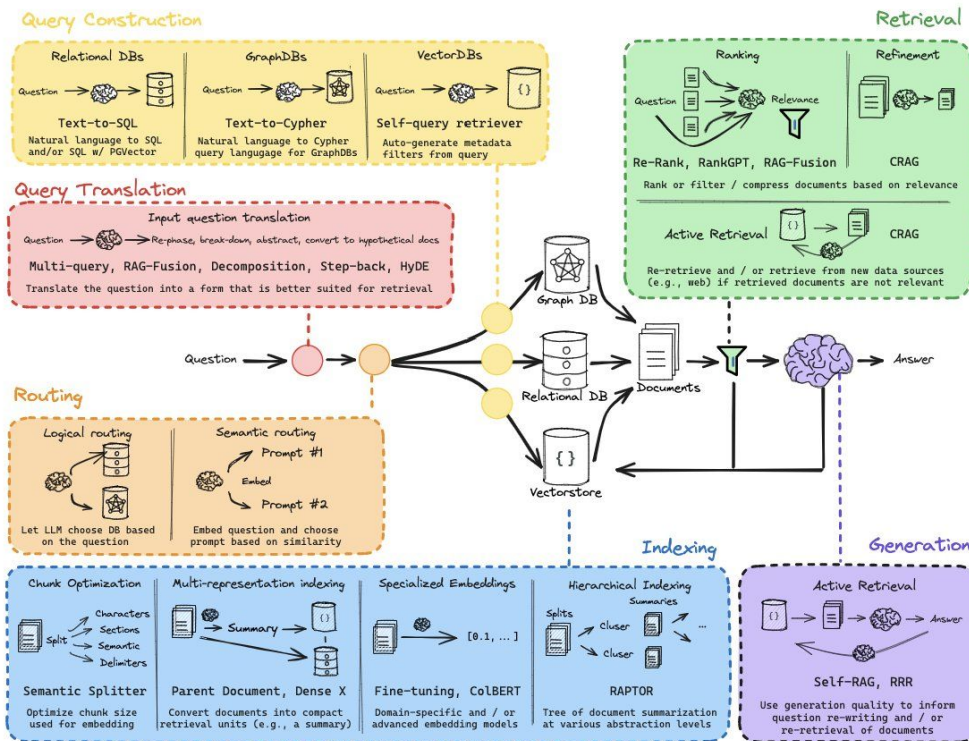


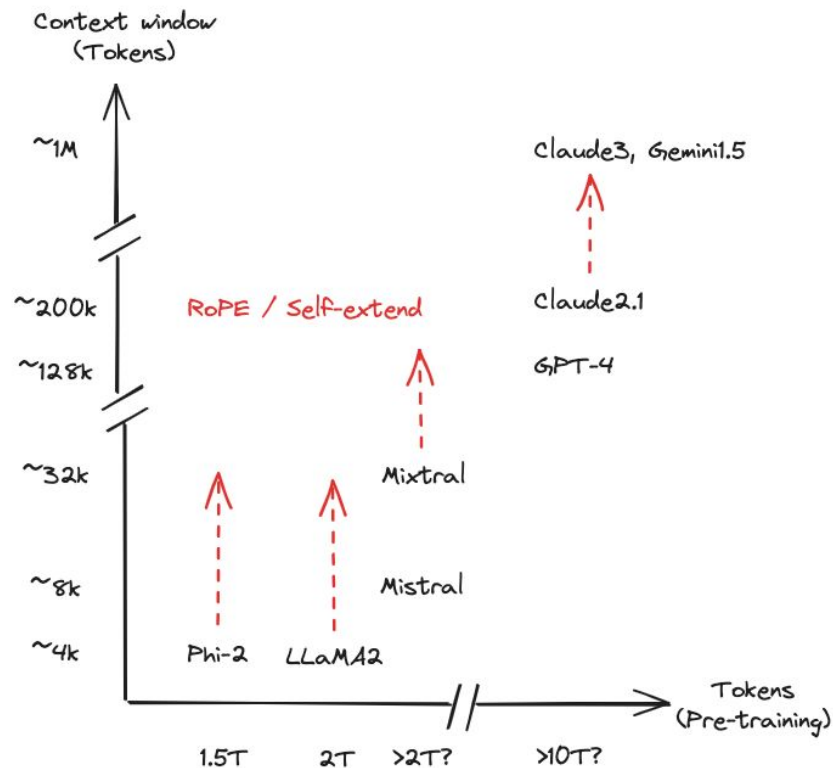
# Unifying RAG and long context LLMs

Lance Martin  
Software Engineer, LangChain  
[@RLanceMartin](https://twitter.com/RLanceMartin)

# Preface: RAG Course



## Context windows are getting larger




## Do we need RAG anymore?

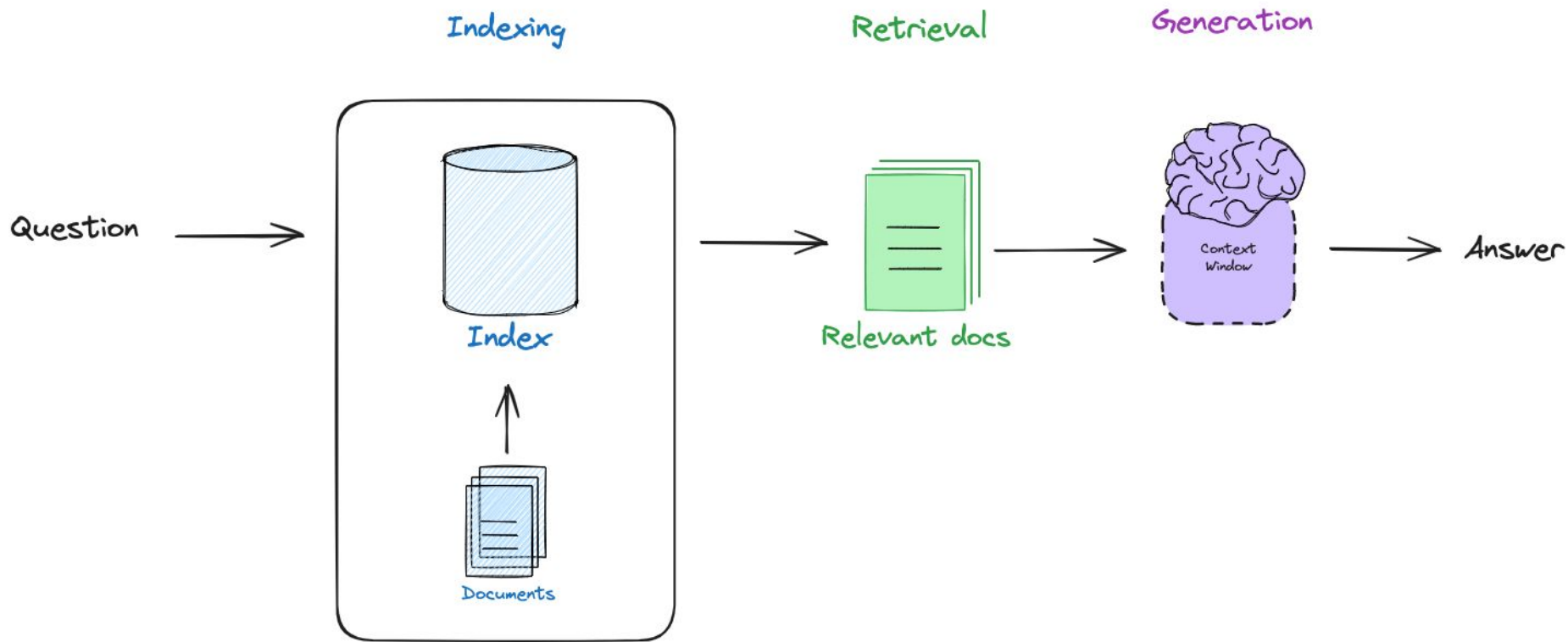


...

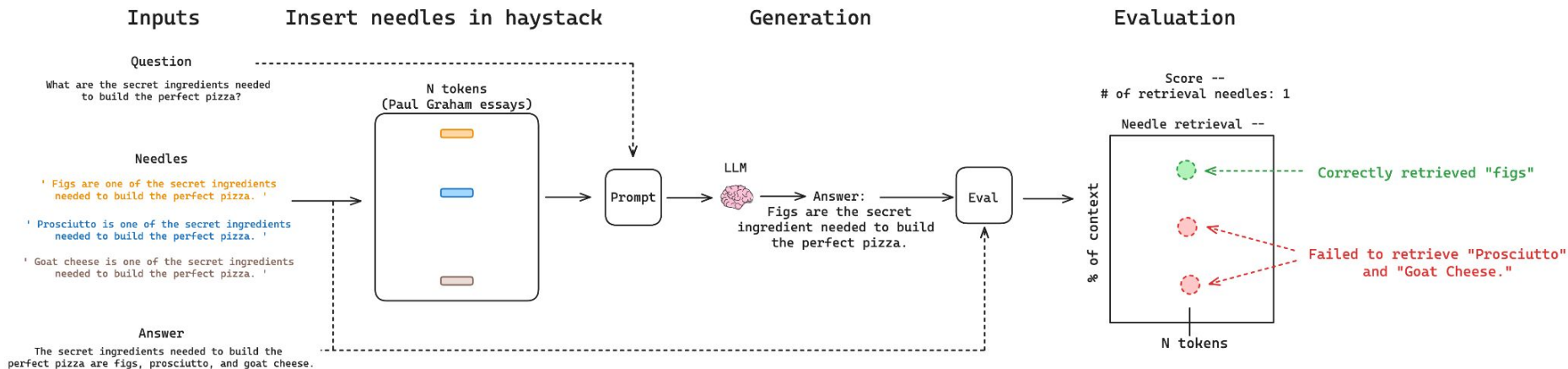
⚠️ RAG might be dead, after reading 58 pages of Gemini 1.5 Pro tech report. Here's my thoughts as AI founder,

1. Simple RAG system like similarity search with vector db will be dead. But more customized RAG will still live. The goal of RAG is mostly on retrieval relevant information. After reading the report, I am convinced LLM can do retrieval really really well.
2. RAG itself may not be dead totally, but 90% of people won't need it anymore. Most dataset can fit in 1M tokens. Just like OpenAI's assistant API, once Gemini API can handle large files, the only thing matters is the cost. However based on the report, 1.5 Pro's training cost and inference cost is much much lower than Gemini 1.0 

## RAG: reasoning & retrieval on multiple chunks of information

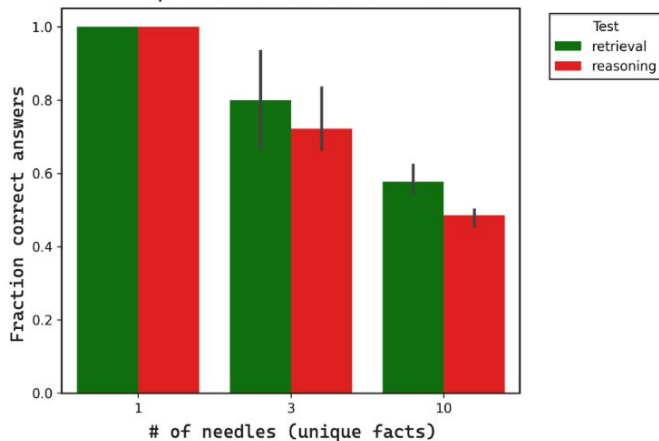


# Needle In A Haystack: test reasoning & retrieval in long context LLMs

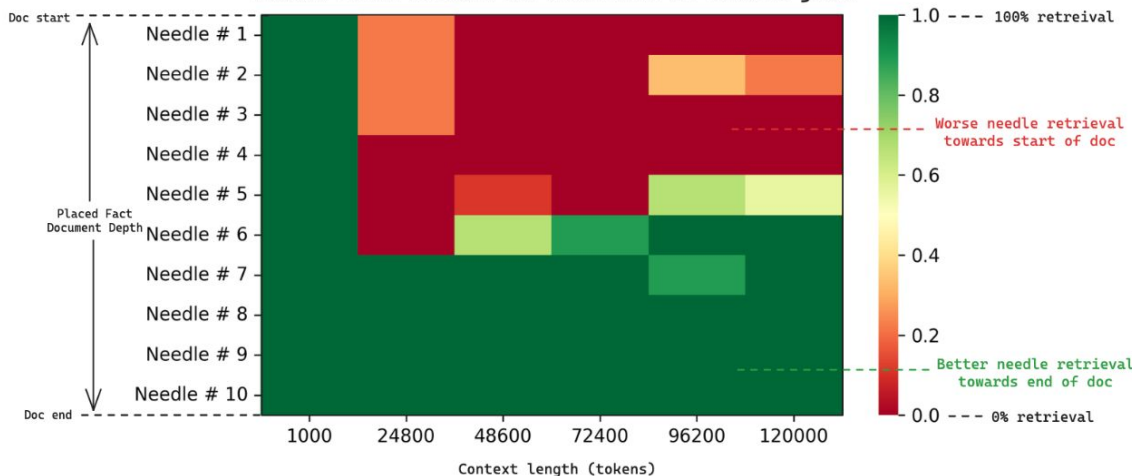


# Retrieval is not guaranteed, reasoning harder than retrieval

Asking GPT-4 to retrieve or retrieve & reason  
1, 3, or 10 needles (facts) in a single turn  
120,000 token context window



Asking GPT-4 to retrieve 10 unique facts in 1 turn  
Assess which needles are retrieved as context grows



<https://youtu.be/UlmyyYQGhzc>

<https://blog.langchain.dev/multi-needle-in-a-haystack/>

## Challenge may be recency bias in LLMs

A likely culprit for this phenomenon is a mismatch between the task LLMs are trained on and context-augmented generation tasks. Among the documents typically used to pre-train LLMs such as web pages, books, articles and code, the most informative tokens for predicting a particular token are typically the most recent ones. During pre-training, this induces a learned bias to attend to recent tokens. In addition, the rotary positional embedding (RoPE) scheme used in the open source models we investigate has an inductive bias towards reduced attention at long distances [27] that may make it even easier for these models to learn to attend preferentially to recent tokens. Extreme recency bias is not a good prior for context augmented generation tasks where far away tokens may, in fact, contain very relevant information.



(1) Be wary of context stuffing. No retrieval guarantees.

(2) Single needle may be misleadingly easy (no reasoning, only 1 fact)



Thomas Ahle    
@thomasahle

...

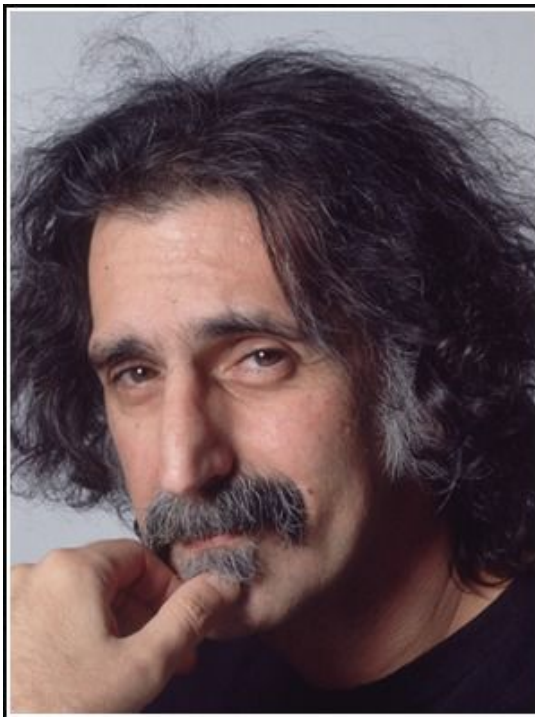
Needle In A Haystack tests are flawed.

Did you know that the Long-Context Attention in Gemini and GPT-4 is based on inserting the sentence *"The best thing to do in San Francisco is eat a sandwich and sit in Dolores Park on a sunny day"* at a random location in a text?

We've seen Claude even **anticipating** that it was being tested because of the awkwardness of the test: [twitter.com/alexalbert\\_/s....](https://twitter.com/alexalbert_/s....) The models know to keep track of the sentence before they are even prompted to retrieve it!

New tests from @NormalComputing demonstrate that models do much worse **if the inserted sentence is subtle**. Like changing the name of one historical person to another.

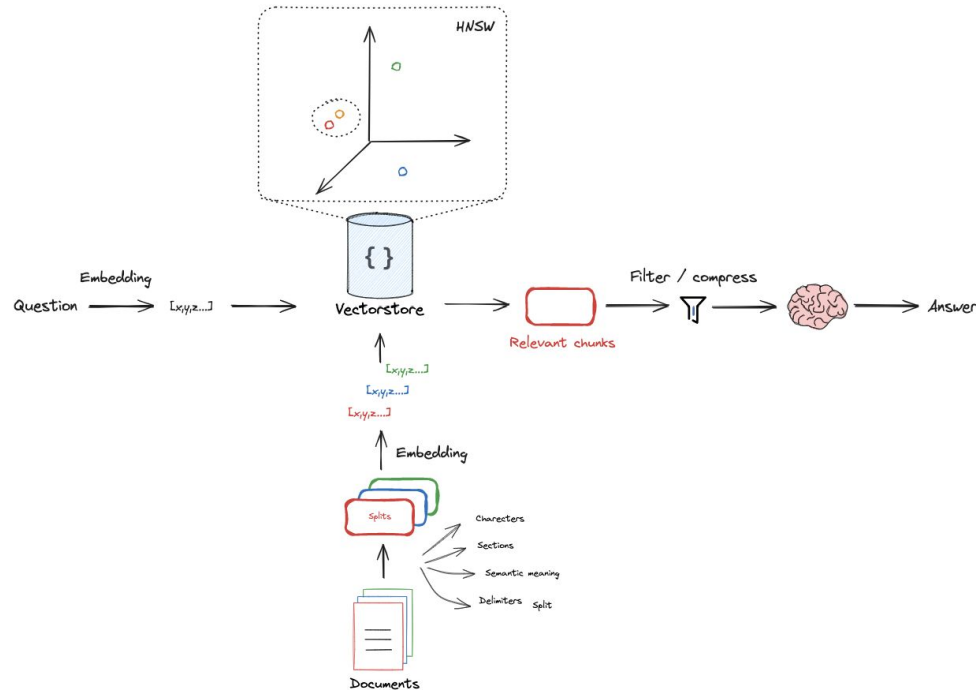
This will obviously get better and RAG will change



Jazz isn't dead. It just smells funny.

— Frank Zappa —

## RAG today focused on **precise retrieval** of relevant doc **chunks**



## Need to balance **system complexity** vs **latency & token usage**



- You need the exact relevant chunk
- risk of over-engineering
  - higher complexity
  - can suffer lower recall
  - sensitivity to chunk size, k, etc

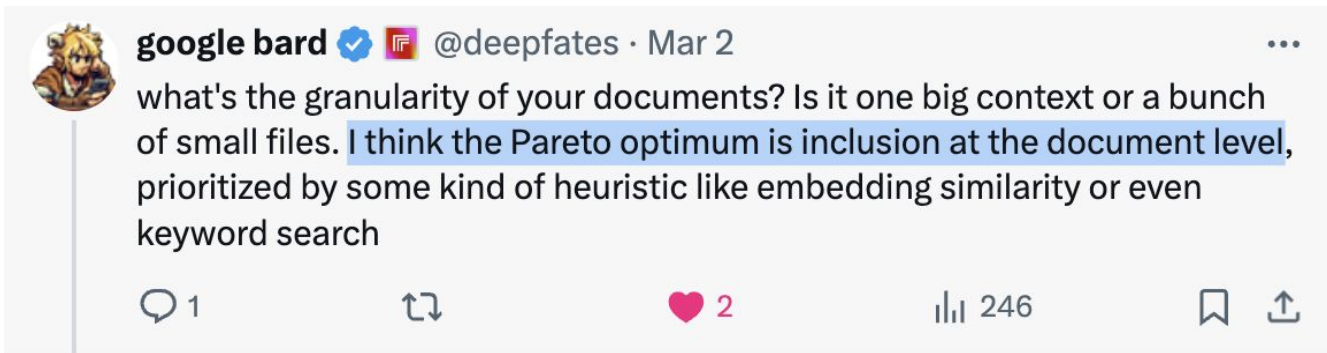


Pareto Optimal

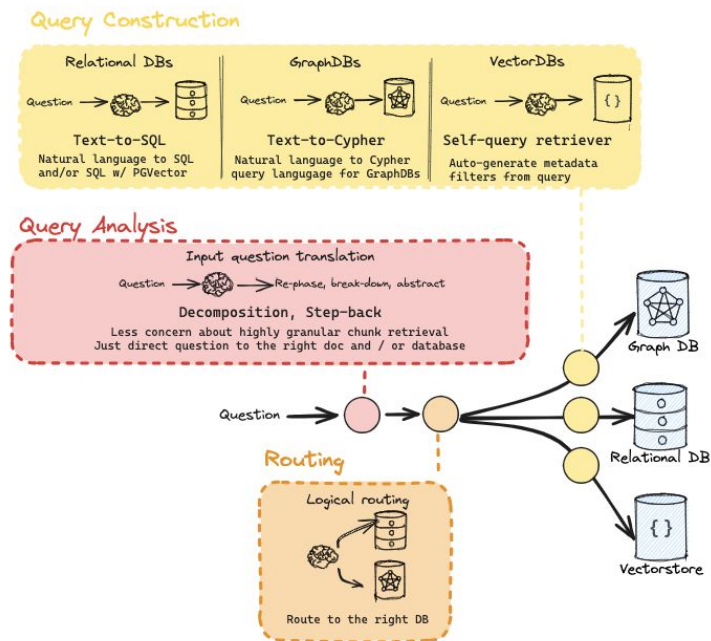


- Just throw all your docs into context
- higher latency
  - higher token usage
  - can't audit retrieval
  - security / auth

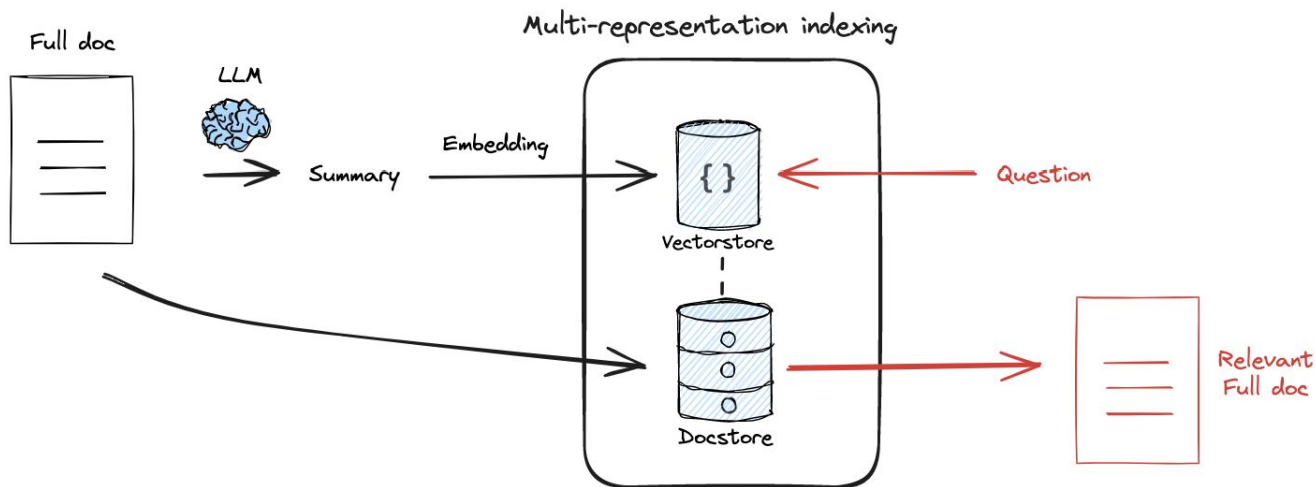
## Some ideas ...



# Query analysis: Connect questions to the right document



## Indexing: Use representations to simplify document retrieval

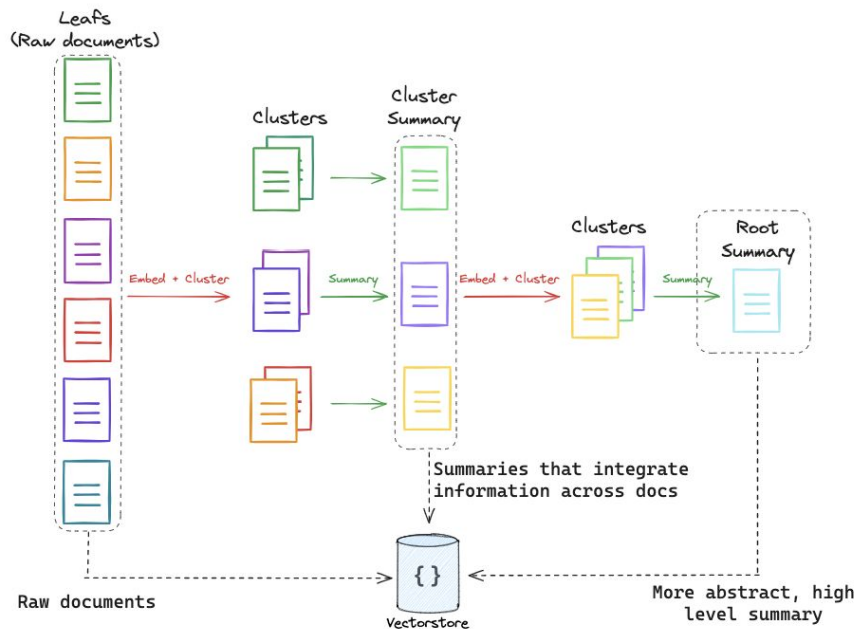


<https://arxiv.org/pdf/2312.06648.pdf>

<https://blog.langchain.dev/semi-structured-multi-modal-rag/>

[https://python.langchain.com/docs/modules/data\\_connection/retrievers/parent\\_document\\_retriever](https://python.langchain.com/docs/modules/data_connection/retrievers/parent_document_retriever)

## Indexing: Use trees to consolidate info across many documents



<https://arxiv.org/abs/2401.18059>

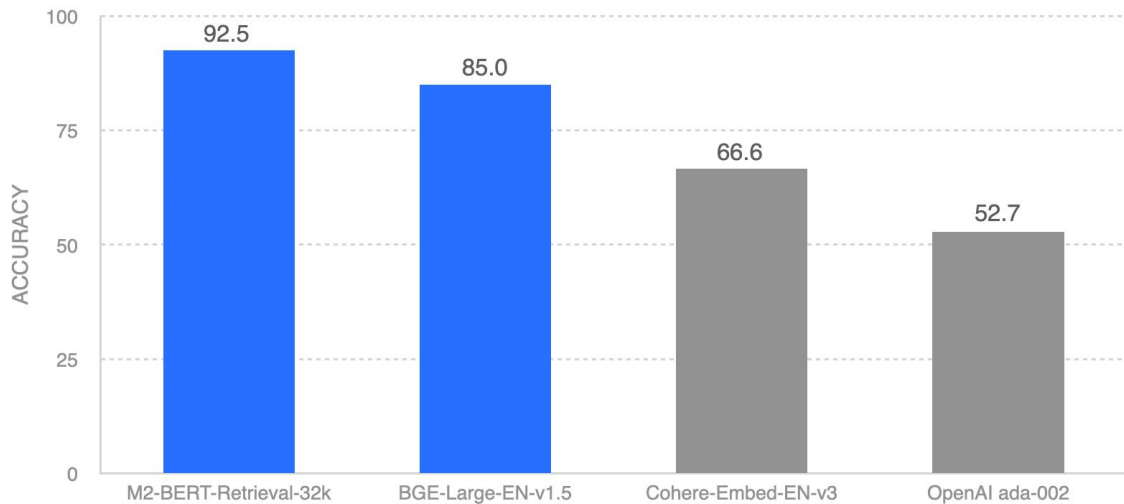
<https://www.youtube.com/watch?v=jbGchdTL7d0>



## Indexing: Use long context embeddings

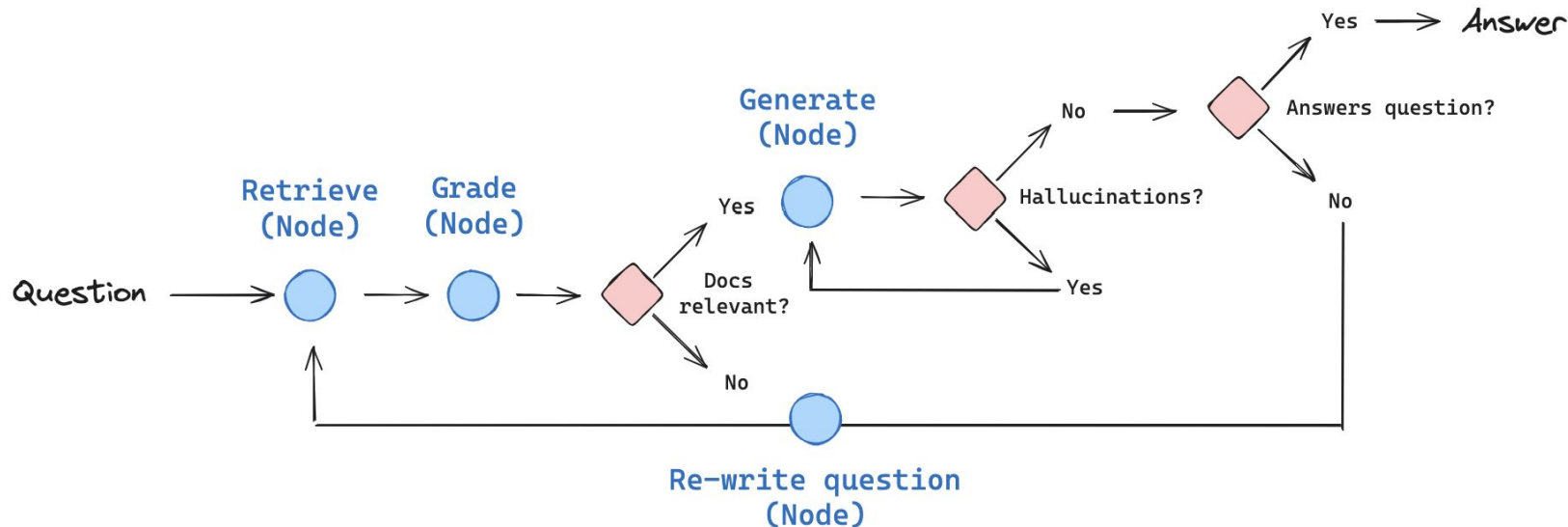
### LONG CONTEXT QUALITY

AVERAGE ACCURACY OF LOCO BENCHMARK



<https://www.together.ai/blog/embeddings-endpoint-release>  
<https://hazyresearch.stanford.edu/blog/2024-01-11-m2-bert-retrieval>

## Reasoning: Use reasoning / self-reflection around RAG

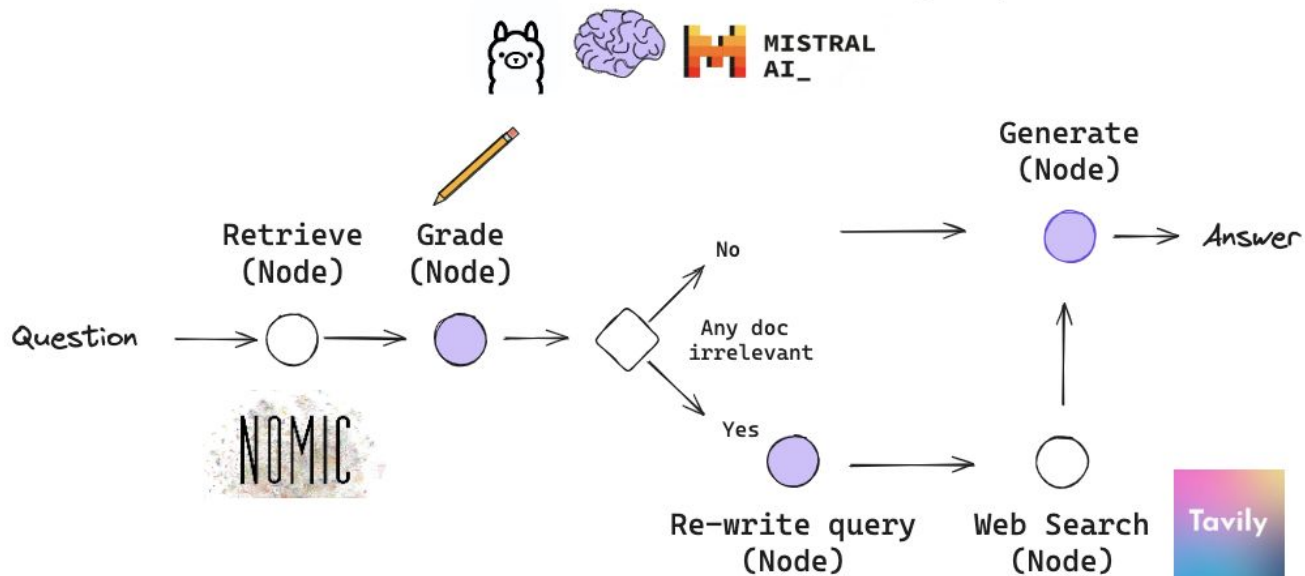


<https://arxiv.org/abs/2310.11511>

<https://www.youtube.com/watch?v=E2shqsYwxck>

[https://github.com/langchain-ai/langgraph/blob/main/examples/rag/langgraph\\_self\\_rag.ipynb](https://github.com/langchain-ai/langgraph/blob/main/examples/rag/langgraph_self_rag.ipynb)

## Reasoning: Use reasoning / self-reflection around RAG



<https://arxiv.org/abs/2401.15884>

<https://www.youtube.com/watch?v=E2shqsYwxck>

[https://github.com/langchain-ai/langgraph/blob/main/examples/rag/langgraph\\_crag.ipynb](https://github.com/langchain-ai/langgraph/blob/main/examples/rag/langgraph_crag.ipynb)

# Overall picture: Doc-centric, no splits or compression, reason pre/post retrieval

