e... > w... > get_text... > _get_t... > ... > e... > f... > _wra... > ... > fo... > ...
```

research\_agent\_databricks

```
[ ] We first query the agent.  
  
[ ] response = agent.query(
    "Tell me about the evaluation dataset used in Self-RAG, and then tell me about the evaluation results"
)  
  
Thought: The current language of the user is: English. I need to use a tool to help me answer the question.
Action: summary_tool_selffrag
Action Input: {'input': 'evaluation dataset used in Self-RAG'}  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.9971250610078541 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.857028781313381 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.1347305779126 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.861534066685435019 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.440033187216304 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.8316333362876902 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.688663848888767 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.4371998234106036 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.4963714134431002 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.001386880039630807 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.5194930129606028 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.3731566711292673 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.986740271139635 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.20581476747561422 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.0194762023180108 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.045137252218768964 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.6807583261753372 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.9198524828467187 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 1.297539029849517 seconds as it raised RateLimitError: Error code: 429 - {'  
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.8941308094886549 seconds as it raised RateLimitError: Error code: 429 - {'  
Observation: Error: Error code: 429 - {'error_code': 'REQUEST_LIMIT_EXCEEDED', 'message': 'REQUEST_LIMIT_EXCEEDED: Exceeded workspace rate limit for databricks-meta-llm'}
```

Thought: It seems like I've hit a rate limit. Let me try again with a different tool.

```
Action: vector_tool_selffrag
Action Input: {'query': 'evaluation dataset used in Self-RAG', 'page_numbers': None}
Observation: The evaluation dataset used in Self-RAG is not explicitly mentioned in the provided context. However, it is mentioned that experiments show that SELF-RAG (
```

Thought: I have some information about the evaluation dataset, but it's not explicit. I need to ask another question to get more information about the evaluation result.

The screenshot shows a Jupyter Notebook interface with the following details:

- Header:** colab.research.google.com/drive/1B8Ukf8ipHVSJF-rDhsCoJQJ6UwQnfh?usp=sharing#scrollTo=wHcm1v6TKgQ
- Toolbar:** Comment, Share, Gemini, RAM Disk, Finish update.
- Code Cell:** `[16] log.debug(`  
[16] APIConnectionError: Connection error.
- Next steps:** Explain error
- Output:** The evaluation dataset used in Self-RAG is not explicitly mentioned, but the model is evaluated on a diverse set of tasks, including Open-domain QA, reasoning, and fact...
- Code Cell:** response = agent.query("What are the MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper?")  
print(str(response))
- Output:** Thought: The current language of the user is: English. I need to use a tool to help me answer the question.  
Action: vector\_tool\_metaapt  
Action Input: {'query': 'What are the MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper?', 'page\_numbers': [8])  
Observation: The comparisons between MetaGPT and ChatDev are described in terms of several metrics, including executability, running times, token usage, code statistics  
Thought: I can answer without using any more tools. I'll use the user's language to answer  
Answer: The MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper are in terms of several metrics, including executability, running times, token usage, code...
- Code Cell:** response = agent.query(  
 "Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?"  
)
- Output:** Thought: The current language of the user is: English. I need to use a tool to help me answer the question.  
Action: summary\_tool\_selfrag  
Action Input: {'input': 'Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?'})  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_chat in 0.6189579289428381 seconds as it raised RateLimitError: Error code: 429 - ...  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_chat in 0.4460745636595751 seconds as it raised RateLimitError: Error code: 429 - ...  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_chat in 0.0029317465256685216 seconds as it raised RateLimitError: Error code: 429 - ...

```
[ ] print(str(response))

[+] The evaluation dataset used in Self-RAG is not explicitly mentioned, but the model is evaluated on a diverse set of tasks, including Open-domain QA, reasoning, and fact

[ ] response = agent.query("What are the MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper?")
print(str(response))

Thought: The current language of the user is: English. I need to use a tool to help me answer the question.
Action: vector_tool_metaapt
Action Input: {'query': 'What are the MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper?', 'page_numbers': [8]}
Observation: The comparisons between MetaGPT and ChatDev are described in terms of several metrics, including executability, running times, token usage, code statistics
Thought: I can answer without using any more tools. I'll use the user's language to answer
Answer: The MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper are in terms of several metrics, including executability, running times, token usage, code statistics

[ ] response = agent.query("What are the MetaGPT comparisons with ChatDev described on page 8 of the MetaGPT paper?")
print(str(response))

Thought: The current language of the user is: English. I need to use a tool to help me answer the question.
Action: summary_tool_selfrag
Action Input: {'input': 'Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?'}

WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.6189579289428381 seconds as it raised RateLimitError: Error code: 429 - {
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.4460745636595751 seconds as it raised RateLimitError: Error code: 429 - {
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.0029317465256685216 seconds as it raised RateLimitError: Error code: 429
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.6111833642944608 seconds as it raised RateLimitError: Error code: 429 - {
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.159664755556576 seconds as it raised RateLimitError: Error code: 429 - {
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.014727069513607138 seconds as it raised RateLimitError: Error code: 429 -
WARNING:llama_index.llms.openai.utils:Retrying llama_index.llms.openai.base.OpenAI._achat in 0.87832889780765717 seconds as it raised RateLimitError: Error code: 429 -
```

colab.research.google.com/drive/nBRUkfBpHVSJF-rDh6c0j0QJ6UwQorfh?usp=sharing#scrollTo=ehjoKT9pTkqQ

research\_agent\_databricks ☆

File Edit View Insert Runtime Tools Help All changes saved

Comment Share J

RAM Disk Gemini

+ Code + Text

VC Metrics: Comparisons with ChatGPT described on page 6 of the Victoria paper. By default, it uses the DALL-E method, and includes calculate latency, training latency, context usage, code size.

response = agent.query()  
"Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?"

Thought: The current language of the user is: English. I need to use a tool to help me answer the question.  
Action: summary\_tool\_selfrag  
Action Input: {'input': 'Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?'}

WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.6189579289428381 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.4460745636595751 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.0029317465256685216 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.6111883642944668 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.159664755556576 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.014727069513607138 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.07832889788765717 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.567737610079809 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.4909845814884308 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.2982103524881386 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.08970114895048653 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.656936944357149 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.42567342788570694 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.29546555626701343 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.6740557219510865 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 1.99969766622177308 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 1.909071481323557 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 1.5203088317667537 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 1.3280017031572666 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.4044926633600179 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 1.9640423910998424 seconds as it raised RateLimitError: Error code: 429 - {  
Observation: Error: Error code: 429 - {'error\_code': 'REQUEST\_LIMIT\_EXCEEDED', 'message': 'REQUEST\_LIMIT\_EXCEEDED: Exceeded workspace rate limit for databricks-meta-llm'}

Thought: It seems like I've hit a rate limit. Let me try a different tool to get the information I need.  
Action: summary\_tool\_metalgpt  
Action Input: {'input': 'Compare the complexity of the approaches in Self-RAG and MetaGPT. Which approach uses more tokens?'}

WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.7844992042396307 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.4866682492741915 seconds as it raised RateLimitError: Error code: 429 - {  
WARNING:llama\_index.llms.openai.utils:Retrying llama\_index.llms.openai.base.OpenAI.\_achat in 0.0768037514413085 seconds as it raised RateLimitError: Error code: 429 - {

21s completed at 5:20PM