

Welcome to ServerlessToronto.org

An Evening with Mark Ryan and Jerry Liu

- 6:00 - 6:10 Networking & Opening remarks
- 6:10 - 6:35 Mark Ryan: The LLM Landscape
- 6:35 - 7:15 Jerry Liu: Solving Core Challenges in RAG Pipelines
- 7:15 - 7:45 Q&A
- 7:45 - 8:00 Manning Publications raffle

Hello
my name is

Introduce Yourself:

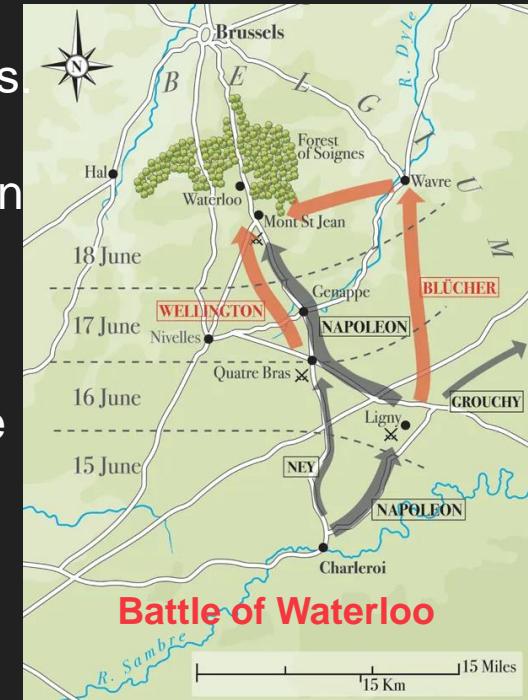
- Where from? Why are you here?
- Looking for work or Offering work?

Help us serve you better: bit.ly/slsto



Why this Generative AI Talk?

- 1. Navigating the Tsunami:** Understand the sweeping changes the "GenAI tsunami" brings to industries and jobs.
- 2. Situational Awareness:** Learn from AI leaders Mark Ryan & Jerry Liu to gain a strategic view of the LLM and RAG landscape.
- 3. Career Transformation:** Learn to position yourself as the architect of automation rather than its subject.
- 4. Practical Advice:** Acquire actionable strategies to apply Generative AI within your enterprise.
- 5. Interactive Learning:** Engage in live Q&A to discuss and clarify your AI dilemmas with experts.

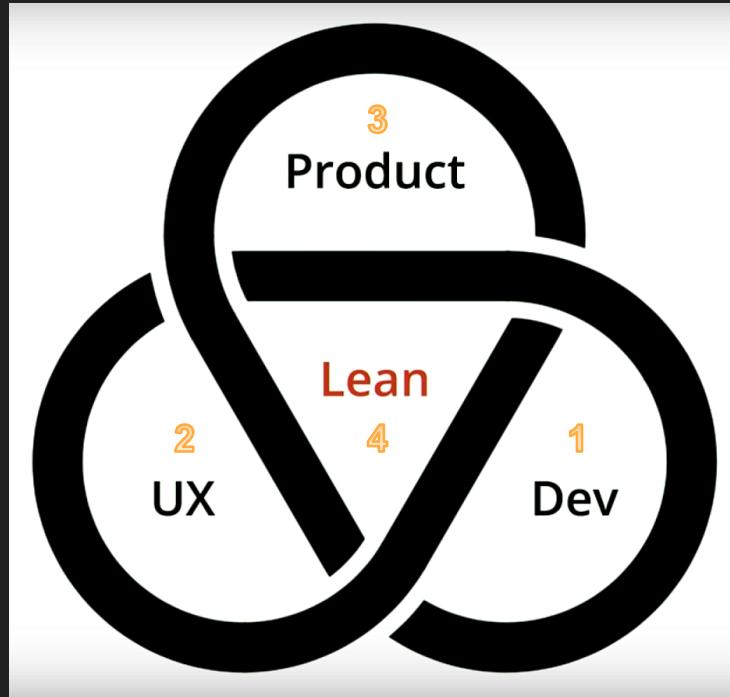


What is Serverless Toronto about?

#1 We started as **Back-end FaaS** Developers who enjoyed 'gluing together' other people's APIs and Managed Services



#2 We build bridges between Serverless Community ("Dev leg"), and Front-end, Voice-First & UX folks ("UX leg")



#3 We're obsessed with creating business value (meaningful Products), focusing on Outcomes/Impact – **NOT** Outputs



#4 Achieve agility **NOT** by "sprinting" faster but working **smarter** (by using bigger building blocks & **less** Ops)

Serverless became New Agile & Mindset

Serverless is a State of Mind...

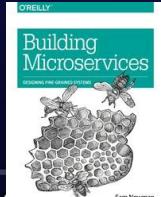
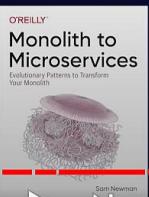
GOTO 2020 • When To Use Microservices (And When Not To!) • Sam Newman & Martin Fowler

WHEN TO USE MICROSERVICES



SAM NEWMAN

Author of "Monolith to Microservices"



5:02 / 38:4

We focus on the tech tool, not the thing
that the tech tool let's you do

Way too often, we – the IT folks,
have obsession with “pimping up
our cars” (infrastructure / code /
pipelines) instead of “driving
business” forward & taking them
places ☺



goto,

serverless

user group

... It is a way to focus on business value.



Jared Short:

1. If the platform has it, use it
2. If the market has it, buy it
3. If you can reconsider requirements, do it
4. If you have to build it, own it.

Ben Kehoe: Serverless is about how you make decisions, not about your choices.

It can be applied to any Tech stack, even On-Prem

Upcoming ServerlessToronto.org Meetups



Master serverless.

theburningmonk.com

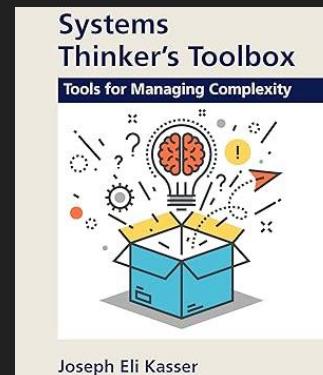
serverless user group

Yan Cui's Money-Saving Tips for the Frugal Serverless Developer

Friday Lunch & Learn, April 19



**From Solo to Success:
Leveraging AI for Your Startup
with Pankaj Upreti**



Summer 2024

serverless
user group

Knowledge Sponsor



MANNING PUBLICATIONS



1. Go to www.manning.com
2. Select *any* e-Book, Video course, or liveProject you want!
3. Add it to your shopping cart (no more than 1 item in the cart)
4. Raffle winners will send me the emails (used in Manning portal),
5. So the publisher can move it to your Dashboard – as if purchased.

Fill out the Survey to win: bit.ly/slsto

serverless
user group

Feature Presentations:

ALL IN AI with Googler Mark Ryan, and LlamaIndex Creator Jerry Liu



serverless
user group

LLM Landscape

Mark Ryan
Developer Knowledge Platform AI Lead, Google Cloud
ryanmark2014@gmail.com



Generative AI Milestones

Major Generative AI Milestones: Part 1

Attention Is All You Need:

Seminal paper from Google that introduced transformers

Oct 2018

GPT-2: OpenAI LLM

DALLE: OpenAI image model

Jun 2017

BERT: Google transformer-based language model

Feb 2019

May 2020

GPT-3: OpenAI LLM

Jan 2021

May 2021

LaMDA: Google LLM

Imagen: Google image model

PaLM: Google LLM

Codex: OpenAI code model

Gato: DeepMind multimodal model

Apr 2022

Aug 2022

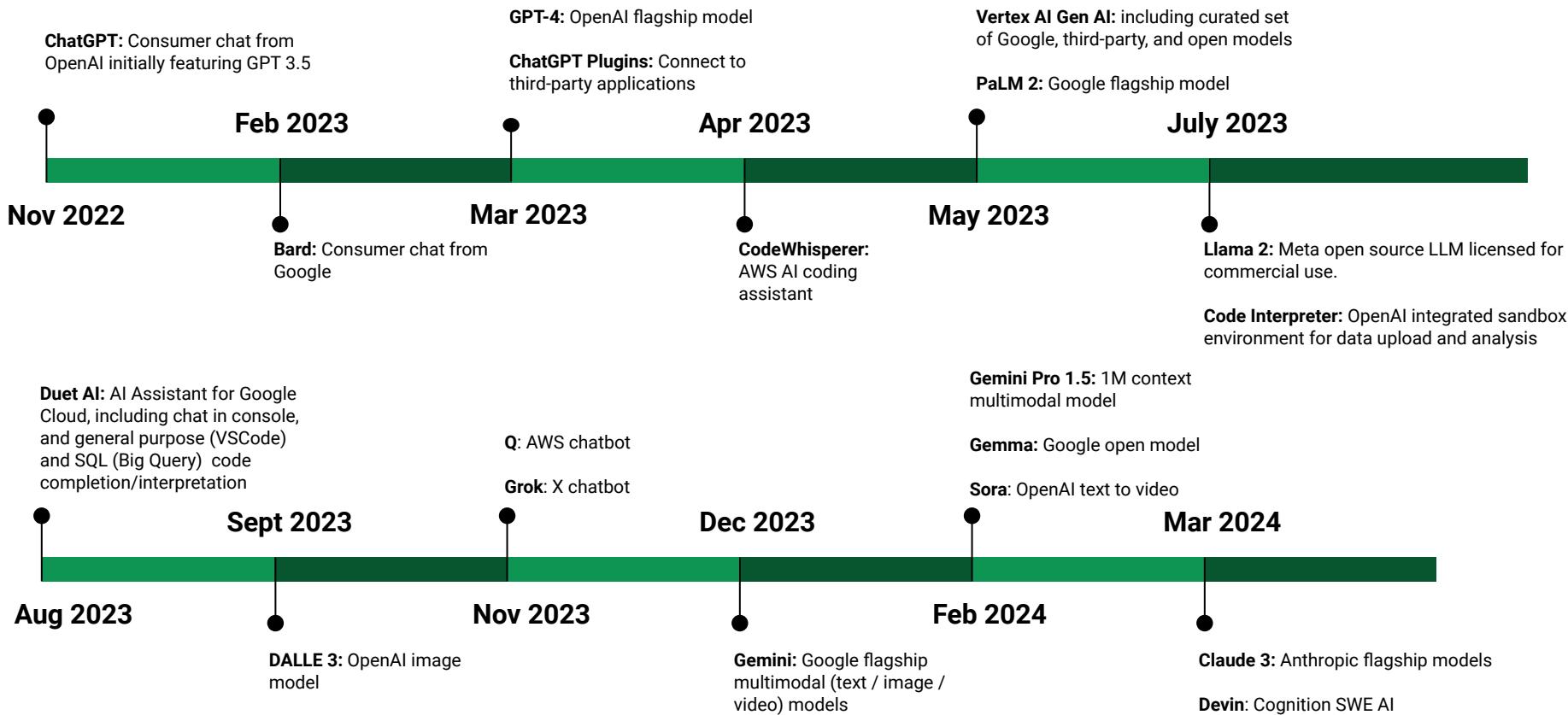
Aug 2021

DALLE 2: OpenAI image model

May 2022

Stable Diffusion: Image model

Major Generative AI Milestones: Part 2



Ecosystem and Vendor Landscape

The Emerging LLM Ecosystem

	Examples	Description	Use Case
Vector databases	<ul style="list-style-type: none">• Pinecone• Chroma• Vertex AI Vector Search	Store and find associations between embeddings, high-dimensional vector representations of data	Grounding LLM responses in a set of documents (example of RAG)
Encapsulated coding environments	OpenAI Code Interpreter / Advanced Data Analysis	Upload datasets & ask questions to get visualizations and code running in a limited Python instance	Ad hoc data analysis
Plugins / extensions	<ul style="list-style-type: none">• ChatGPT plugins / GPTs• Vertex AI extensions	Connect LLMs to third-party / external applications	Access current data / query & modify data that is external to the LLM
LLM app development frameworks	<ul style="list-style-type: none">• LangChain• Llamaindex• Autogen	LLM-centric framework to manage workflow (data sources, agents, models, etc)	Assembling LLM-based applications

Generative AI Landscape by Vendor

Vendor	Prod. Suite Assistance	Developer / Ops Assistant	Consumer Chat	Enterprise Gen AI	Dev / Hobbyist Gen AI	Open Foundation Models
Google	Gemini for Google Workspace	Duet AI for Google Cloud	Gemini	Vertex AI	Google AI for Developers	Gemma
Microsoft	CoPilot 365	Github Copilot	Bing Chat	Azure OpenAI		
OpenAI		ChatGPT	ChatGPT	ChatGPT Enterprise	ChatGPT	
AWS		• Q • CodeWhisperer		Bedrock / Titan		
Anthropic			Claude 3*	Claude 3*		
Meta						Llama 2*
Mistral						Mixtral 8x7B*

Twitter: [@MarkRyanMkm](https://twitter.com/@MarkRyanMkm)

LinkedIn:

www.linkedin.com/in/mark-ryan-31826743

YouTube:

[@markryan2475](https://www.youtube.com/@markryan2475)



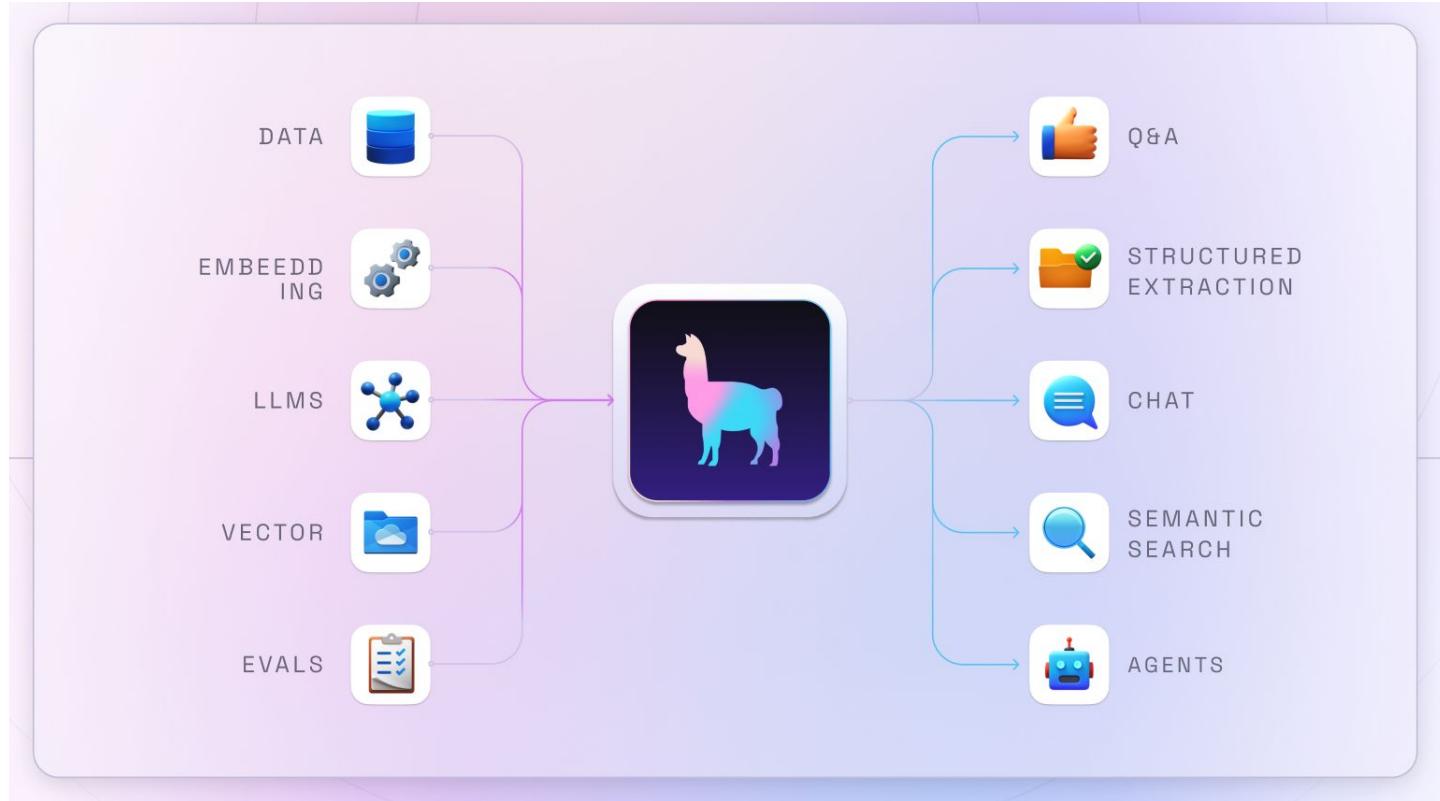


RAG in 2024

Jerry Liu, LlamaIndex co-founder/CEO



Llamaindex: Context Augmentation for your LLM app





Paradigms for inserting knowledge

Retrieval Augmentation - Fix the model, put context into the prompt



Before college the two main things I worked on, outside of school, were writing and programming. I didn't write essays. I wrote what beginning writers were supposed to write then, and probably still are: short stories. My stories were awful. They had hardly any plot, just characters with strong feelings, which I imagined made them deep...



Input Prompt

Here is the context:
Before college the two main things...



Given the context,
answer the following
question:
{query_str}

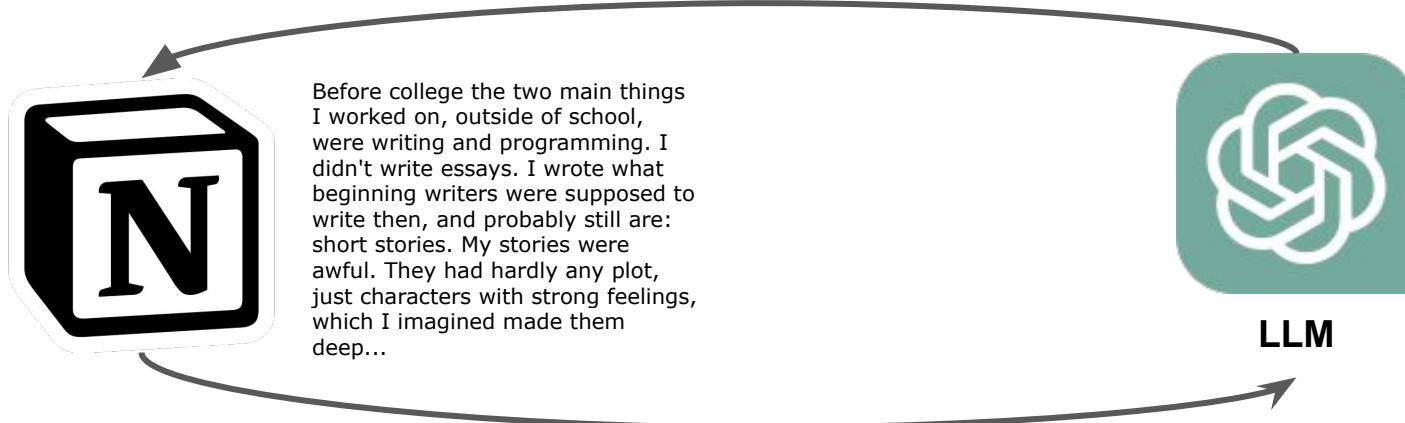


LLM



Paradigms for inserting knowledge

Fine-tuning - baking knowledge into the weights of the network

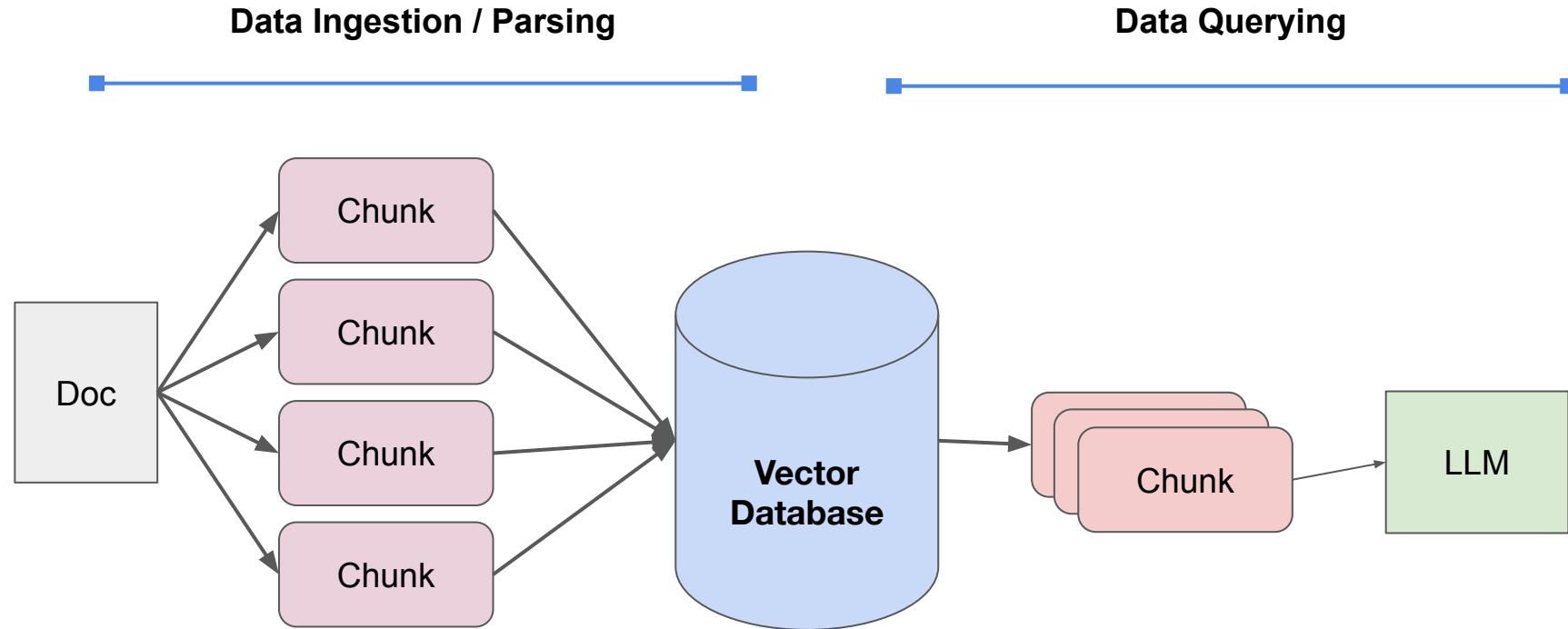




RAG Stack



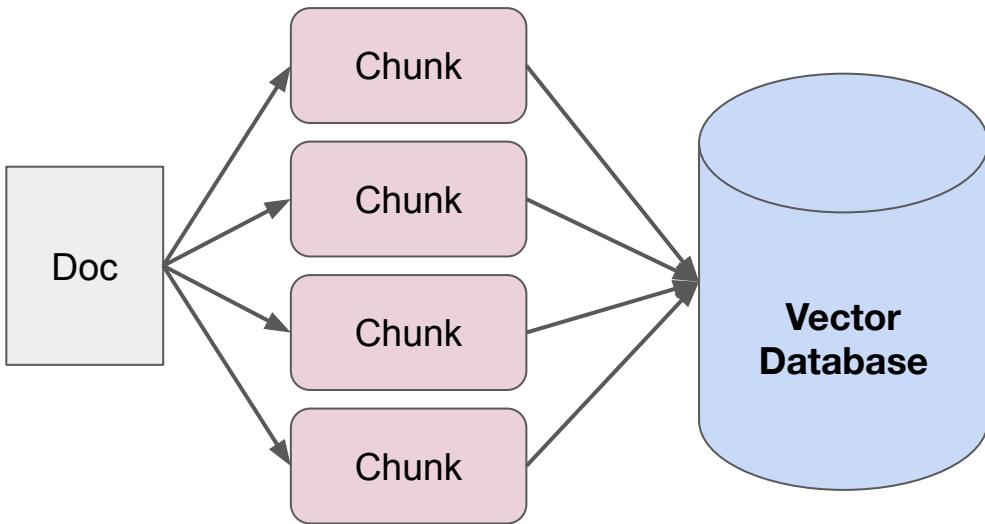
Current RAG Stack for building a QA System



5 Lines of Code in LlamalIndex!



Current RAG Stack (Data Ingestion/Parsing)



Process:

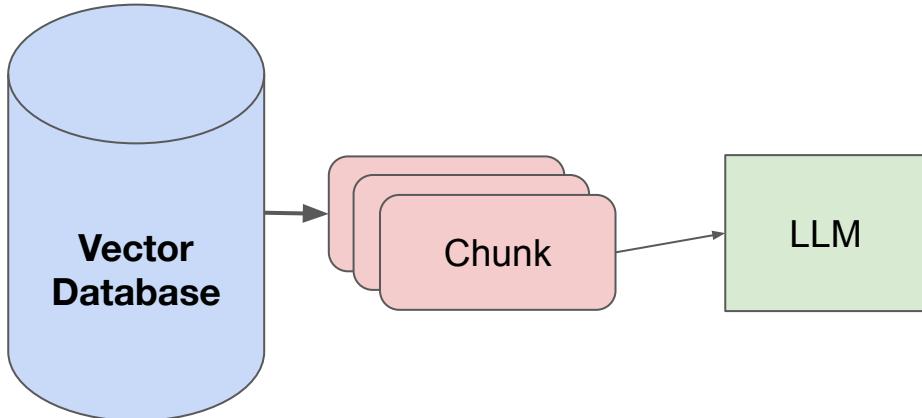
- Split up document(s) into even chunks.
- Each chunk is a piece of raw text.
- Generate embedding for each chunk (e.g. OpenAI embeddings, sentence_transformer)
- Store each chunk into a vector database



Current RAG Stack (Querying)

Process:

- Find top-k most similar chunks from vector database collection
- Plug into LLM response synthesis module

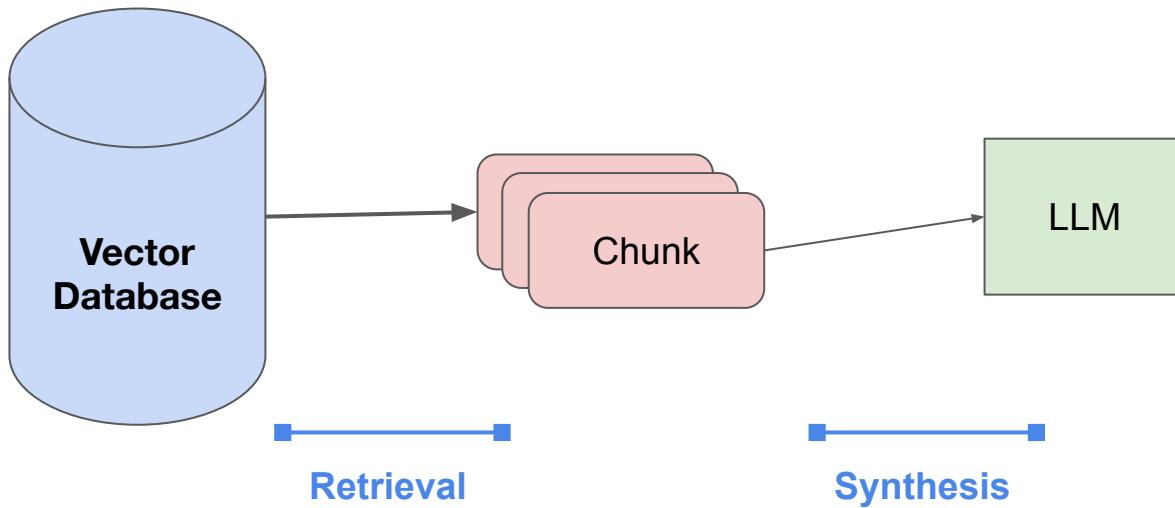




Current RAG Stack (Querying)

Process:

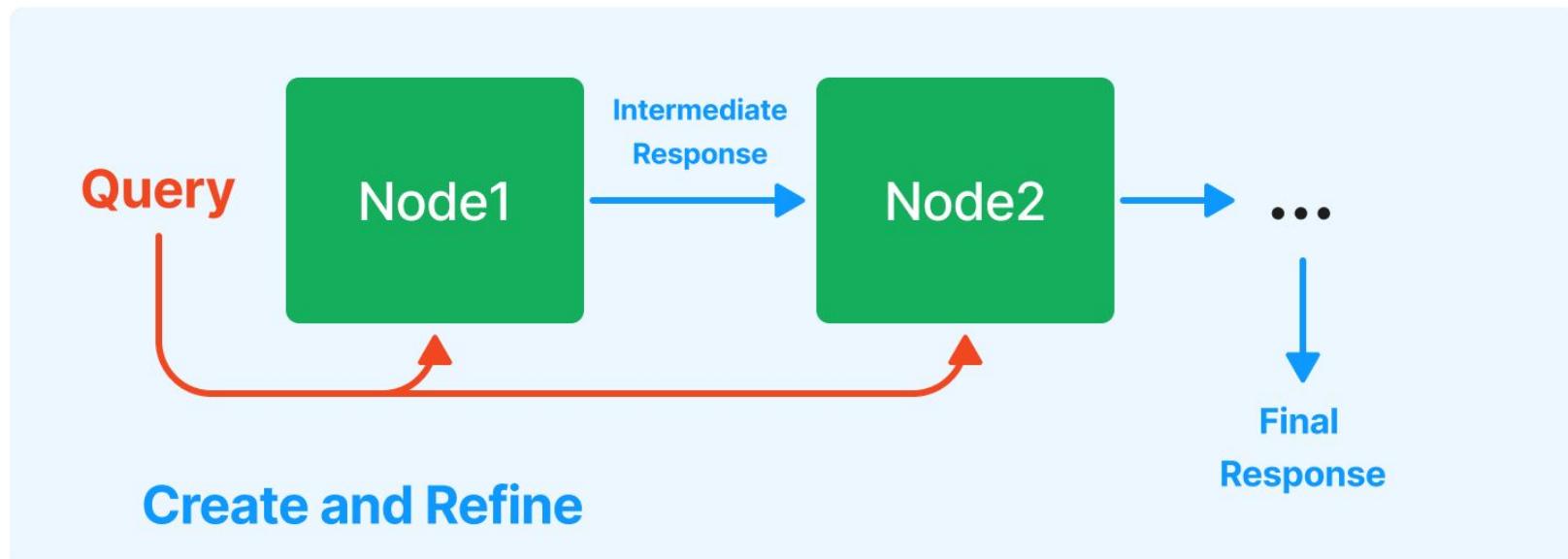
- Find top-k most similar chunks from vector database collection
- Plug into LLM **response synthesis module**





Response Synthesis

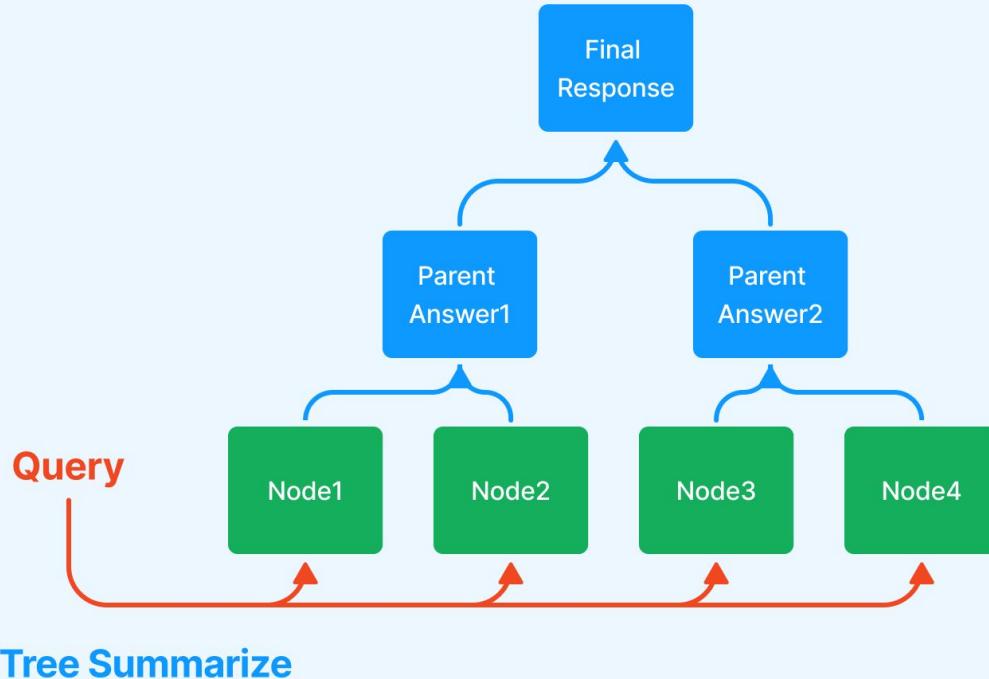
Create and refine





Response Synthesis

Tree Summarize





Quickstart

https://colab.research.google.com/drive/1knQpGJLHj-LTTHqlZhgcjDH5F_nJliY0?usp=sharing



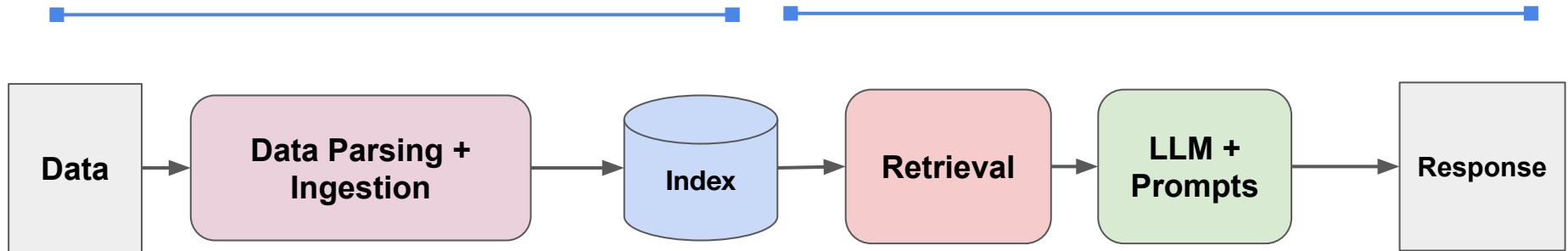


Challenges with “Naive” RAG



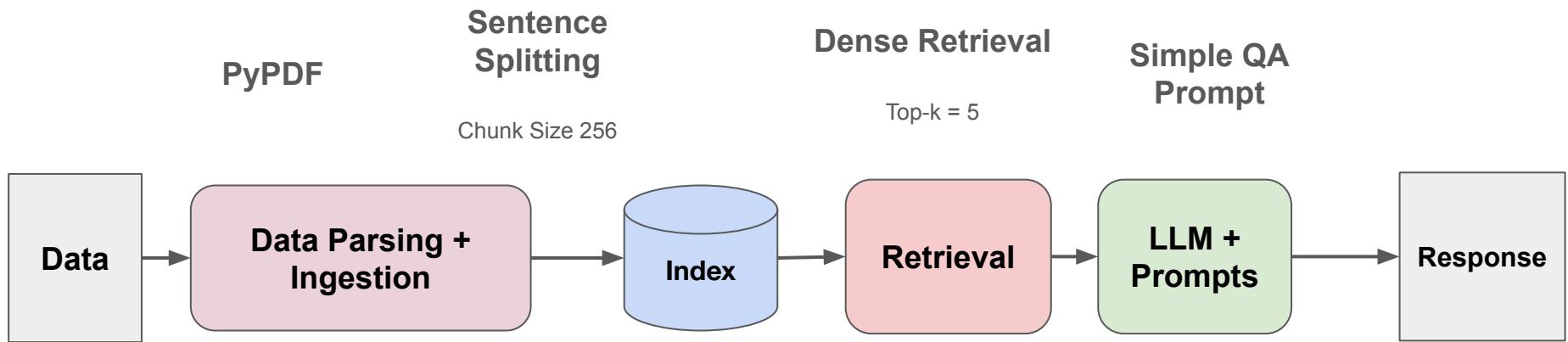
RAG

Data Parsing & Ingestion Data Querying





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)



Easy to Prototype, Hard to Productionize

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

Failure Modes:

- Response Quality: Bad Retrieval, Bad Response Generation
- Hard to Improve: Too many parameters to tune
- Systems: Latency, Cost, Security



Easy to Prototype, Hard to Productionize

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

Failure Modes:

- **Response Quality:** Bad Retrieval, Bad Response Generation
- **Hard to Improve:** Too many parameters to tune
- **Systems:** Latency, Cost, Security



Challenges with Naive RAG (Response Quality)

- Bad Retrieval
 - **Low Precision:** Not all chunks in retrieved set are relevant
 - Hallucination + Lost in the Middle Problems
 - **Low Recall:** Not all relevant chunks are retrieved.
 - Lacks enough context for LLM to synthesize an answer
 - **Outdated information:** The data is redundant or out of date.



Challenges with Naive RAG (Response Quality)

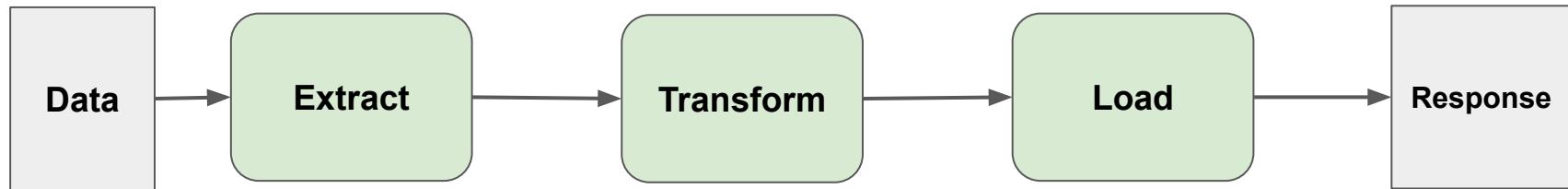
- Bad Retrieval
 - **Low Precision:** Not all chunks in retrieved set are relevant
 - Hallucination + Lost in the Middle Problems
 - **Low Recall:** Not all relevant chunks are retrieved.
 - Lacks enough context for LLM to synthesize an answer
 - **Outdated information:** The data is redundant or out of date.
- Bad Response Generation
 - **Hallucination:** Model makes up an answer that isn't in the context.
 - **Irrelevance:** Model makes up an answer that doesn't answer the question.
 - **Toxicity/Bias:** Model makes up an answer that's harmful/offensive.



Difference with Traditional Software

Traditional software is defined by a set of programmatic rules.

Given an input, you can easily reason about the expected output.



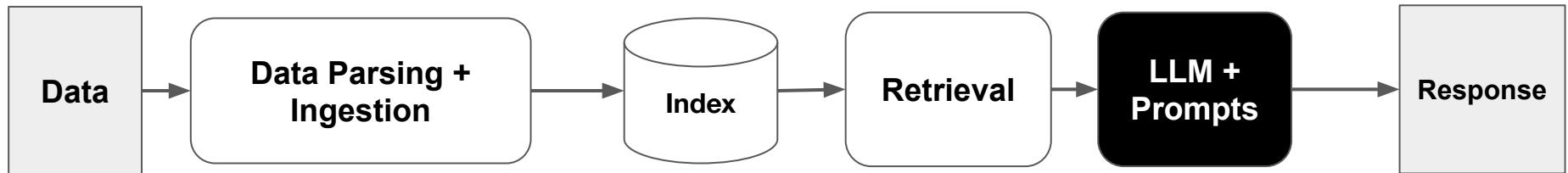
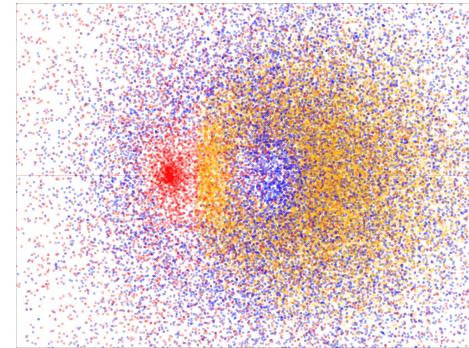


Difference with Traditional Software

AI-powered software is defined by a **black-box set of parameters**.

It is really hard to reason about what the function space looks like.

The model parameters are tuned, the surrounding parameters (prompt templates) are not.

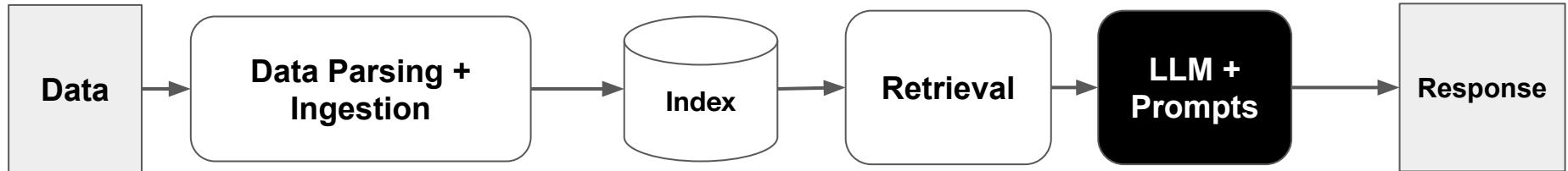




Difference with Traditional Software

If one component of the system is a black-box, all components of the system become black boxes.

The more components, the more parameters you have to tune.

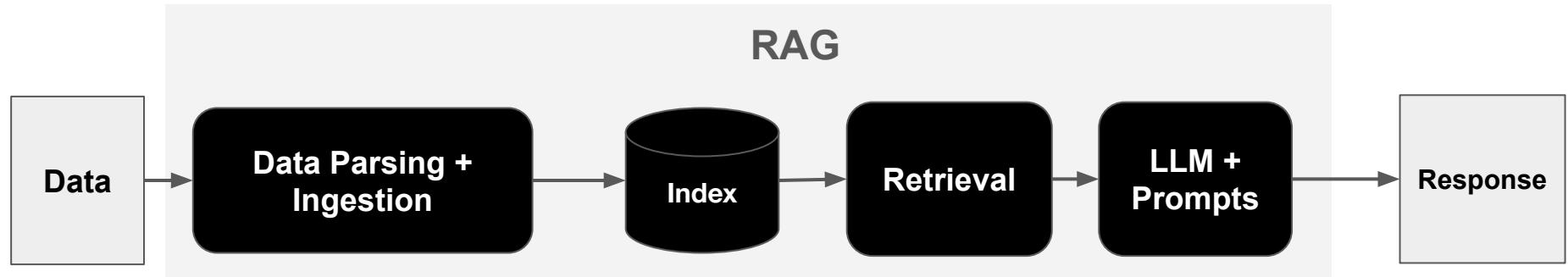




Difference with Traditional Software

If one component of the system is a black-box, all components of the system become black boxes.

Every parameter affects the performance of the end system.



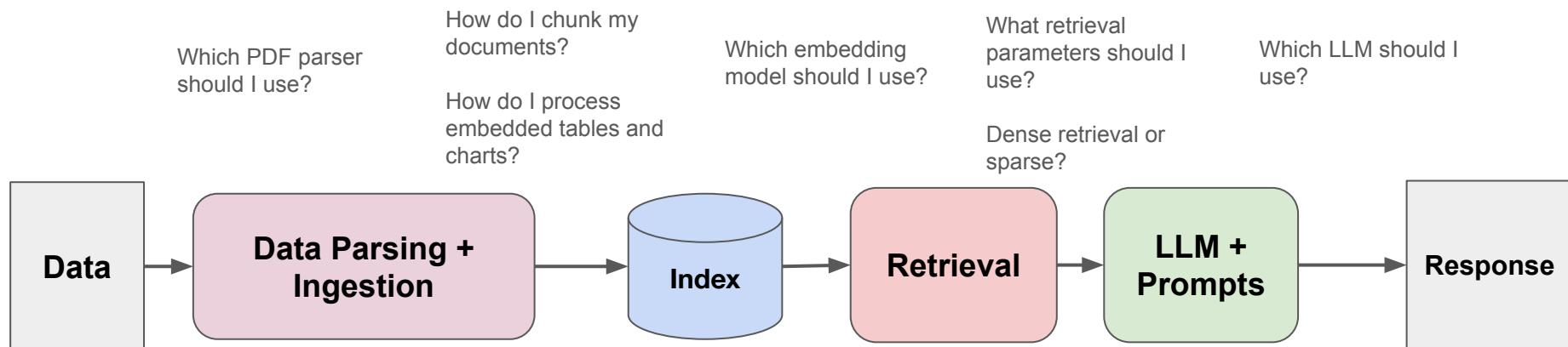


There's Too Many Parameters

Every parameter affects the performance of the entire RAG pipeline.

Which parameters should a user tune?

There's too many options!





Mapping Pain Points to Solutions



Solution

Categorize by pain point, and establish best practices



Solution

Categorize by pain point, and establish best practices

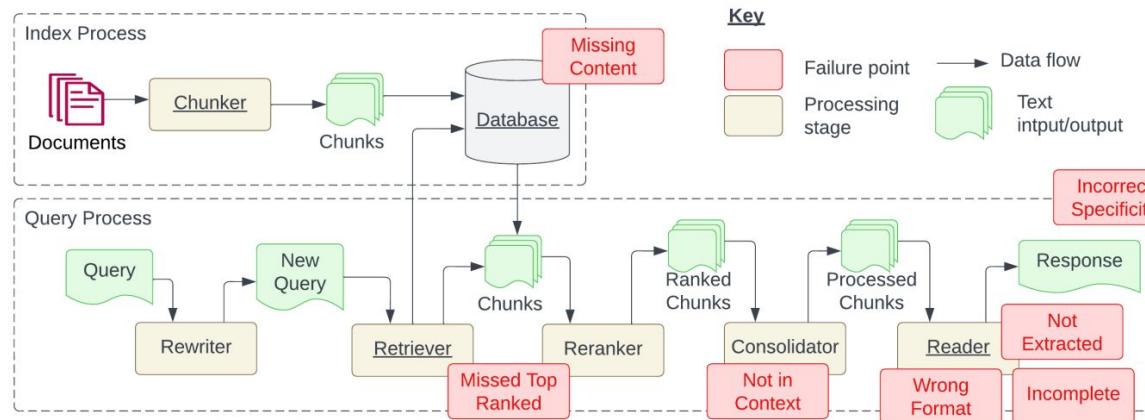


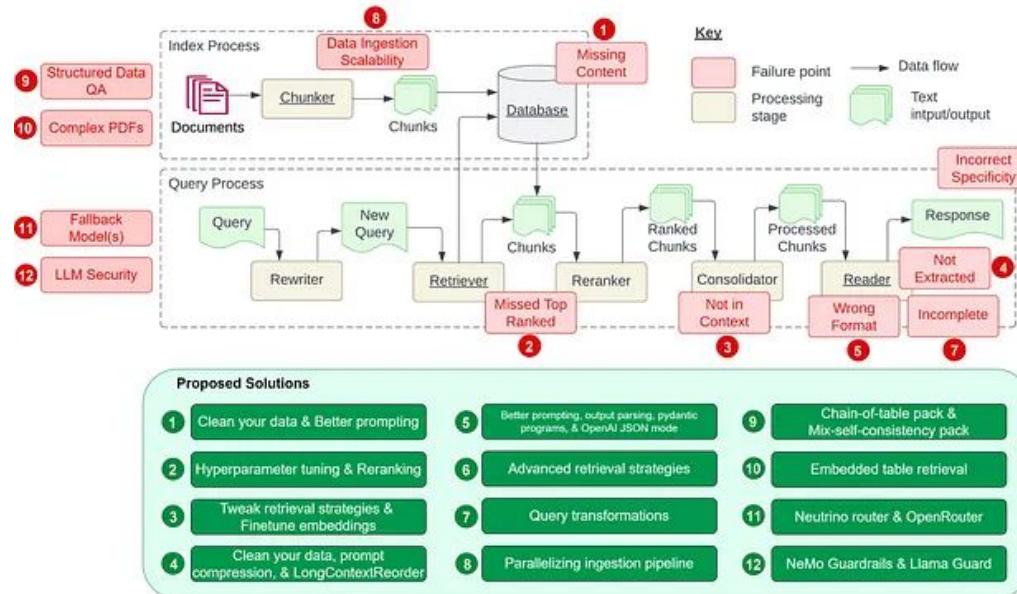
Figure 1: Indexing and Query processes required for creating a Retrieval Augmented Generation (RAG) system. The indexing process is typically done at development time and queries at runtime. Failure points identified in this study are shown in red boxes. All required stages are underlined. Figure expanded from [19].

[“Seven Failure Points When Engineering a Retrieval Augmented Generation System”](#), Barnett et al.



Solution

Categorize by pain point, and establish best practices



[“12 RAG Pain Points and Proposed Solutions”, by Wengi Glantz](#)



Pain Points

Response Quality Related

1. Context Missing in the Knowledge Base
2. Context Missing in the Initial Retrieval Pass
3. Context Missing After Reranking
4. Context Not Extracted
5. Output is in Wrong Format
6. Output has Incorrect Level of Specificity
7. Output is Incomplete

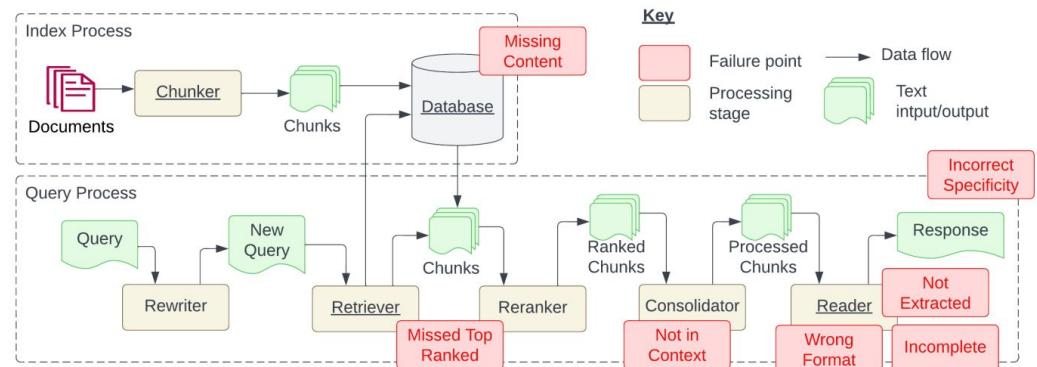


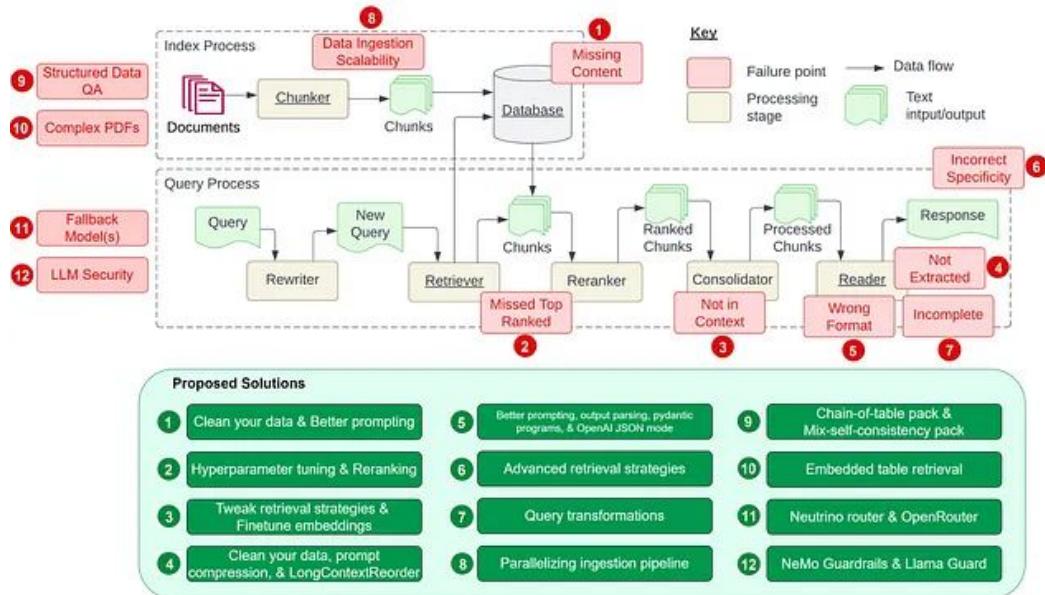
Figure 1: Indexing and Query processes required for creating a Retrieval Augmented Generation (RAG) system. The indexing process is typically done at development time and queries at runtime. Failure points identified in this study are shown in red boxes. All required stages are underlined. Figure expanded from [19].



Pain Points

Scalability

8. Can't Scale to Larger Data Volumes



Security

12. LLM Security

Use Case Specific

9. Ability to QA Tabular Data

10. Ability to Parse PDFs



Pain Points

Scalability

8. Can't Scale to Larger Data Volumes

11. Rate-Limit Errors

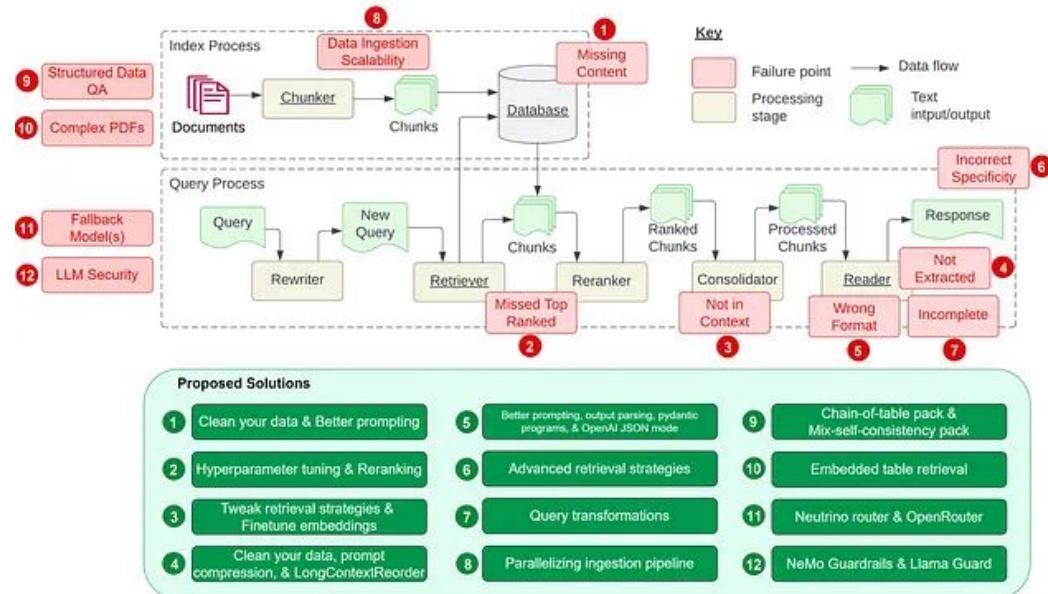
Security

12. LLM Security

Use Case Specific

9. Ability to QA Tabular Data

10. Ability to Parse PDFs





Let's figure out solutions

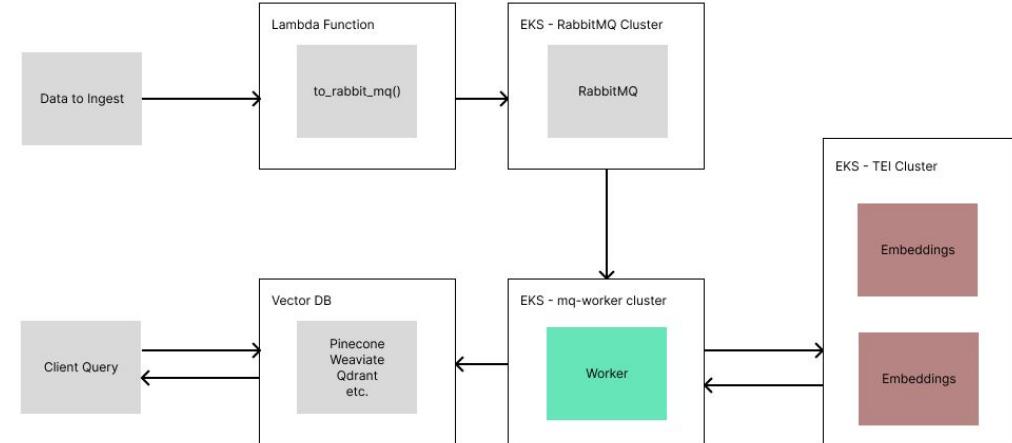


1. Context Missing in the Knowledge Base

Clean your data: Pick a good document parser (more on this later!)

Add in Metadata: inject global context to each chunk

Keep your data updated: Setup a recurring data ingestion pipeline. Upsert documents to prevent duplicates.

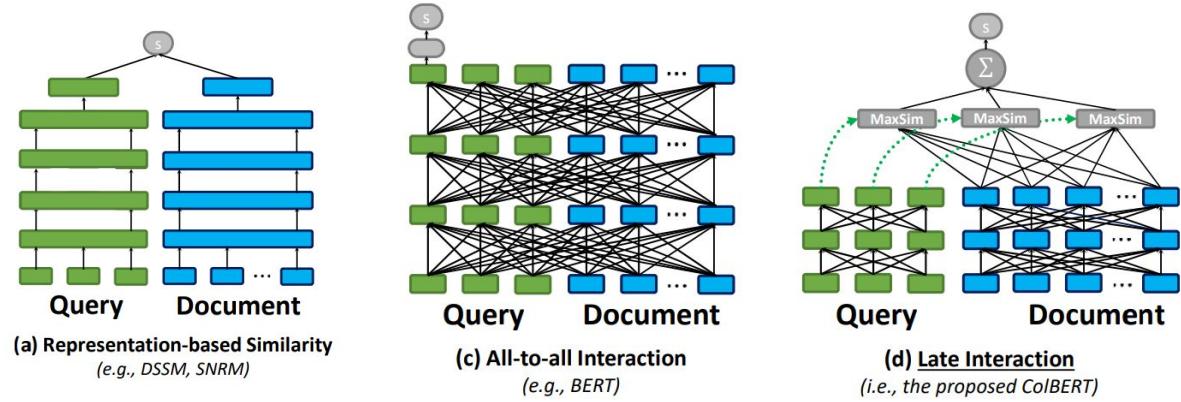




2. Context Missing in the Initial Retrieval Pass

Solution: Hyperparameter tuning for chunk size and top-k

Solution: Reranking



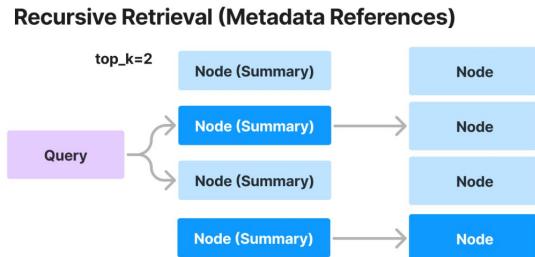
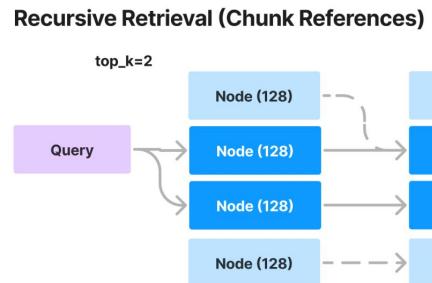
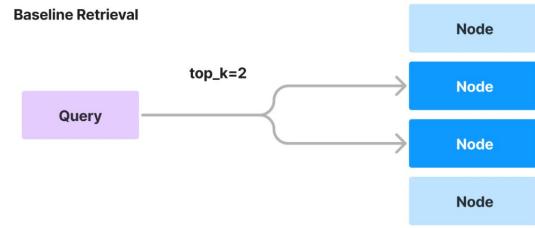
Source: ColBERT



3. Context Missing After Reranking

Solution: try out fancier retrieval methods
(small-to-big, auto-merging, auto-retrieval,
ensembling, ...)

Solution: fine-tune your embedding models
to task-specific data

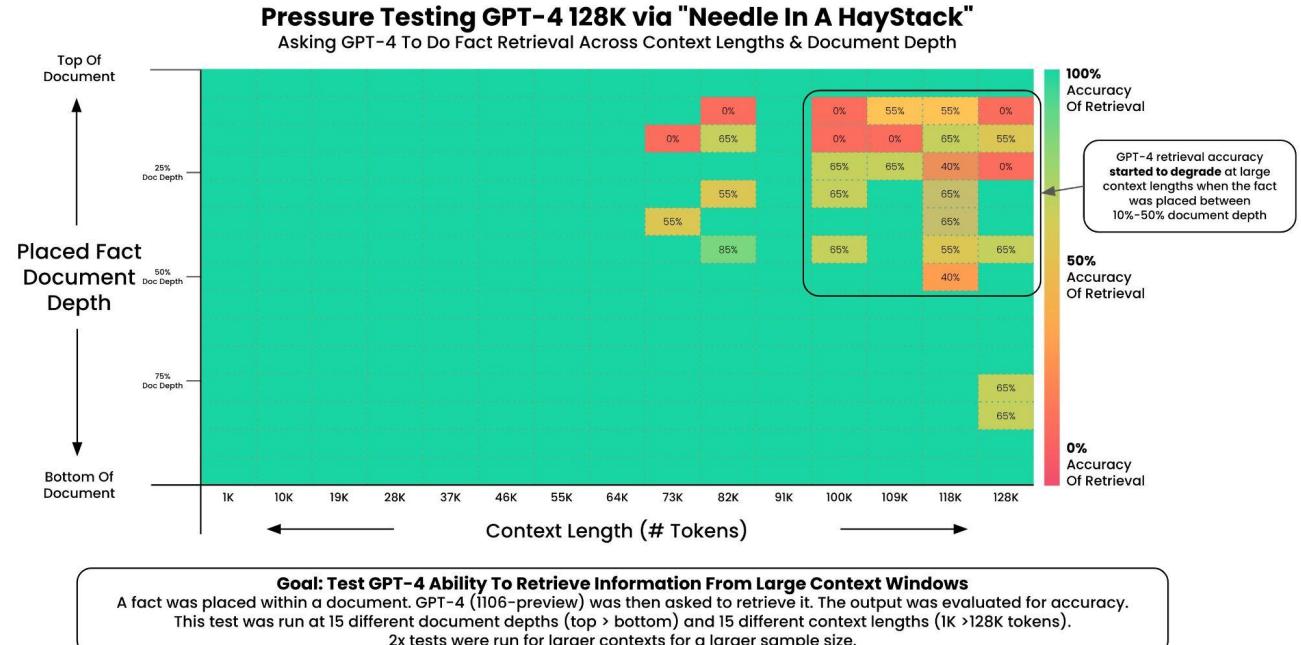




4. Context is there, but not extracted by the LLM

The context is there,
but the LLM doesn't
understand it.

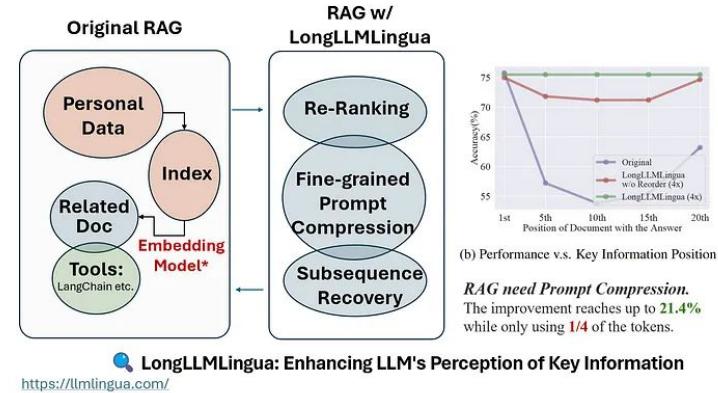
“Lost in the middle”
Problems.





4. Context is there, but not extracted by the LLM

Solution: Prompt Compression
(LongLLMLingua)



Solution: LongContextReorder

Retrieved Set (Reverse Order)

Node 1: (0.98)
Node 2: (0.93)
Node 3: (0.84)
Node 4 (0.81)
Node 5
...
...
Node N

Long Context Re-order

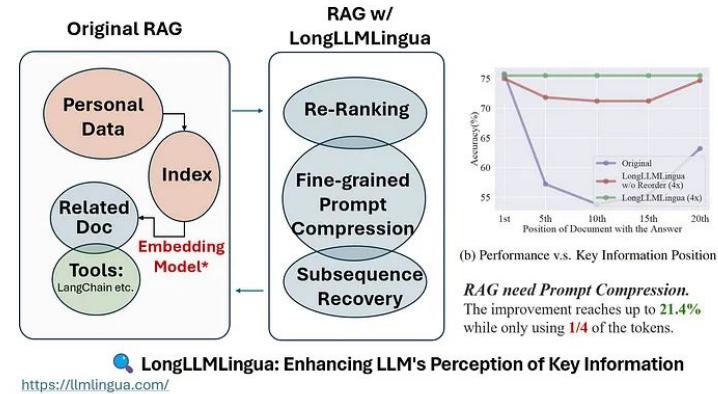
Node 1: (0.98)
Node 3: (0.84)
Node 5
...
Node N
...
Node 6
Node 4 (0.81)
Node 2: (0.93)

The ends matter the most



4. Context is there, but not extracted by the LLM

Solution: Prompt Compression
(LongLLMLingua)



Solution: LongContextReorder

Retrieved Set (Reverse Order)

Node 1: (0.98)
Node 2: (0.93)
Node 3: (0.84)
Node 4 (0.81)
Node 5
...
...
Node N

Long Context Re-order

Node 1: (0.98)
Node 3: (0.84)
Node 5
...
Node N
...
Node 6
Node 4 (0.81)
Node 2: (0.93)

The ends matter the most



5. Output is in Wrong Format

A lot of use cases require outputting the answer in JSON format.

Solutions:

Better text prompting/output parsing

Use OpenAI function calling + JSON mode

Use token-level prompting (LMQL, Guidance)

The following is a character profile for an RPG game in JSON format.

```
```json
{
 "description": "A quick and nimble fighter.",
 "name": "Ranger",
 "age": 20,
 "armor": "plate",
 "weapon": "sword",
 "class": "fighter",
 "mantra": "I am the ranger.",
 "strength": 10,
 "items": [
 "dagger",
 "shield",
 "bow",
]
}```
```

Source: Guidance



## 7. Incomplete Answer

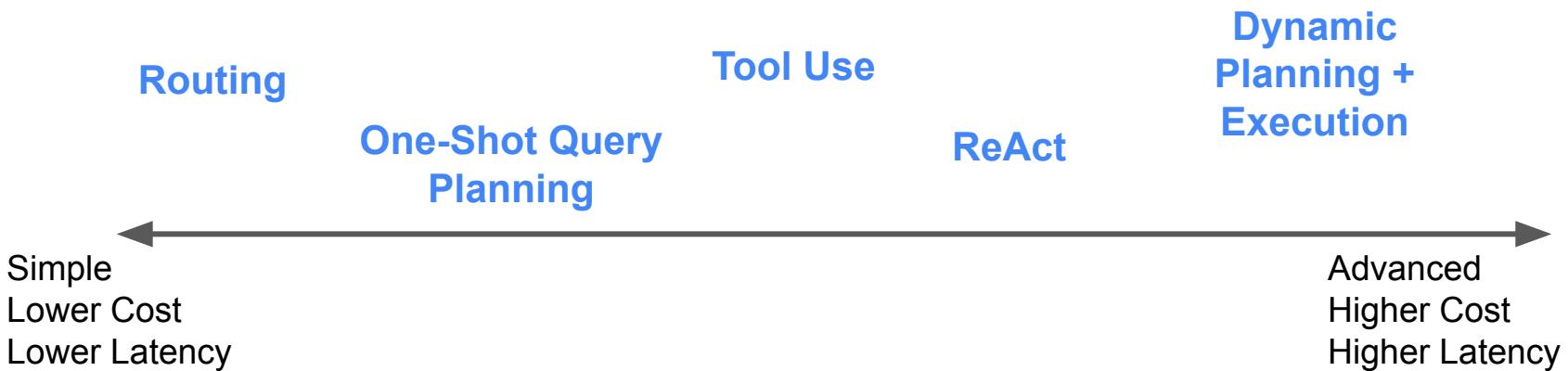
What if you have a complex multi-part question?

Naive RAG is primarily good for answering simple questions about specific facts.



## 7. Incomplete Answer

**Solution:** Add Agentic Reasoning





## 8. Scaling your Data Pipeline

Pain points:

- Processing thousands/millions of docs is slow
- How do we efficiently handle document updates?



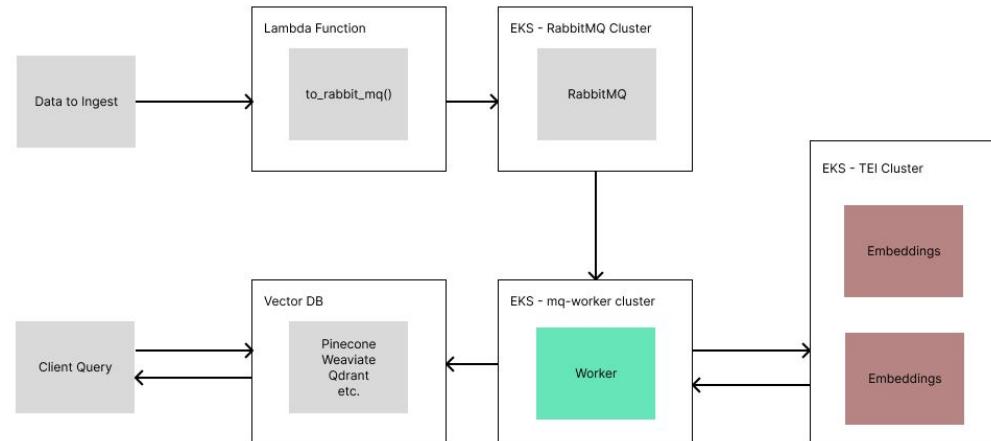
# 8. Scaling your Data Pipeline

Pain points:

- Processing thousands/millions of docs is slow
- How do we efficiently handle document updates?

Reference Production Ingestion Stack

- Parallelize document processing
- HuggingFace TEI
- RabbitMQ Message Queue
- AWS EKS clusters



[https://github.com/run-llama/llamaindex\\_aws\\_ingestion](https://github.com/run-llama/llamaindex_aws_ingestion)



# 10. Proper RAG over Complex Documents



# Advanced Retrieval: Embedded Tables

How do we model complex docs  
with embedded tables?

RAG with naive chunking +  
retrieval → leads to hallucinations!

Embedded Table →

## Annual rankings

The rankings are published annually in March, so the net worths listed are snapshots taken at that time. These lists only show the top 10 wealthiest billionaires for each year.

### Legend

Icon	Description
—	Has not changed from the previous ranking.
▲	Has increased from the previous ranking.
▼	Has decreased from the previous ranking.

### 2023

In the 37th annual *Forbes* list of the world's billionaires, the list included 2,640 billionaires with a total net wealth of \$12.2 trillion, down 28 members and \$500 billion from 2022. Over half of the list is poorer than the previous year, including [Elon Musk](#), who fell from No. 1 to No. 2.<sup>[2]</sup> The list also marks for the first time a French citizen was in the top position as well as a non-American for the first time since 2013 when the Mexican [Carlos Slim Helú](#) was the world's richest person. The list, like in 2022, counted 15 under 30 billionaires with the richest of them being [Red Bull](#) heir [Mark Mateschitz](#) with a net worth of \$34.7 billion. The youngest of the lot were Clemente Del Vecchio, heir to the [Luxottica](#) fortune shared with his six siblings and stepmother, and Kim Jung-yang, whose fortune lies in Japanese-South Korean gaming giant [Nexon](#), both under-20s.<sup>[10]</sup>

No. ♦	Name ♦	Net worth (USD) ♦	Age ♦	Nationality ♦	Primary source(s) of wealth ♦
1 ▲	Bernard Arnault & family	\$211 billion ▲	74	🇫🇷 France	LVMH
2 ▼	Elon Musk	\$180 billion ▼	51	🇺🇸 United States	Tesla, SpaceX
3 ▼	Jeff Bezos	\$114 billion ▼	59	🇺🇸 United States	Amazon
4 ▲	Larry Ellison	\$107 billion ▲	78	🇺🇸 United States	Oracle Corporation
5 —	Warren Buffett	\$106 billion ▼	92	🇺🇸 United States	Berkshire Hathaway
6 ▼	Bill Gates	\$104 billion ▼	67	🇺🇸 United States	Microsoft
7 ▲	Michael Bloomberg	\$94.5 billion ▲	81	🇺🇸 United States	Bloomberg L.P.
8 ▲	Carlos Slim & family	\$93 billion ▲	83	🇲🇽 Mexico	Telmex, América Móvil, Grupo Carso
9 ▲	Mukesh Ambani	\$83.4 billion ▼	65	🇮🇳 India	Reliance Industries
10 ▼	Steve Ballmer	\$80.7 billion ▼	67	🇺🇸 United States	Microsoft

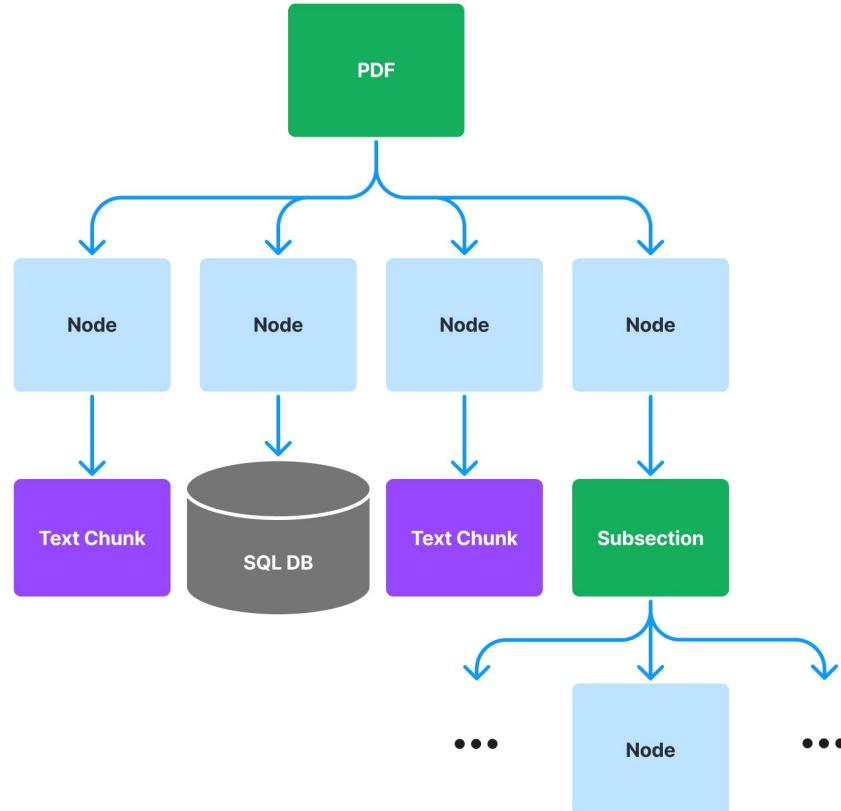


# Advanced Retrieval: Embedded Tables

Instead: model data hierarchically.

Index tables/figures by their summaries.

**The only missing component:**  
how do I parse out the tables from  
the data?





# 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		Life & Savings		Sources
	Ranking	Market share (in %)	Ranking	Market share (in %)	
Main Developed Markets	France	2	12.9	3	8.4 "France Assureurs" as of December 31, 2022. Market share based on statutory premiums and market estimations by SIA (Swiss Insurance Association) figures as of January 31, 2023.
	Switzerland	1	13.3	4	7.8 GDV (German association of Insurance companies) as of December 31, 2021.
	Germany	6	4.8	8	3.4 Assuralia (Belgium Professional Union of Insurance companies) based on gross written premium as of September 30, 2022.
	Belgium	1	17.7	4	8.7 UK General Insurance: Competitor Analytics 2021, Global Data, n/a as of December 31, 2021.
	United Kingdom	4	8.2	n/a	n/a as of December 31, 2021.
	Ireland	1	31.9	n/a	n/a Insurance Ireland P&C Statistics 2021 as of December 31, 2021.
	Spain	5	4.9	9	3.1 Spanish Association of Insurance Companies. ICEA as of January 31, 2022.
	Italy	5	5.8	9	3.9 Associazione Nazionale Imprese Assicuratrici (ANIA) as of December 31, 2021.
	Japan	13	0.6	9	5.0 Disclosed financial reports (excluding Kampo Life) for the 12 months ended September 30, 2022.
	Hong Kong	1	7.0	7	5.0 Insurance Authority statistics based on gross written premiums as of September 30, 2022.
Main Emerging Markets	XL Insurance in the United States	16	1.8	n/a	AM Best 2021 as of December 31, 2021, in the United States in Commercial lines.
	XL Reinsurance worldwide	14	2.3	n/a	n/a AM Best 2021 as of December 31, 2021.
	Thailand	18	1.8	5	7.2 TGIA (Thai General Insurance Association) as of December 31, 2022 and TLLA (Thai Life Assurance Association) as of November 30, 2022.
	Indonesia	n/a	n/a	2	8.7 AAJI Statistic measured on Weighted New Business Premium as of September 30, 2022.
	Philippines	n/a	n/a	6	8.6 Insurance Commission measured on total premium income as of June 30, 2022.
(a)	China	n/a	0.4	n/a	n/a CBIRC (China Banking and Insurance Regulatory Commission) as of December 31, 2022 <sup>16</sup> .
	Mexico	3	8.0	12	2.0 AMIS (Asociación Mexicana de Instituciones de Seguros) as of September 30, 2022.
	Brazil	15	1.4	n/a	n/a SUSEP (Superintendência de Seguros Privados) as of September 2022.
	(a) For Property & Casualty insurance market, CBIRC did not disclose information on ranking. For Life & Savings insurance market, CBIRC did not disclose information on market shares and ranking.				

## PyPDF

1 Please find below AXA's rankings and market shares in the main countries where it operates:  
 2 Property & Casualty Life & Savings  
 3 Market  
 4 share  
 5 (in %) Market  
 6 share  
 7 (in %) Ranking Ranking Sources  
 8 France 2 12.9 3 8.4 "France Assureurs" as of December 31, 2022.  
 9 Market share based on statutory premiums and market  
 10 estimations by SIA (Swiss Insurance Association) figures  
 11 as of January 31, 2023. Switzerland 1 13.3 4 7.8  
 12 GOV (German association of Insurance companies)  
 13 as of December 31, 2021. Germany 6 4.8 8 3.4  
 14 Assuralia (Belgium Professional Union of Insurance  
 15 companies) based on gross written premium  
 16 as of September 30, 2022.s t e k. Mar topped Main Dev Belgium 1 17.7 4 8.7  
 17 UK Genera l Insurance: Competitor Analytics 2021, Global Data,  
 18 as of December 31, 2021. United Kingdom 4 8.2 n/a n/a  
 19 Ireland 1 31.9 n/a n/a Insurance Ireland P&C Statistics 2021 as of December 31, 2021.  
 20 Spanish Association of Insurance Companies. ICEA  
 21 as of December 31, 2022. Spain 5 4.9 9 3.1  
 22 Associazione Nazionale Imprese A ssicuratrici (ANIA)  
 23 as of December 31, 2021. Italy 5 5.8 9 3.9  
 24 Disclosed financial r eports (ex cluding Kampo Life)  
 25 for the 12 months ended September 30, 2022. Japan 13 0.6 9 5.0  
 26 Insur ance Authority statistics based on gross written premiums  
 27 as of September 30, 2022. Hong Kong 1 7.0 7 5.0  
 28 XL Insurance in  
 29 the United S tates AM Best 2021 as of December 31, 2021, in the United States in Commercial lines. 1  
 30 6 1.8 n/a n/a  
 31 XL Reinsurance worldwide 14 2.3 n/a AM Best 2021 as of December 31, 2021.  
 32 Thailand 18 1.8 5.7 2 TGIA (Thai General Insurance Association) as of December 31, 2022 and TL  
 33 AA (Thai Life  
 34 Assurance Association) as of November 30, 2022. ts e rk ing Ma g Emer Main  
 35 Indonesia n/a n/a 2 8.7 AAJI Statistic measured on Weighted New Business Premium as of Sep tember 30, 2022.  
 36 Philippines n/  
 37 a n/a 6 8.6 Insurance Commission measured on total premium income as of June 30, 2022.  
 38 China n/a 0.4 n/a CBIRC (China Banking and Insurance Regulatory Commission) as of Dec ember 31, 2022  
 39 (a).  
 40 Mexico 3 8.0 12 2.0 AMIS (Asociación Mexicana de Instituciones de Seguros) as of Sept ember 30, 2022.  
 41 Brazil 15 1.4 n/a n/a SUSEP (Superintendência de Seguros Privados) as of September 2022.



# Introducing LlamaParse

A genAI-native parser  
designed to let you build  
RAG over complex  
documents

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

The screenshot shows the GitHub repository page for `llama_parse`. The repository is public and has 5 branches and 0 tags. The main branch has 39 commits from `logan-markewich`. Recent commits include:

- even more rename (44 minutes ago)
- rename (47 minutes ago)
- rename (47 minutes ago)
- remove extra files (2 days ago)
- Initial commit (2 days ago)
- rename (47 minutes ago)
- add print statement to print job\_id for easier debugging (1 hour ago)
- rename (47 minutes ago)

The repository details sidebar includes:

- About**: Parse files for optimal RAG, www.llamaindex.ai, Readme, MIT license, Activity, Custom properties, 2 stars, 4 watching, 0 forks, Report repository.
- Releases**: No releases published, Create a new release.
- Packages**: No packages published, Publish your first package.
- Contributors**: 4 contributors: logan-markewich, hatianzhang, sourabhdesai, jerryjiu.

The repository page also features a preview section for LlamaParse, stating it's an API for efficient file parsing and representation for retrieval and context augmentation using LlamaIndex frameworks. It integrates with LlamaIndex and is currently available in preview mode for free. It supports PDF files only.



# Introducing LlamaParse

## Capabilities

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

Screenshot of the GitHub repository for `llama_parse`. The repository has 5 branches and 0 tags. The main branch shows several commits by `logan-markewich` including renames and initial commits. The repository page includes sections for About, Releases, Packages, and Contributors.

**About**  
Parse files for optimal RAG  
[www.llamaindex.ai](#)  
Readme  
MIT license  
Activity  
Custom properties  
2 stars  
4 watching  
0 forks  
Report repository

**Releases**  
No releases published  
[Create a new release](#)

**Packages**  
No packages published  
[Publish your first package](#)

**Contributors** 4

Avatar	Name	Role
	<a href="#">logan-markewich</a>	Logan
	<a href="#">hatianzhang</a>	Haotian Zhang
	<a href="#">sourabhdesai</a>	Sourabh Desai
	<a href="#">jerryjiu</a>	Jerry Liu

**llama\_parse** Public

main 5 Branches 0 Tags

Go to file Add file Code

logan-markewich even more rename a84bec3 - 44 minutes ago 39 Commits

examples even more rename 44 minutes ago

llama\_parse rename 47 minutes ago

tests rename 47 minutes ago

.gitignore remove extra files 2 days ago

LICENSE Initial commit 2 days ago

README.md rename 47 minutes ago

poetry.lock add print statement to print job\_id for easier debugging 1 hour ago

pyproject.toml rename 47 minutes ago

README MIT license

### 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](#).

Currently available in preview mode for free. Try it out today!

NOTE: Currently, only PDF files are supported.

#### Getting Started

# Current PDFReader



# Llama Parse

	Asian	Mexican	Muslim	Physical disability	Jewish	Middle Eastern	Chinese	Mental disability	Latino	Native American	Women	Black	LGBTQ
<b>Pretrained</b>													
MPT	7B	0.53	0.34	0.25	0.29	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	30B	0.38	0.28	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
Falcon	7B	0.54	0.35	0.26	0.22	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	40B	0.52	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
LLAMA 1	7B	0.60	0.47	0.38	0.30	0.22	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	13B	0.52	0.45	0.35	0.28	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	33B	0.52	0.45	0.35	0.28	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	65B	0.52	0.45	0.35	0.28	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
LLAMA 2	7B	0.63	0.41	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	13B	0.52	0.41	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	34B	0.52	0.41	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	70B	0.52	0.41	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
<b>Fine-tuned</b>													
ChatGPT		0.53	0.34	0.25	0.29	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
MPT-instruct	7B	0.53	0.34	0.25	0.29	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
Falcon-instruct	7B	0.52	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
LLAMA 2-CHAT	7B	0.51	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	13B	0.51	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	34B	0.51	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21
	70B	0.51	0.33	0.25	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21	0.21

Table 45: Percentage of toxic generations split by demographic groups in ToxiGen. A small percentage indicates low toxicity in model generations. Demographic group labels are adopted from ToxiGen.

	Asian Americans	African Americans	European Americans	Hispanic and Latino Americans	
<b>Pretrained</b>					
MPT	7B	0.53	0.34	0.25	0.29
	30B	0.38	0.28	0.21	0.21
Falcon	7B	0.54	0.35	0.26	0.21
	40B	0.52	0.33	0.26	0.21
LLAMA 1	7B	0.41	0.32	0.28	0.21
	13B	0.40	0.32	0.26	0.21
	33B	0.38	0.32	0.26	0.21
	65B	0.41	0.34	0.27	0.21
LLAMA 2	7B	0.45	0.33	0.27	0.21
	13B	0.42	0.31	0.28	0.21
	34B	0.40	0.31	0.28	0.21
	70B	0.42	0.34	0.28	0.21
<b>Fine-tuned</b>					
ChatGPT		0.18	0.16	0.15	0.15
MPT-instruct	7B	0.32	0.32	0.29	0.21
Falcon-instruct	7B	0.40	0.34	0.30	0.21
LLAMA 2-CHAT	7B	0.55	0.43	0.39	0.29
	13B	0.51	0.40	0.38	0.29
	34B	0.46	0.40	0.35	0.29
	70B	0.51	0.43	0.39	0.29

Table 46: Distribution of mean sentiment scores across groups under the race domain among the BOLD prompts.

	Asian	Mexican	Muslim	Physical disability	Jewish	Middle Eastern	Chinese	Mental disability	Latino	Native American	Women	Black	LGBTQ
<b>Pretrained</b>													
MPT	7B	0.50	0.35	0.24	0.19	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	30B	0.54	0.39	0.24	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
Falcon	7B	0.05	0.20	0.24	0.29	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	40B	0.59	0.31	0.24	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
LLAMA 1	7B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	13B	0.60	0.32	0.25	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	33B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	65B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
LLAMA 2	7B	0.65	0.31	0.24	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	13B	0.60	0.32	0.25	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	34B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	70B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
<b>Fine-tuned</b>													
ChatGPT		0.03	0.22	0.08	0.09	0.06	0.06	0.06	0.07	0.07	0.07	0.07	0.07
MPT-instruct	7B	0.56	0.27	0.11	0.04	0.04	0.03	0.03	0.03	0.03	0.03	0.03	0.03
Falcon-instruct	7B	0.63	0.35	0.22	0.19	0.19	0.19	0.19	0.19	0.19	0.19	0.19	0.19
LLAMA 2-CHAT	7B	0.65	0.31	0.24	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	13B	0.60	0.32	0.25	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	34B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
	70B	0.65	0.37	0.28	0.21	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20

Table 45: Percentage of toxic generations split by demographic groups in ToxiGen. A small percentage indicates low toxicity in model generations. Demographic group labels are adopted from ToxiGen.

	Asian Americans	African Americans	European Americans	Hispanic and Latino Americans	
<b>Pretrained</b>					
MPT	7B	0.38	0.34	0.25	0.39
	30B	0.38	0.32	0.23	0.33
Falcon	7B	0.36	0.29	0.26	0.47
	40B	0.36	0.29	0.26	0.48
LLAMA 1	7B	0.41	0.32	0.28	0.46
	13B	0.40	0.32	0.26	0.45
	33B	0.39	0.32	0.26	0.46
	65B	0.41	0.34	0.27	0.44
LLAMA 2	7B	0.38	0.33	0.27	0.43
	13B	0.42	0.31	0.28	0.45
	34B	0.40	0.34	0.28	0.42
	70B	0.42	0.34	0.28	0.52
<b>Fine-tuned</b>					
ChatGPT		0.18	0.16	0.15	0.19
MPT-instruct	7B	0.38	0.32	0.29	0.32
Falcon-instruct	7B	0.40	0.34	0.30	0.36
LLAMA 2-CHAT	7B	0.55	0.43	0.40	0.49
	13B	0.51	0.40	0.38	0.49
	34B	0.46	0.40	0.35	0.39
	70B	0.51	0.43	0.40	0.49

Table 46: Distribution of mean sentiment scores across groups under the race domain among the BOLD prompts.



# LlamaParse Results

The best parser at table extraction == the only parser for advanced RAG

Expanded: <https://drive.google.com/file/d/1fyQAg7nOtChQzhF2Ai7HEeKYYqdeWsdt/view?usp=sharing>

## LlamaParse

PvPDF

Apple Inc. CONDENSED COMPARATIVE STATEMENTS OF OPERATIONS (Unaudited)					
	Year Ended December 29,		Year Ended December 30,		Year Ended December 31,
	2012	2011	2010	2009	2008
<b>Net sales:</b>					
Products	\$ 87,181	\$ 76,306	\$ 50,082	\$ 37,761	\$ 26,747
Services	10,000	10,000	10,000	10,000	10,000
Total net sales <sup>(1)</sup>	<b>97,181</b>	<b>86,306</b>	<b>60,082</b>	<b>47,761</b>	<b>36,747</b>
Cost of sales:					
Products	62,000	52,000	35,000	27,000	19,000
Services	8,455	8,985	10,000	10,000	10,000
Total cost of sales	<b>70,455</b>	<b>60,985</b>	<b>45,000</b>	<b>37,000</b>	<b>29,000</b>
Gross margin	<b>26,726</b>	<b>25,321</b>	<b>15,082</b>	<b>10,761</b>	<b>7,747</b>
 <b>Operating expenses:</b>					
Research and development	7,107	6,764	5,847	4,241	3,000
Sales, general and administrative	23,955	21,717	18,133	12,720	8,747
Total operating expenses	<b>31,062</b>	<b>28,481</b>	<b>23,980</b>	<b>16,961</b>	<b>11,747</b>
 <b>Operating income:</b>					
Net income from operations	<b>26,663</b>	<b>26,840</b>	<b>11,102</b>	<b>7,800</b>	<b>5,000</b>
Other income, net, less:					
Interest expense, net	3,425	3,655	2,050	1,625	1,000
Interest income, net	(1,000)	(1,000)	(1,000)	(1,000)	(1,000)
Provision for uncertain tax positions	16,945	16,945	10,000	7,000	5,000
Gain on sale of business	(2,612)	(2,612)	(2,741)	(2,741)	(2,741)
Net income	<b>8,203</b>	<b>7,383</b>	<b>4,362</b>	<b>2,075</b>	<b>1,259</b>
 <b>Income per share:</b>					
Basic	\$ 1.47	\$ 1.28	\$ 0.70	\$ 0.55	\$ 0.47
Diluted	\$ 1.49	\$ 1.29	\$ 0.70	\$ 0.55	\$ 0.47
Shares used in computing basic earnings per share					
Basic	59,129,044	49,055,000	37,070,031	27,757,000	19,250,000
Diluted	60,070,000	59,945,000	45,154,000	35,857,000	25,357,000
 <b>Net cash provided by operating activities:</b>					
Net income	\$ 8,203	\$ 7,383	\$ 4,362	\$ 2,075	\$ 1,259
Depreciation	22,423	22,025	14,244	10,254	7,026
Changes in operating assets and liabilities	(1,544)	(1,544)	(1,000)	(1,000)	(1,000)
Net cash provided by operating activities	<b>29,082</b>	<b>28,383</b>	<b>18,362</b>	<b>12,275</b>	<b>8,259</b>
 <b>Net cash used in investing activities:</b>					
Capital expenditures	55,092	50,250	35,357	25,357	17,357
Acquisitions	(1,000)	(1,000)	(1,000)	(1,000)	(1,000)
Net cash used in investing activities	<b>54,092</b>	<b>49,250</b>	<b>34,357</b>	<b>24,357</b>	<b>16,357</b>
 <b>Net cash provided by financing activities:</b>					
Borrowings	\$ 45,000	\$ 45,000	\$ 30,000	\$ 20,000	\$ 10,000
Repurchases	(30,000)	(30,000)	(20,000)	(15,000)	(10,000)
Dividends	8,648	7,655	5,000	3,000	2,000
Net cash provided by financing activities	<b>23,648</b>	<b>17,655</b>	<b>15,000</b>	<b>12,000</b>	<b>1,000</b>
 <b>Total net cash flow</b>	<b>\$ 8,203</b>	<b>\$ 7,383</b>	<b>\$ 18,362</b>	<b>\$ 12,275</b>	<b>\$ 8,259</b>
 <b>Net cash balance at beginning of period:</b>					
U.S.	\$ 10,000	\$ 10,000	\$ 10,000	\$ 10,000	\$ 10,000
Non-U.S.	—	—	—	—	—
<b>Total net cash balance</b>	<b>\$ 10,000</b>	<b>\$ 10,000</b>	<b>\$ 10,000</b>	<b>\$ 10,000</b>	<b>\$ 10,000</b>

By MuBDE

Appended CONSOLIDATED STATEMENT OF OPERATIONS [Unaudited]					
	For the year ended December 31,				
	2023	2022	2021	2020	2019
<b>Total sales:</b>					
Revenues	\$ 27,978	\$ 26,928	\$ 24,862	\$ 21,371	\$ 19,761
Cost of sales	(20,797)	(19,755)	(18,055)	(15,795)	(14,252)
Total sales (\$)	\$ 7,181	\$ 7,173	\$ 6,807	\$ 5,576	\$ 5,509
<b>Cost of sales:</b>					
Direct costs	42,000	41,000	39,000	35,000	32,000
Selling expenses	8,485	8,984	8,168	7,047	6,282
General overhead	88,771	87,061	78,197	71,997	64,000
Total cost of sales	\$ 111,256	\$ 117,045	\$ 116,265	\$ 104,044	\$ 92,282
<b>Gross margin:</b>					
	\$ 73,875	\$ 70,000	\$ 67,998	\$ 60,960	\$ 59,999
<b>Operating expenses:</b>					
Research and development	7,977	6,231	5,645	5	—
Marketing and sales	10,495	10,000	9,000	7,000	5,000
Net operating expenses	18,472	16,231	14,645	7,000	5,000
<b>Operating income:</b>					
	\$ 55,403	\$ 53,769	\$ 53,353	\$ 60,955	\$ 54,999
<b>Interest expense:</b>					
Interest expense - net	34,000	33,000	32,000	31,000	30,000
Interest income - net	(1,000)	(1,000)	(1,000)	(1,000)	(1,000)
Interest expense before taxes	33,000	32,000	31,000	30,000	30,000
Interest tax benefit	8,600	8,200	7,600	7,000	6,000
Interest expense after taxes	24,400	23,800	23,400	23,000	24,000
<b>Income before taxes:</b>					
	\$ 31,003	\$ 30,969	\$ 30,953	\$ 30,955	\$ 30,999
<b>Taxes:</b>					
State	\$ 1,687	\$ 139	\$ 203	\$ 203	\$ 203
Federal	4,465	3,600	3,600	3,600	3,600
Local and other non-deductible taxes	—	—	—	—	—
Total	6,152	4,738	4,803	4,803	4,803
<b>Income before extraordinary items:</b>					
	\$ 24,851	\$ 26,231	\$ 26,150	\$ 26,152	\$ 26,196
<b>Extraordinary items:</b>					
Amortization of intangible assets	—	—	—	—	—
Loss on sale of business	—	—	—	—	—
Other	—	—	—	—	—
Total extraordinary items	—	—	—	—	—
<b>Income before net income tax:</b>					
	\$ 24,851	\$ 26,231	\$ 26,150	\$ 26,152	\$ 26,196
<b>Income tax expense:</b>					
Federal	6,487	6,067	6,067	6,067	6,067
State	720	620	747	747	747
Local	1,687	1,439	1,439	1,439	1,439
Other	—	—	—	—	—
Total income tax expense	8,904	7,925	8,245	8,245	8,245
<b>Net income:</b>					
	\$ 15,947	\$ 18,306	\$ 17,895	\$ 17,907	\$ 17,951
<b>Non-controlling interest:</b>					
Share of net income	—	—	—	—	—
Change in fair value of non-controlling interest	—	—	—	—	—
Total non-controlling interest	—	—	—	—	—
<b>Net income available for common stockholders:</b>					
	\$ 15,947	\$ 18,306	\$ 17,895	\$ 17,907	\$ 17,951
<b>Dividends:</b>					
Preferred dividends	—	—	—	—	—
Common dividends	—	—	—	—	—
Total dividends	—	—	—	—	—
<b>Net income available for common stockholders after dividends:</b>					
	\$ 15,947	\$ 18,306	\$ 17,895	\$ 17,907	\$ 17,951
<b>Net cash provided by operating activities:</b>					
	\$ 23,000	\$ 22,000	\$ 21,000	\$ 18,000	\$ 16,000
<b>Investment in property, plant and equipment:</b>					
Purchase	—	—	—	—	—
Sale	—	—	—	—	—
Depreciation	—	—	—	—	—
Total investment in property, plant and equipment	—	—	—	—	—
<b>Net cash used in investing activities:</b>					
	\$ 23,000	\$ 22,000	\$ 21,000	\$ 18,000	\$ 16,000
<b>Proceeds from borrowings:</b>					
Bank loans	—	—	—	—	—
Notes payable	—	—	—	—	—
Accounts payable	—	—	—	—	—
Total proceeds from borrowings	—	—	—	—	—
<b>Repayments of borrowings:</b>					
Bank loans	—	—	—	—	—
Notes payable	—	—	—	—	—
Accounts payable	—	—	—	—	—
Total repayments of borrowings	—	—	—	—	—
<b>Net cash used in financing activities:</b>					
	\$ 23,000	\$ 22,000	\$ 21,000	\$ 18,000	\$ 16,000
<b>Net increase (decrease) in cash and cash equivalents:</b>					
	\$ 23,000	\$ 22,000	\$ 21,000	\$ 18,000	\$ 16,000
<b>Cash and cash equivalents at beginning of period:</b>					
	\$ 10,000	\$ 10,000	\$ 10,000	\$ 10,000	\$ 10,000
<b>Cash and cash equivalents at end of period:</b>					
	\$ 33,000	\$ 32,000	\$ 31,000	\$ 28,000	\$ 26,000

Texttract

Category	Commodity Combinations by Type of Operation (Unaudited)		
	2002 (\$ millions)	2001 (\$ millions)	2000 (\$ millions)
Net sales	\$ 61,656	\$ 70,516	\$ 50,025
Product			
Sales	61,656	70,516	50,025
Total net sales	61,656	70,516	50,025
Cost of sales			
Sales	41,866	49,146	31,942
Gross margin	19,789	11,370	8,083
Gross margin %	31.3%	16.1%	16.2%
Selling, general and administrative expenses	4,781	5,657	3,807
Research and development	5,875	6,357	5,625
Depreciation and amortization	1,045	1,250	1,000
Total operating expenses	11,601	13,064	10,472
Operating income	\$ 50,058	\$ 20,452	\$ 19,553
Interest expense			
Interest expense related to senior notes	50,658	33,371	35,161
Interest income			
Interest income related to senior notes	4,023	3,750	3,013
Interest income from bank	1,000	1,000	1,000
Operating cash flow	\$ 51,480	\$ 20,101	\$ 17,561
Capital expenditures			
Capital expenditures	\$ 1,177	\$ 126	\$ 636
Capital expenditures related to capital projects	\$ 1,166	\$ 23	\$ 617
Capital expenditures related to acquisitions			
Capital expenditures	16,992,158	16,101,231	10,141,221
Capital expenditures related to acquisitions	16,979,958	16,074,931	10,116,736
Acquisition			
Acquisition of equity interests	\$ 41,051	\$ 41,051	\$ 50,760
Acquisition of assets	22,443	22,929	45,254
Acquisition of business	1,000	1,000	1,000
Joint venture	5,935	5,700	5,457
Other	1,000	1,000	1,000
Total acquisition	\$ 65,428	\$ 65,812	\$ 50,760
Net cash used by acquisitions			
Purchase	\$ 41,051	\$ 41,051	\$ 50,760
Sale	3,200	3,200	2,627
Net	3,641	3,651	3,833
Charterholders' Equity and Accumulated Earnings			
Charterholders' equity	29,394	36,961	45,000
Accumulated earnings	1,000	1,000	1,000
Total increase	\$ 51,480	\$ 20,101	\$ 17,561

PdfMine

Apple Inc.					
CONDENSED CONSOLIDATED STATEMENT OF OPERATIONS (Unaudited)					
	Three Months Ended		Three Months Ended		Per Share Data
	2023	2022	2023	2022	
<b>Total revenue:</b>					
Products	\$ 87,016	\$ 85,765	\$ 285,332	\$ 277,352	
Services	5,000	4,900	15,000	14,900	
Total revenue <sup>1</sup>	\$ 92,016	\$ 85,665	\$ 300,332	\$ 292,252	
<b>Cost of sales:</b>					
Products	52,600	51,000	159,500	153,500	
Services	5,000	4,900	15,000	14,900	
Total cost of sales	\$ 57,600	\$ 55,900	\$ 174,500	\$ 168,400	
<b>Gross margin:</b>					
	\$ 34,416	\$ 29,765	\$ 125,832	\$ 123,852	
<b>Operating expenses:</b>					
Sales and marketing	5,700	5,700	16,945	16,945	
R&D	10,000	10,000	26,900	26,900	
General, administrative and other	10,000	10,700	26,597	26,597	
Total operating expenses	\$ 25,700	\$ 26,400	\$ 60,442	\$ 60,442	
<b>Operating income:</b>					
Net sales	\$ 66,365	\$ 64,965	\$ 162,890	\$ 156,410	
Cost of sales	(52,600)	(51,000)	(159,500)	(153,500)	
SG&A	(25,700)	(26,400)	(60,442)	(60,442)	
Depreciation and amortization	7,000	7,000	18,500	18,500	
Interest expense, net	(6,000)	(6,000)	(15,000)	(15,000)	
Interest income	5,000	5,000	12,500	12,500	
Other net	1,000	1,000	3,000	3,000	
Total operating income	\$ 10,065	\$ 9,565	\$ 16,148	\$ 16,148	
<b>Non-operating income:</b>					
Interest income	5,000	5,000	12,500	12,500	
Other non-operating income	1,000	1,000	3,000	3,000	
Total non-operating income	\$ 6,000	\$ 6,000	\$ 15,500	\$ 15,500	
<b>Non-operating expense:</b>					
Interest expense	(5,000)	(5,000)	(12,500)	(12,500)	
Other non-operating expense	(1,000)	(1,000)	(3,000)	(3,000)	
Total non-operating expense	\$ (6,000)	\$ (6,000)	\$ (15,500)	\$ (15,500)	
<b>Net income:</b>					
Net income	\$ 14,065	\$ 13,565	\$ 36,148	\$ 36,148	
Less accumulated other comprehensive loss	(1,000)	(1,000)	(2,500)	(2,500)	
Net income	\$ 13,065	\$ 12,565	\$ 33,648	\$ 33,648	
<b>Net income per share:</b>					
Basic	\$ 1.49	\$ 1.28	\$ 0.36	\$ 0.36	
Diluted	\$ 1.49	\$ 1.28	\$ 0.36	\$ 0.36	
<b>Shares outstanding:</b>					
Basic	8,829,545	8,500,333	10,914,231	10,781,231	
Diluted	8,829,545	8,500,333	10,914,231	10,781,231	
<b>Reconciliation of non-GAAP financial measures:</b>					
<b>Net Income:</b>					
Net income	\$ 66,365	\$ 64,965	\$ 162,890	\$ 156,410	
Depreciation and amortization	7,000	7,000	18,500	18,500	
Interest expense, net	(6,000)	(6,000)	(15,000)	(15,000)	
Interest income	5,000	5,000	12,500	12,500	
Other net	1,000	1,000	3,000	3,000	
Total net income	\$ 60,365	\$ 58,965	\$ 136,890	\$ 130,410	
<b>Net Income Per Share:</b>					
Basic	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
Diluted	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
<b>Shares Outstanding:</b>					
Basic	8,829,545	8,500,333	10,914,231	10,781,231	
Diluted	8,829,545	8,500,333	10,914,231	10,781,231	
<b>Reconciliation of non-GAAP financial measures:</b>					
<b>Net Income:</b>					
Net income	\$ 66,365	\$ 64,965	\$ 162,890	\$ 156,410	
Depreciation and amortization	7,000	7,000	18,500	18,500	
Interest expense, net	(6,000)	(6,000)	(15,000)	(15,000)	
Interest income	5,000	5,000	12,500	12,500	
Other net	1,000	1,000	3,000	3,000	
Total net income	\$ 60,365	\$ 58,965	\$ 136,890	\$ 130,410	
<b>Net Income Per Share:</b>					
Basic	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
Diluted	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
<b>Shares Outstanding:</b>					
Basic	8,829,545	8,500,333	10,914,231	10,781,231	
Diluted	8,829,545	8,500,333	10,914,231	10,781,231	
<b>Reconciliation of non-GAAP financial measures:</b>					
<b>Net Income:</b>					
Net income	\$ 66,365	\$ 64,965	\$ 162,890	\$ 156,410	
Depreciation and amortization	7,000	7,000	18,500	18,500	
Interest expense, net	(6,000)	(6,000)	(15,000)	(15,000)	
Interest income	5,000	5,000	12,500	12,500	
Other net	1,000	1,000	3,000	3,000	
Total net income	\$ 60,365	\$ 58,965	\$ 136,890	\$ 130,410	
<b>Net Income Per Share:</b>					
Basic	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
Diluted	\$ 6.81	\$ 6.61	\$ 15.21	\$ 14.81	
<b>Shares Outstanding:</b>					
Basic	8,829,545	8,500,333	10,914,231	10,781,231	
Diluted	8,829,545	8,500,333	10,914,231	10,781,231	



# Steerability

## 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  $x = a$  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  $f'(a)$ .

- 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 by definition.

- 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

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

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

[Without instruction]

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

```
|1.|Let's find the derivative of constant function f(x) = α. The differential coefficient of f(x) at x = a is|
|---|---|
| |lim α → 0 (f(a + ε) - f(a)) / ε = lim α → 0 (α - α) = lim α → 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) = αx + β. The derivative of f(x) at x = a is|
|---|---|
| |lim ε → 0 (f(a + ε) - f(a)) = lim ε → 0 (α(a + ε) + β - (αa + β)) = lim ε → 0 α = α|
| |Thus, the derivative of f(x) is f'(x) = α, 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 = a is|
|---|---|
| |lim ε → 0 ((a + ε)^2 - a^2) / ε = lim ε → 0 (a^2 + 2aε + ε^2 - a^2) = lim ε → 0 (2aε + ε^2) = lim ε → 0 (2a + ε) = 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

- 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.

## ## Chapter 1 Let's Differentiate a Function!

# Default (no instructions)



# Steerability

## 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

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

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  $f'(a)$ .

2. 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 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  $x = a$  is

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

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

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 \$\$)"""

[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\}} 0 = 0 \\
&\text{Thus, the derivative of } f(x) \text{ is } f'(x) = 0. \text{ This makes sense, since our function is constant—the rate of change is 0.}
\end{aligned}
```

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

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

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  $f(x) = x^2$ . The differential coefficient of  $f(x)$  at  $x = a$  is

## With Instructions



# What's next for RAG: Agents?

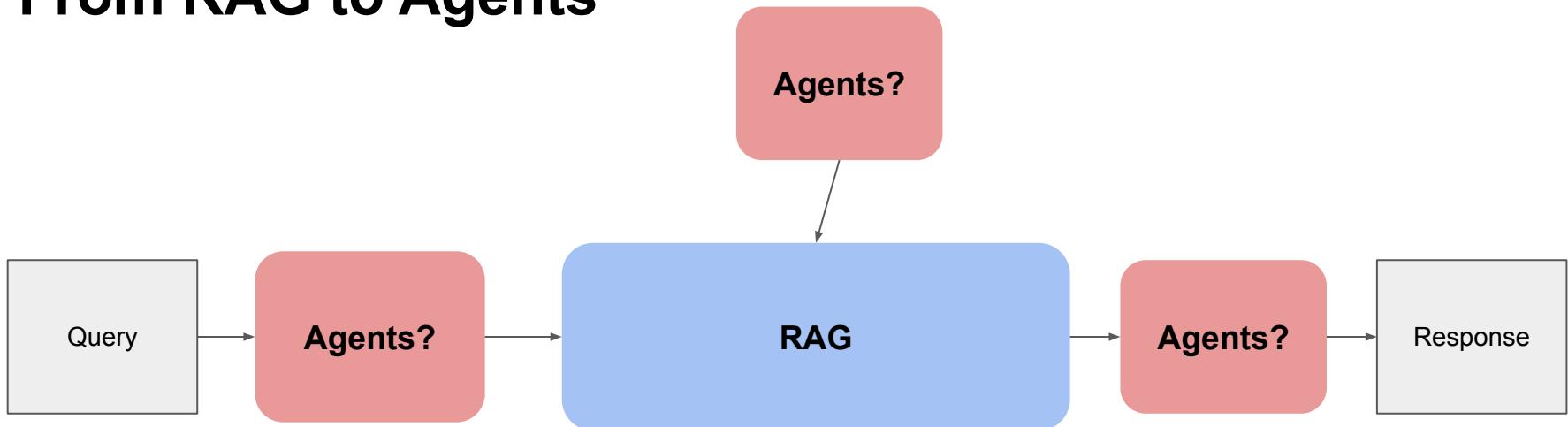


# From RAG to Agents



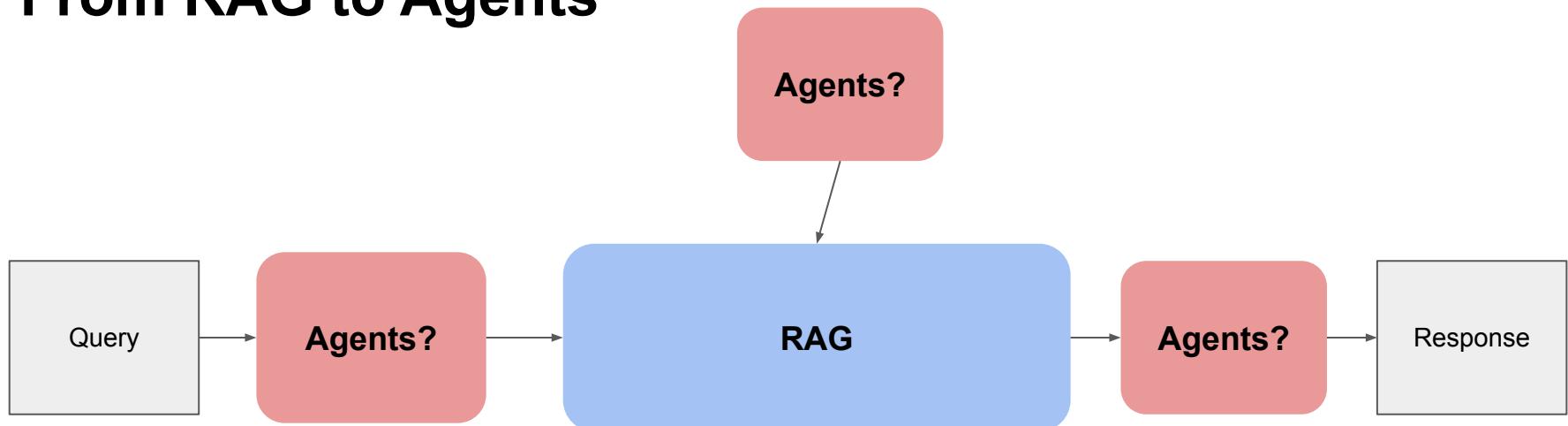


# From RAG to Agents





# From RAG to Agents

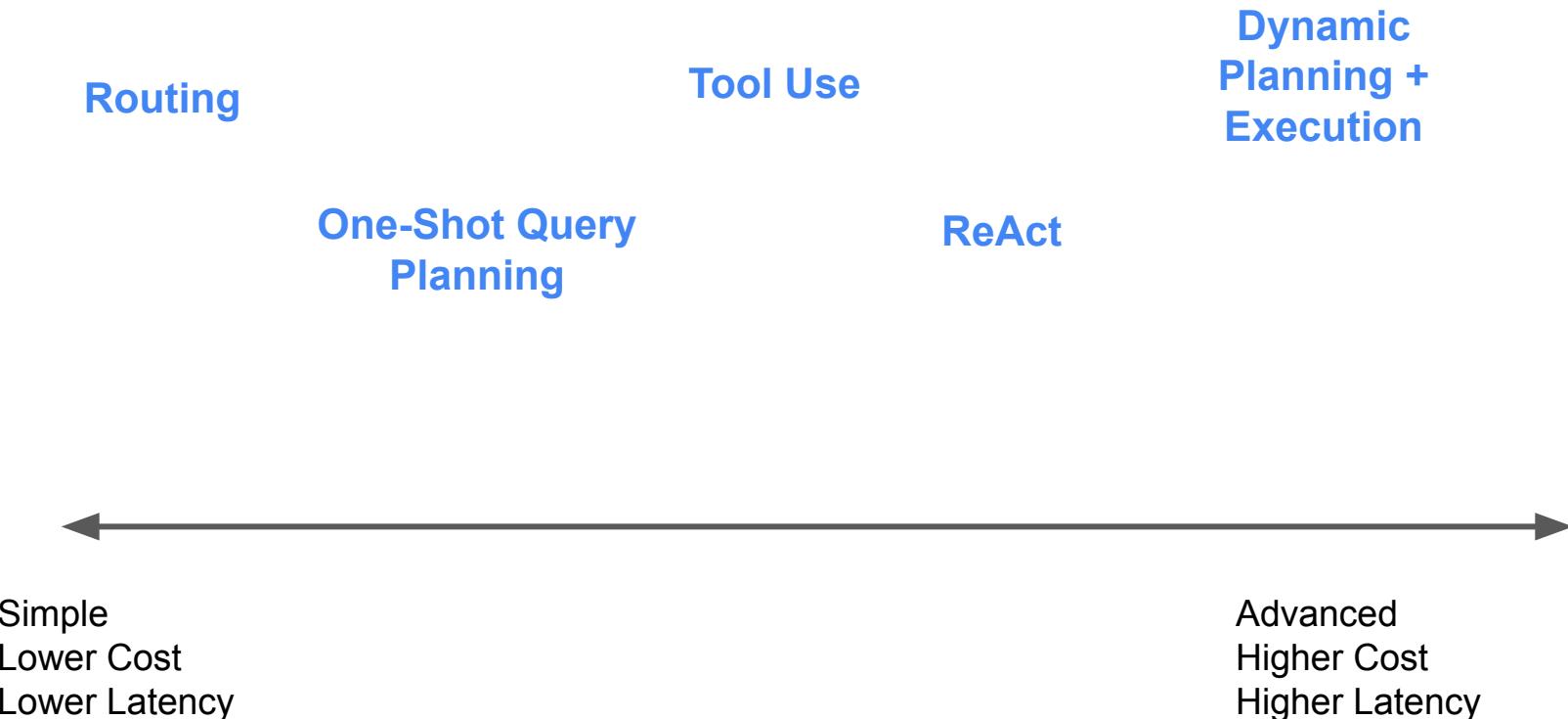


**Agent Definition:** Using LLMs for automated reasoning and tool selection

**RAG is just one Tool:** Agents can decide to use RAG with other tools



# From Simple to Advanced Agents

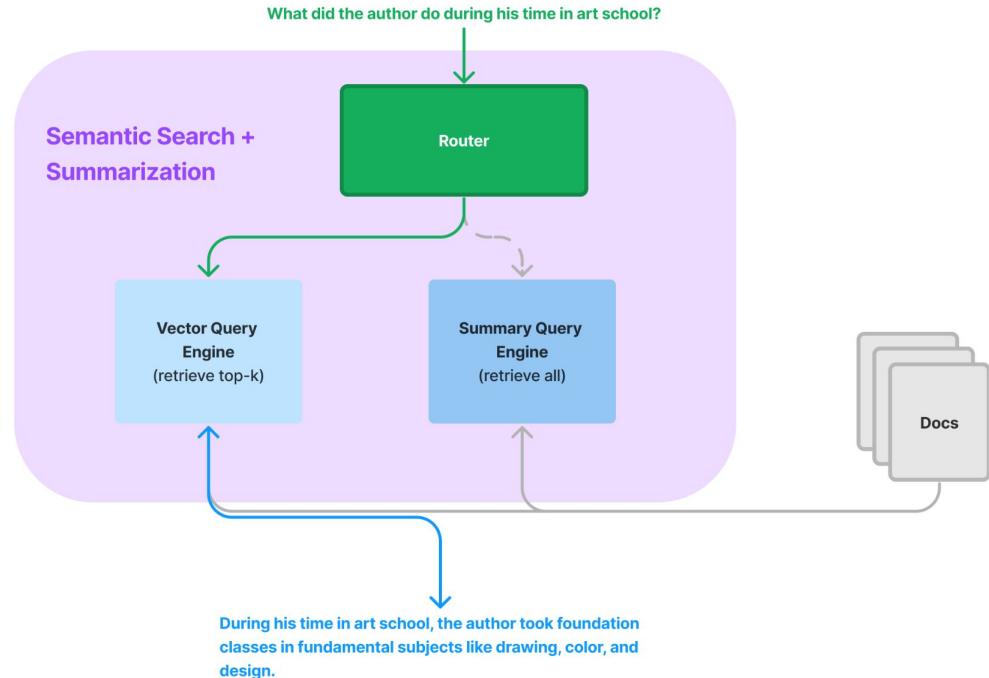




# Routing

Simplest form of agentic reasoning.

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

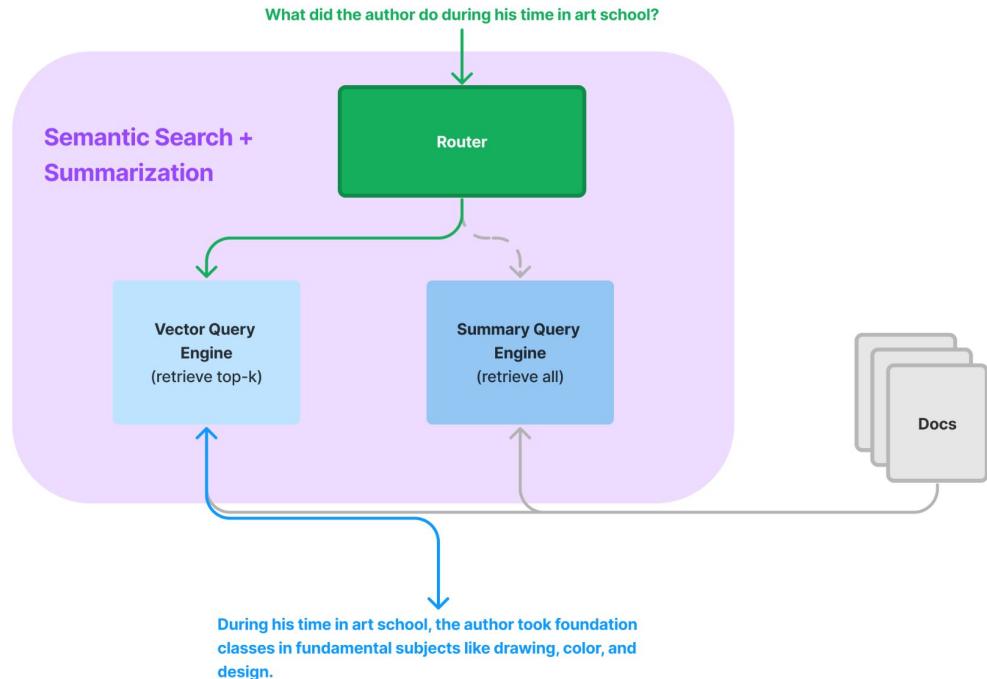




# Routing

**Use Case:** Joint QA and Summarization

Guide

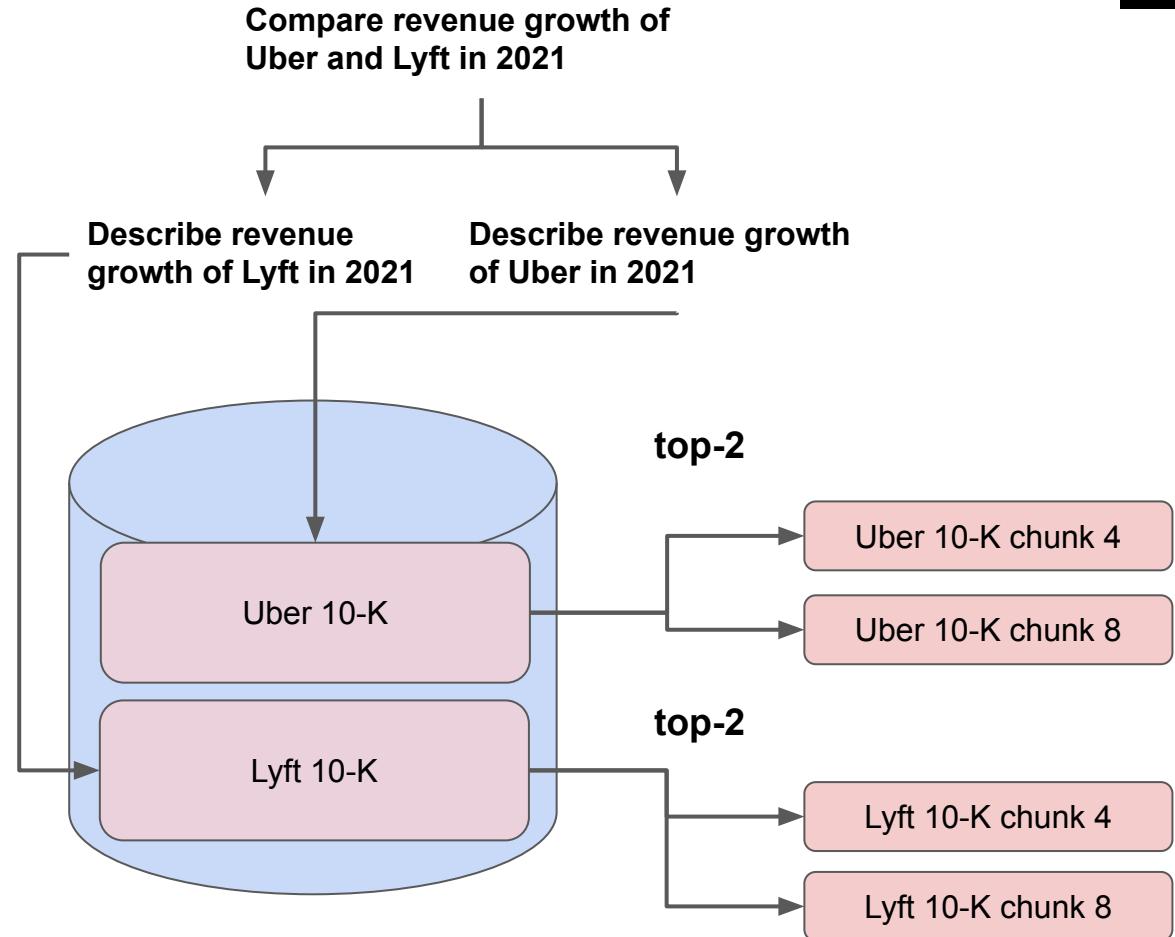




# Query Planning

Break down query into parallelizable sub-queries.

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

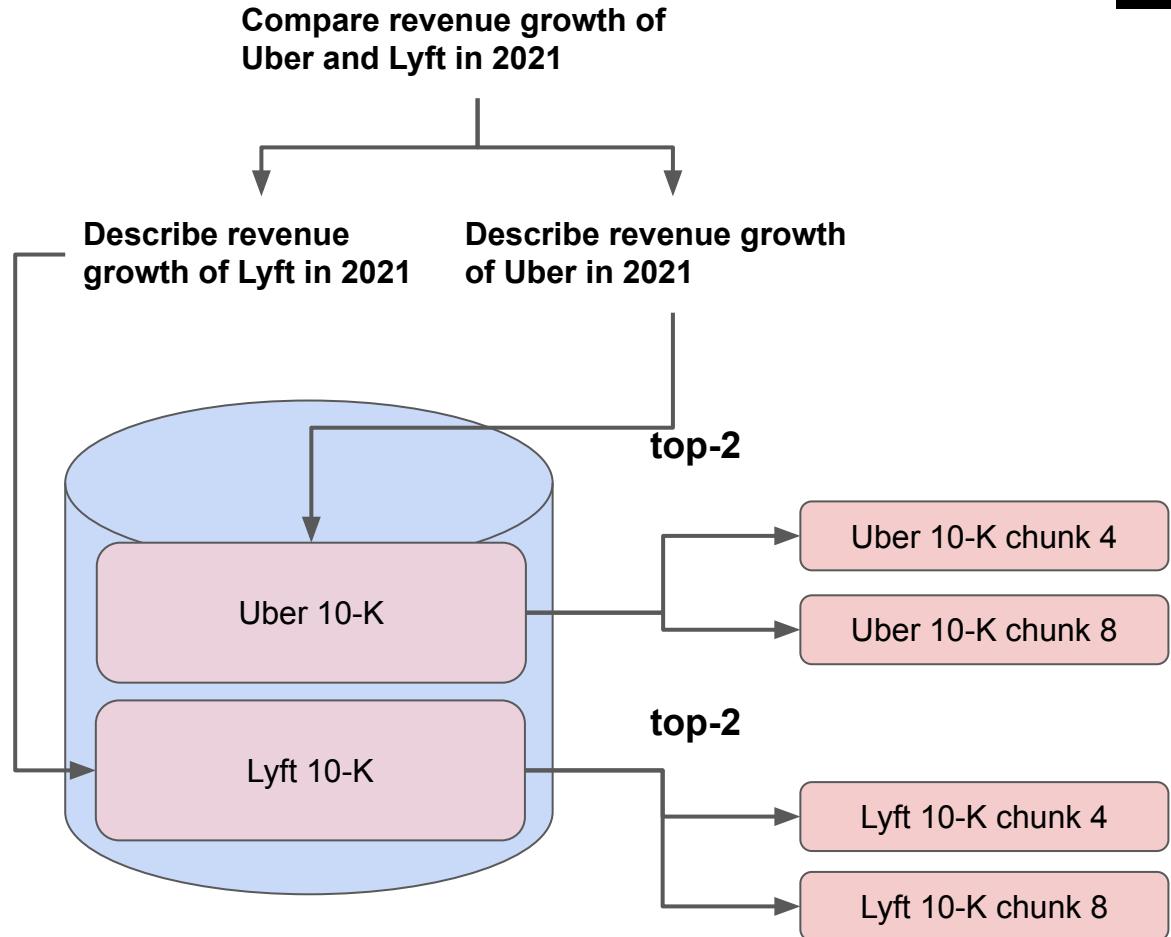




# Query Planning

**Example:** Compare revenue of Uber and Lyft in 2021

[Query Planning Guide](#)



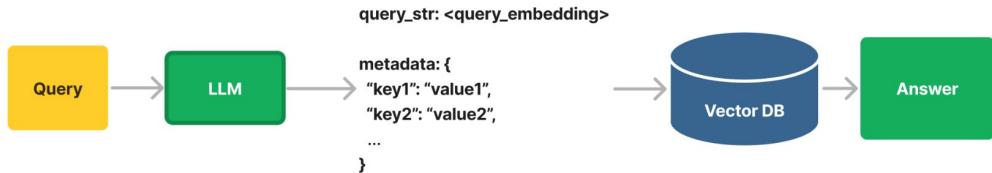


# Tool Use

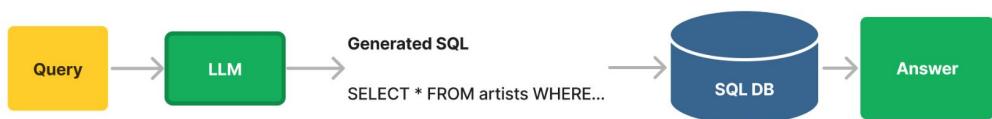
Use an LLM to call an API

Infer the parameters of that API

## Auto-Retrieval



## Text-to-SQL



## Calendar





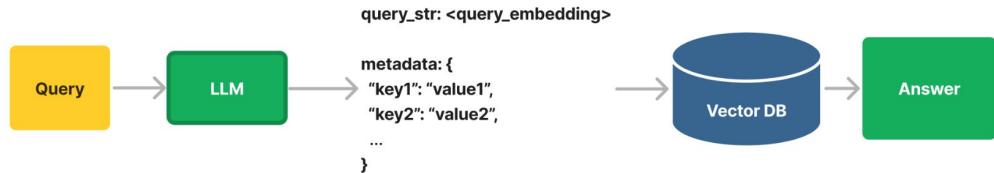
# Tool Use

In normal RAG you just pass through the query.

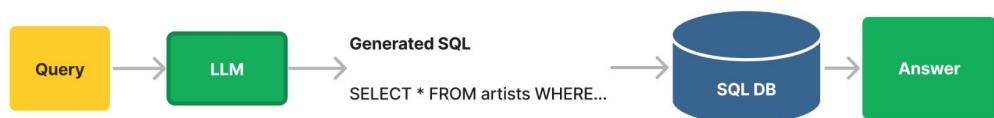
But what if you used the LLM to infer all the parameters for the API interface?

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

## Auto-Retrieval



## Text-to-SQL



## Calendar





# This is cool but

- How can an agent tackle sequential multi-part problems?
- How can an agent maintain state over time?

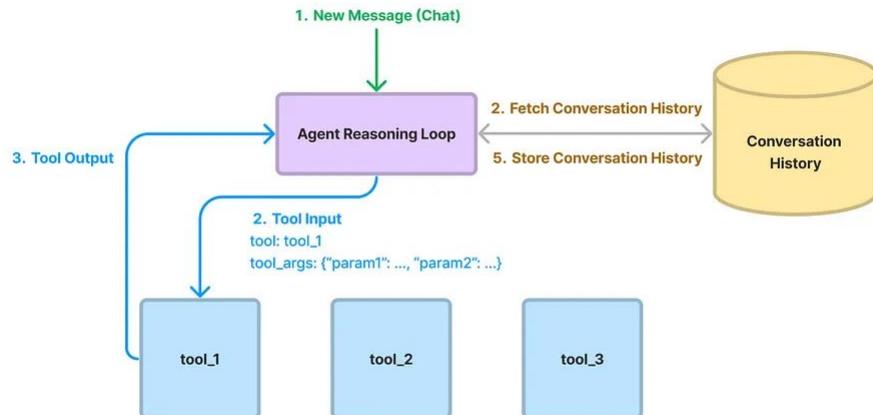


# This is cool but

- How can an agent tackle sequential multi-part problems?
  - Let's make it loop
- How can an agent maintain state over time?
  - Let's add basic memory



# Data Agents - Core Components



## Agent Reasoning Loop

- [ReAct Agent](#) (any LLM)
- [OpenAI Agent](#) (only OAI)

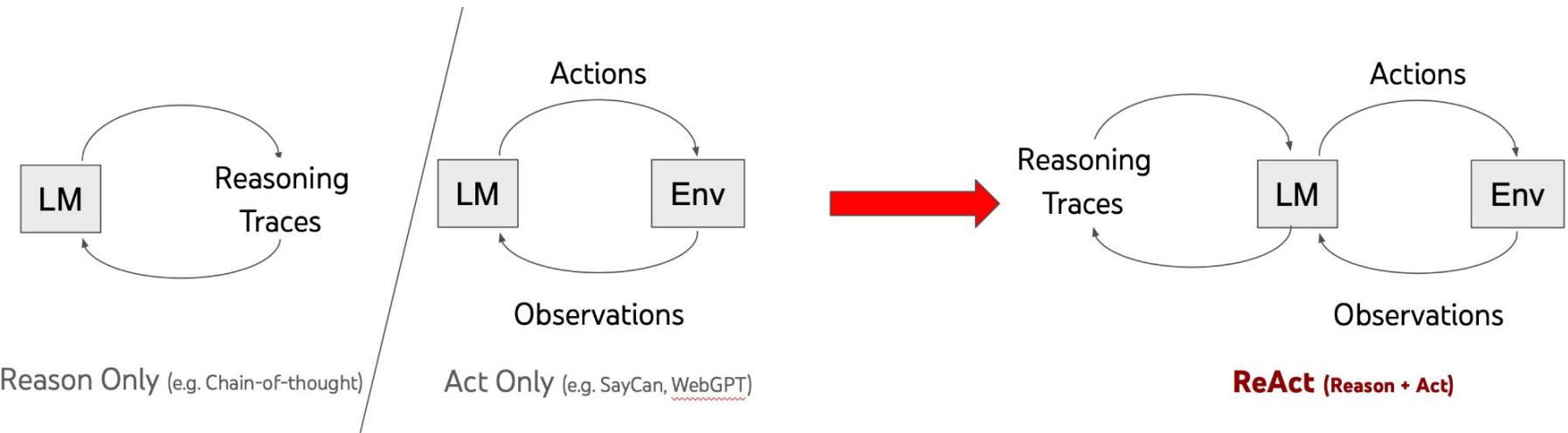
## Tools

[Query Engine Tools \(RAG pipeline\)](#)

[LlamaHub Tools](#) (30+ tools to external services)



# ReAct: Reasoning + Acting with LLMs



Source: <https://react-lm.github.io/>



# ReAct: Reasoning + Acting with LLMs

Add a loop around  
query  
decomposition + tool  
use

```
[]: response = agent.chat(
 "Tell me about the demographics of Houston, and compare that with the demographics of Chicago"
)

Added user message to memory: Tell me about the demographics of Houston, and compare that with the demographics o
f Chicago
==> Calling Function ==
Calling function: vector_tool_Houston with args: {
 "input": "demographics"
}
Got output: Houston has a population of 2,304,580 according to the 2020 U.S. census. In 2017, the estimated popul
ation was 2,312,717, and in 2018 it was 2,325,502. The city has a diverse demographic makeup, with a significant
number of undocumented immigrants residing in the Houston area, comprising nearly 9% of the city's metropolitan p
opulation in 2017. The age distribution in Houston includes a significant number of individuals under 15 and betw
een the ages of 20 to 34. The median age of the city is 33.4. The city has a mix of homeowners and renters, with
an estimated 42.3% of Houstonians owning housing units. The median household income in 2019 was $52,338, and 20.
1% of Houstonians lived at or below the poverty line.
=====

==> Calling Function ==
Calling function: vector_tool_Chicago with args: {
 "input": "demographics"
}
Got output: Chicago experienced rapid population growth during its first hundred years, becoming one of the faste
st-growing cities in the world. From its founding in 1833 with fewer than 200 people, the population grew to over
4,000 within seven years. By 1890, the population had surpassed 1 million, making Chicago the fifth-largest city
in the world at the time. The city's population continued to grow, reaching its highest recorded population of 3.
6 million in 1950. However, in the latter half of the 20th century, Chicago's population declined, dropping to un
der 2.7 million by 2010. The city experienced a rise in population for the 2000 census, followed by a decrease in
2010, and then another increase for the 2020 census. According to U.S. census estimates as of July 2019, the larg
est racial or ethnic groups in Chicago are non-Hispanic White (32.8%), Blacks (30.1%), and Hispanics (29.0%). Add
itionally, Chicago has the third-largest LGBTQ population in the United States, with an estimated 7.5% of the adul
t population identifying as LGBTQ in 2018.
=====
```



# ReAct: Reasoning + Acting with LLMs

Superset of query planning + routing capabilities.

## ReAct + RAG Guide

```
[]: response = agent.chat(
 "Tell me about the demographics of Houston, and compare that with the demographics of Chicago"
)

Added user message to memory: Tell me about the demographics of Houston, and compare that with the demographics o
f Chicago
==> Calling Function ==
Calling function: vector_tool_Houston with args: {
 "input": "demographics"
}
Got output: Houston has a population of 2,304,580 according to the 2020 U.S. census. In 2017, the estimated popul
ation was 2,312,717, and in 2018 it was 2,325,502. The city has a diverse demographic makeup, with a significant
number of undocumented immigrants residing in the Houston area, comprising nearly 9% of the city's metropolitan p
opulation in 2017. The age distribution in Houston includes a significant number of individuals under 15 and betw
een the ages of 20 to 34. The median age of the city is 33.4. The city has a mix of homeowners and renters, with
an estimated 42.3% of Houstonians owning housing units. The median household income in 2019 was $52,338, and 20.
1% of Houstonians lived at or below the poverty line.
=====

==> Calling Function ==
Calling function: vector_tool_Chicago with args: {
 "input": "demographics"
}
Got output: Chicago experienced rapid population growth during its first hundred years, becoming one of the faste
st-growing cities in the world. From its founding in 1833 with fewer than 200 people, the population grew to over
4,000 within seven years. By 1890, the population had surpassed 1 million, making Chicago the fifth-largest city
in the world at the time. The city's population continued to grow, reaching its highest recorded population of 3.
6 million in 1950. However, in the latter half of the 20th century, Chicago's population declined, dropping to un
der 2.7 million by 2010. The city experienced a rise in population for the 2000 census, followed by a decrease in
2010, and then another increase for the 2020 census. According to U.S. census estimates as of July 2019, the larg
est racial or ethnic groups in Chicago are non-Hispanic White (32.8%), Blacks (30.1%), and Hispanics (29.0%). Add
itionally, Chicago has the third-largest LGBTQ population in the United States, with an estimated 7.5% of the adul
t population identifying as LGBTQ in 2018.
=====
```



# Can we make this even better?

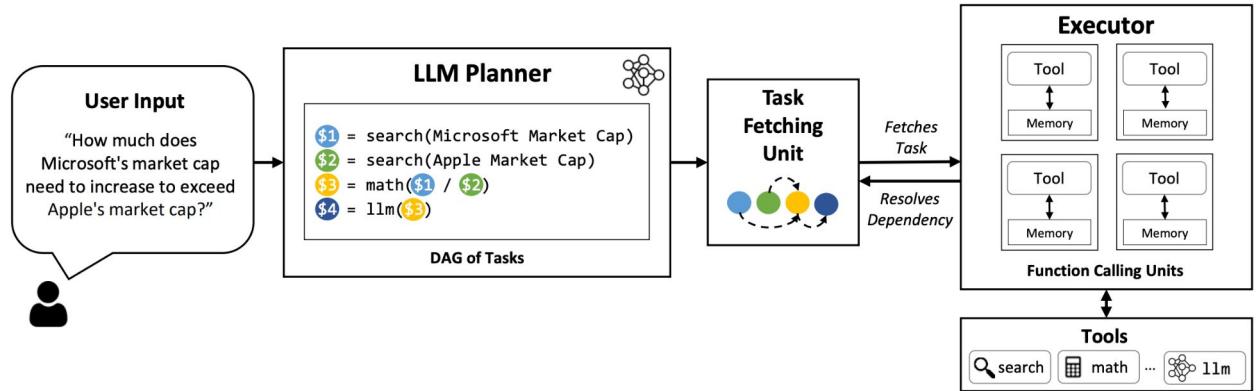
- Stop being so short-sighted - plan ahead at each step
- Parallelize execution where we can



# LLMCompiler

Kim et al. 2023

An agent compiler  
for parallel  
multi-function  
planning +  
execution.



**Figure 2:** Overview of the LLMCompiler framework: the workflow from initial user input to task execution. Beginning with user input, the LLM Planner generates a sequence of tasks with their inter-dependencies. These tasks are then dispatched by the Task Fetching Unit to the Executor based on their dependencies, thus allowing for their parallel executions. For instance, in this example, Task \$1 and \$2 are fetched together for parallel execution of two independent search tasks. After each task is performed, the results (i.e., observations) are forwarded back to the Task Fetching Unit to unblock the dependent tasks after replacing their placeholder variables (e.g., the variable \$1 and \$2 in Task \$3) with actual values. Once all tasks have been executed, the final answer is delivered to the user.



# LLMCompiler

Plan out steps  
beforehand, and  
replan as necessary

## LLMCompiler Agent

```
[17]: response = agent.chat(
 "Is the climate of Chicago or Seattle better during the wintertime?"
)
print(str(response))

> Running step f8fdf4cb-9dde-4aba-996d-edbcee53c4c2 for task 20df27d7-cc27-4311-bb57-b1a6f4ad5799.
> Step count: 0
> Plan: 1. vector_tool_Chicago("climate during wintertime")
2. vector_tool_Seattle("climate during wintertime")
3. join()<END_OF_PLAN>
Ran task: vector_tool_Seattle. Observation: During wintertime, Seattle experiences cool, wet conditions. Extreme cold temperatures, below about 15 °F or -9 °C, are rare due to the moderating influence of the adjacent Puget Sound, the greater Pacific Ocean, and Lake Washington. The city is often cloudy due to frequent storms and lows moving in from the Pacific Ocean, and it has many "rain days". However, the rainfall is often a light drizzle.
Ran task: vector_tool_Chicago. Observation: During wintertime, the city experiences relatively cold and snowy conditions. Blizzards can occur, as they did in winter 2011. The normal winter high from December through March is about 36 °F (2 °C). January and February are the coldest months. A polar vortex in January 2019 nearly broke the city's cold record of -27 °F (-33 °C), which was set on January 20, 1985. Measurable snowfall can continue through the first or second week of April. The city's proximity to Lake Michigan tends to keep the lakefront somewhat cooler in summer and less brutally cold in winter than inland parts of the city and suburbs away from the lake. Northeast winds from wintertime cyclones departing south of the region sometimes bring the city lake-effect snow.
Ran task: join. Observation: None
> Thought: Comparing the two climates, Seattle seems to have a milder winter climate than Chicago.
> Answer: Seattle
Seattle
```



# 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)

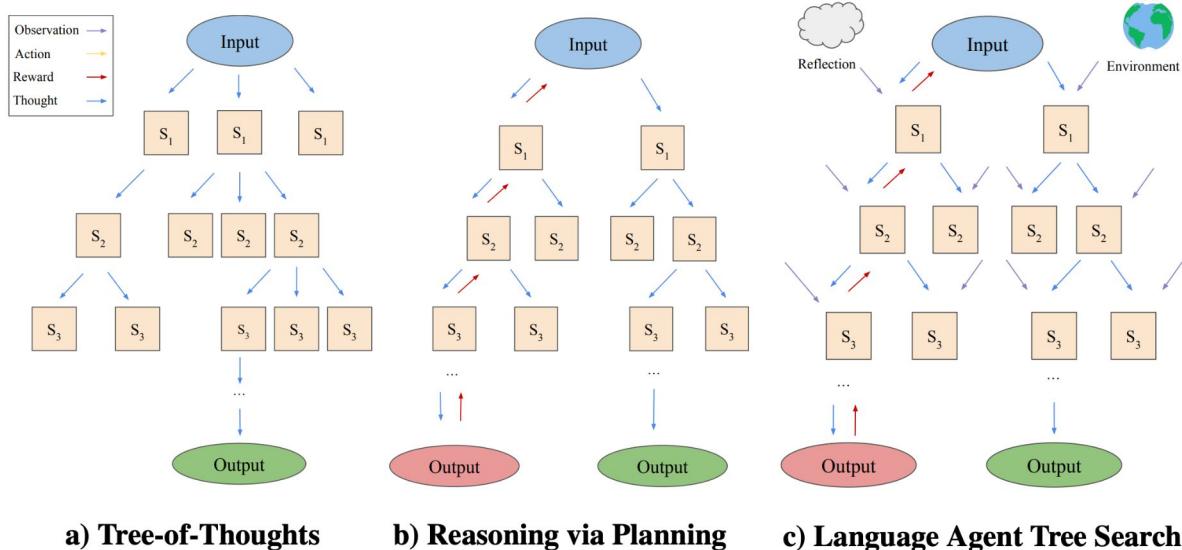


Figure 2: An overview of the differences between LATS and recently proposed LM search algorithms ToT (Yao et al., 2023a) and RAP (Hao et al., 2023). LATS leverages environmental feedback and self-reflection to further adapt search and improve performance.



# Additional Requirements

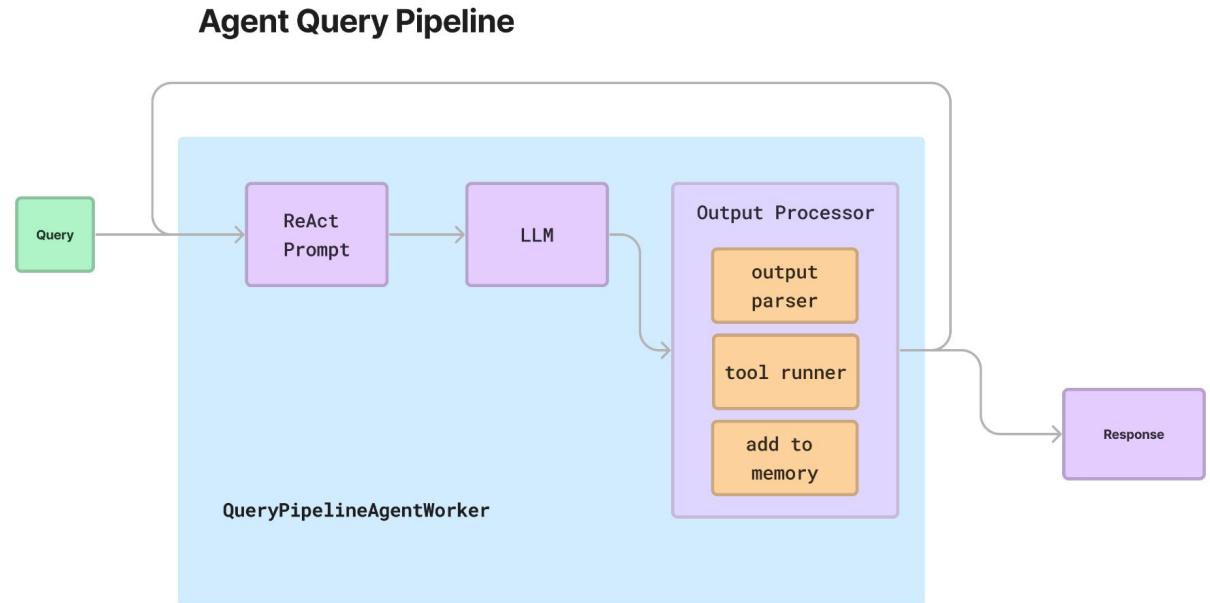
- **Observability:** see the full trace of the agent
  - [Observability Guide](#)
- **Control:** Be able to guide the intermediate steps of an agent *step-by-step*
  - [Lower-Level Agent API](#)
- **Customizability:** Define your own agentic logic around any set of tools.
  - [Custom Agent Guide](#)
  - [Custom Agent with Query Pipeline Guide](#)



# Additional Requirements

Possible through our query pipeline syntax

## Query Pipeline Guide





# What's next for RAG: Long Contexts?



# Is RAG Dead?

Gemini 1.5 Pro has a 1-10M context window.

What does this mean for RAG?

[https://x.com/Francis\\_YAO\\_/status/1759962812229800012?s=20](https://x.com/Francis_YAO_/status/1759962812229800012?s=20)

 **Yao Fu**   
@Francis\_YAO\_ 

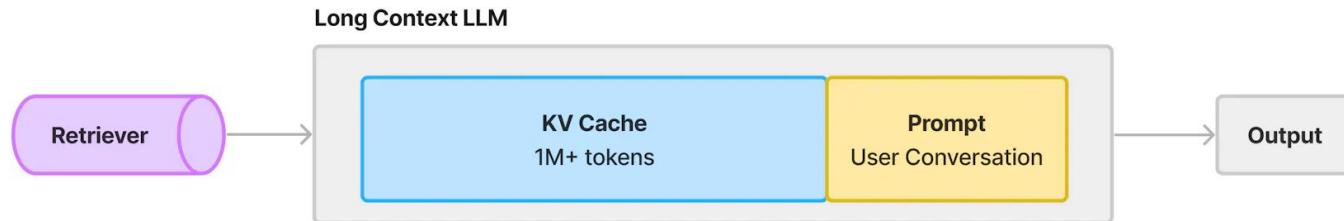
Over the last two days after my claim "long context will replace RAG", I have received quite a few criticisms (thanks and really appreciated!) and many of them stand a reasonable point. Here I have gathered the major counterargument, and try to address them one-by-one (feels like a paper rebuttal):

- **RAG is cheap, long context is expensive.** True, but remember, compared to LLM, BERT-small is also cheap, and n-gram is even cheaper, but they are not used today, because we want the model to be smart first, then makes smart models cheaper -- history of AI tells **it is much easier to make smart models cheaper than making cheap model smart** -- when it is cheap, it's never smart.
- **Long context can mix retrieval and reasoning during the whole decoding processing.** RAG only does the retrieval at the very beginning. Typically, given a question, RAG retrieves the paragraphs that are related to the question, then generates. Long-context does the retrieval for every layer and every token. In many cases the model needs to **do on-the-fly per-token interleaved retrieval and reasoning**, and only knows what to retrieve after getting the results of the first reasoning step. Only long-context can do such cases.
- **RAG supports trillion level tokens, long-context is 1M.** True, but there is a natural distribution of the input document, and I tend to believe **most of the cases that requires retrieval is under million level**. For example, imagine a layer working on a case whose input is related legal documents, or a student learning machine learning whose input are three ML books -- does not feel as long as 1B right?



# Our Position

1. Frameworks are valuable whether or not RAG lives or dies
2. Certain RAG concepts will go away, but others will remain and evolve





# Long Context LLMs will Solve the Following

1. Developers will worry less about tuning chunking algorithms
2. Developers will need to spend less time tuning retrieval and chain-of-thought over single documents
3. Summarization will be easier
4. Personalized memory will be better and easier to build



# Some Challenges Remain

1. 10M tokens is not enough for large document corpuses (hundreds of MB, GB)
2. Embedding models are lagging behind in context length
3. Cost and Latency
4. A KV Cache takes up a significant amount of GPU memory, and has sequential dependencies



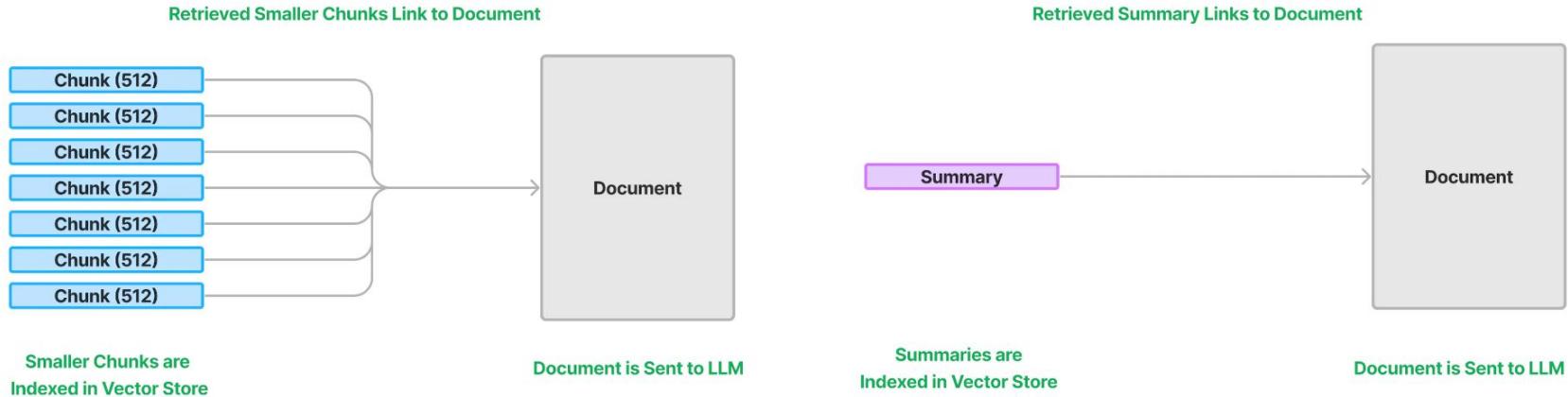
# New RAG Architectures

1. Small to Big Retrieval over Documents
2. Intelligent Routing for Latency/Cost Tradeoffs
3. Retrieval Augmented KV Caching



# Small to Big Retrieval over Documents

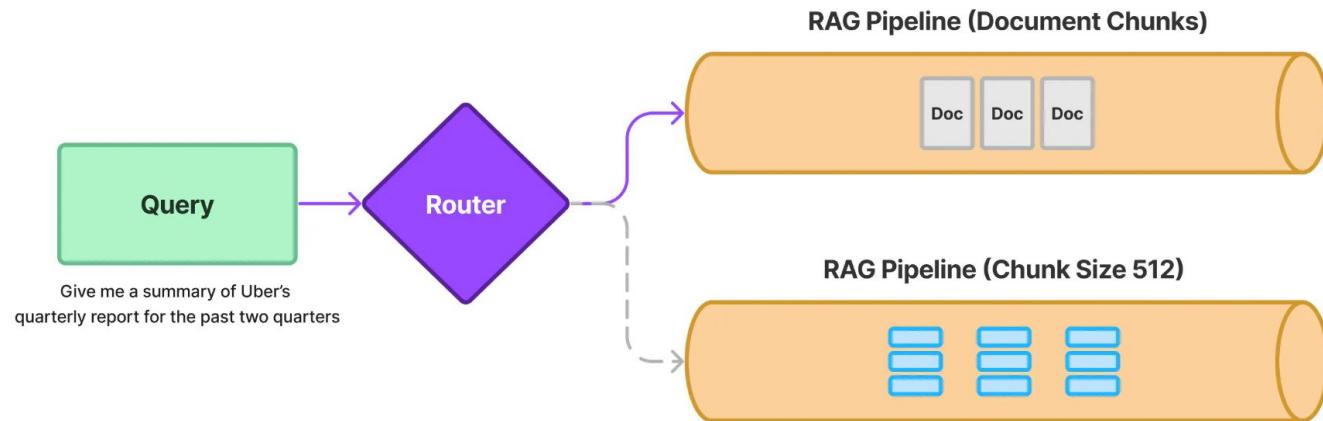
## Small-to-Big Retrieval over Documents





# Intelligent Routing for Latency/Cost Tradeoffs

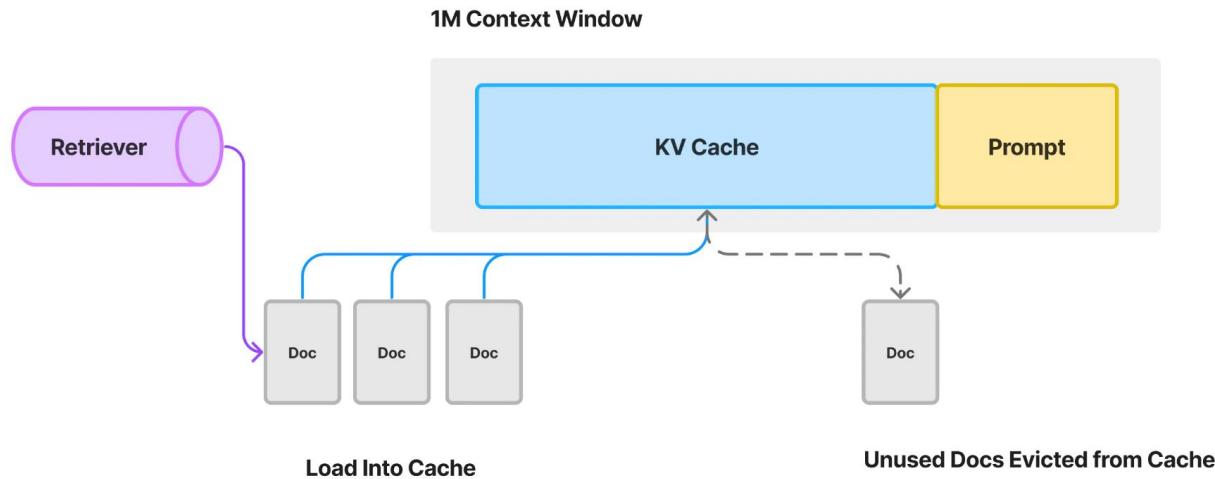
## Intelligent Routing for Latency/Cost Tradeoffs





# Retrieval Augmented KV Caching

## Retrieval for KV Caching





# serverless user group toronto

[www.ServerlessToronto.org](http://www.ServerlessToronto.org)

Reducing the gap between IT and Business needs