# Unit AI-2
## Retrieval-Augmented Generation

Codesmith

# Roadmap

- Overview
- How RAG works
- RAG optimizations
- RAG evaluations

# Overview

# What is RAG?

RAG combines information retrieval with generative AI models to enhance the accuracy and relevance of AI-generated content.

It provides an effective and efficient way of programmatically scaffolding queries with additional context.

# Why focus on RAG?

LLMs are not designed to retrieve data - RAG is key to unlocking their potential.

The RAG layer is especially relevant for AI engineering.
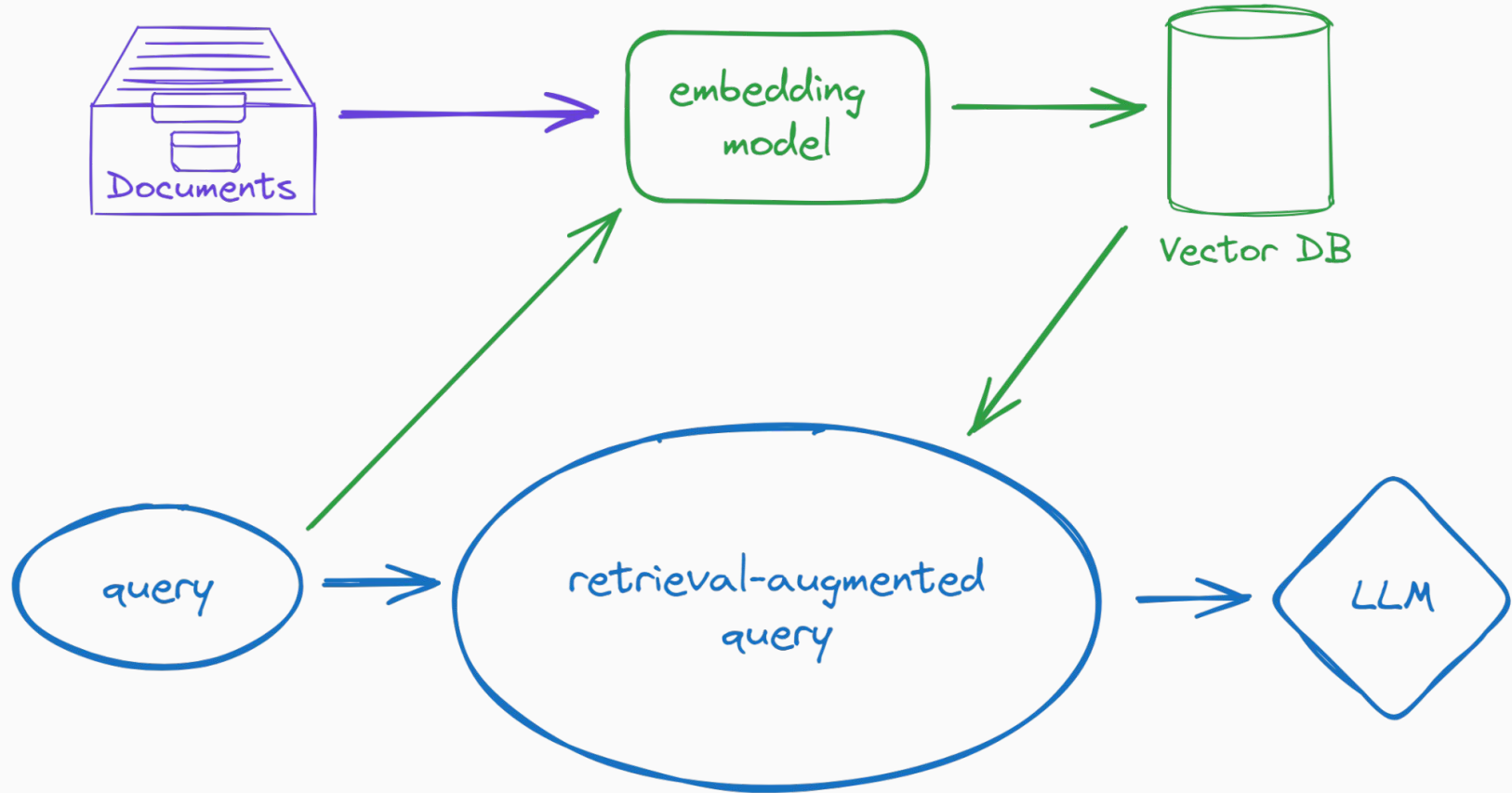
# Benefits of RAG

- Groundedness
  - Relevance
  - Accuracy
  - Currentness
  - Domain / proprietary specificity
- Interpretability
- Efficiency

# Prompting vs RAG vs fine-tuning

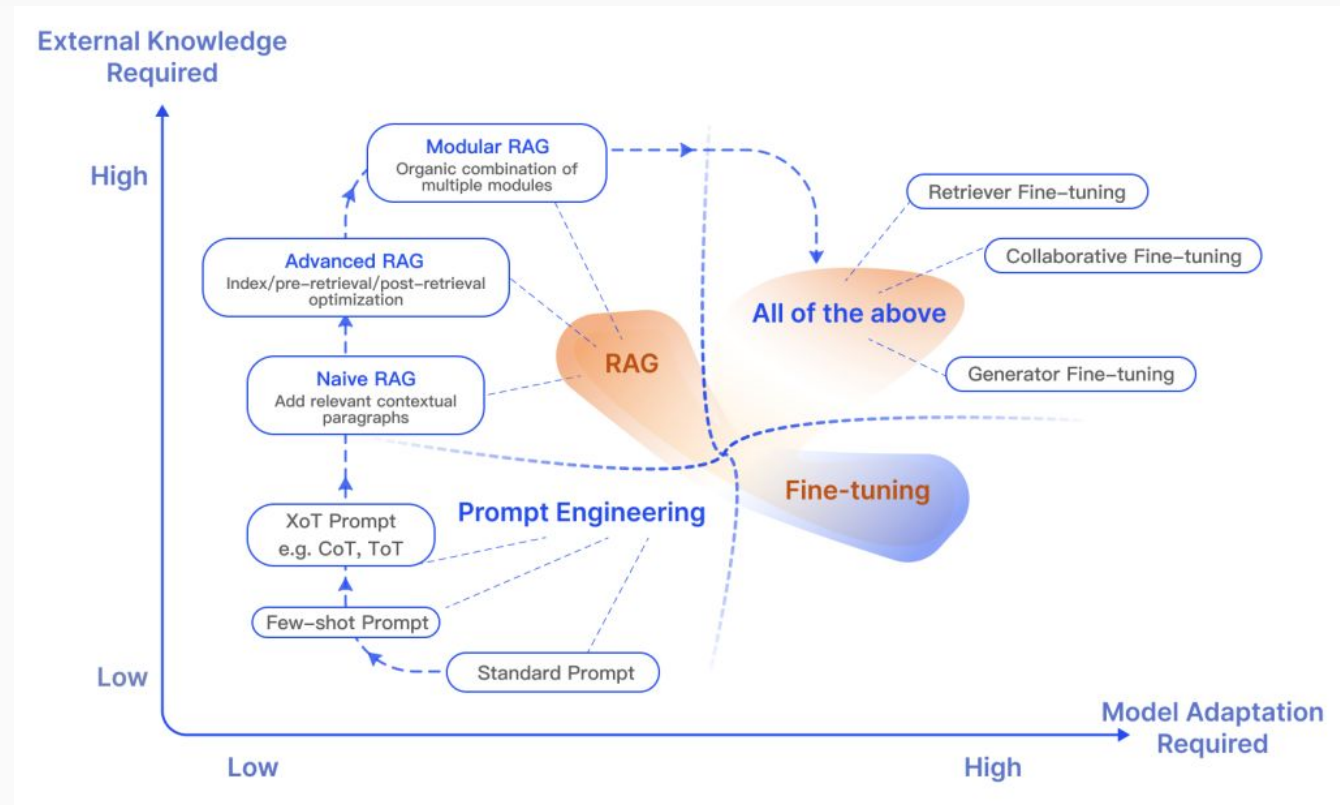| Goal | Prompting | RAG | Fine-Tuning |
|---|---|---|---|
| Grounding | 😐 | ✅ | 😐 |
| Consistency | ✅ | 😐 | ✅ |
| Confidence | ✅ | 😐 | ✅ |
| Interpretability | ❌ | ✅ | ❌ |
| Alignment | 😐 | 😐 | ✅ |
| Robustness | 😐 | ❌ | ✅ |
| Latency | ❌ | 😐 | ✅ |
| Cost | ❌ | 😐 | 😐 |

# How RAG works

# RAG workflow

# Naive RAG challenges

- Precision
- Recall
- Faithfulness
- Answer relevance
- Accuracy

# Data ingestion

- Discovery
- Acquisition
- Validation
- **Transformation**
- **Loading**

# RAG vs fine-tuning



Retrieval-Augmented Generation for Large Language Models: A Survey

# Retrieval optimization

# Chunk optimization

Balance retrieval specificity against generation context.
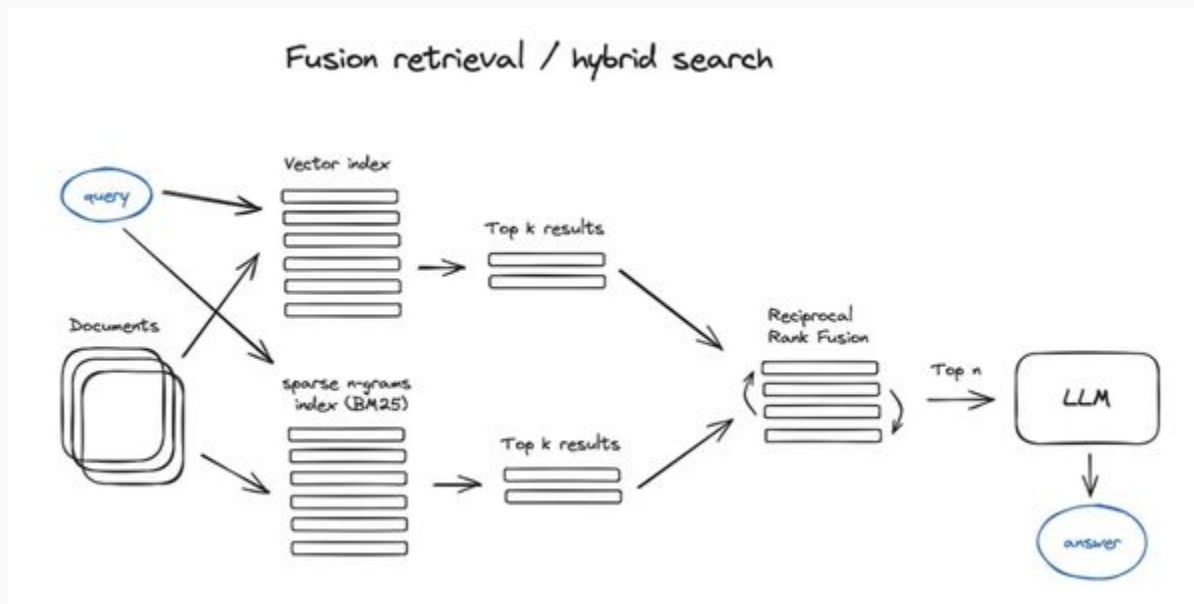
- Fixed-size

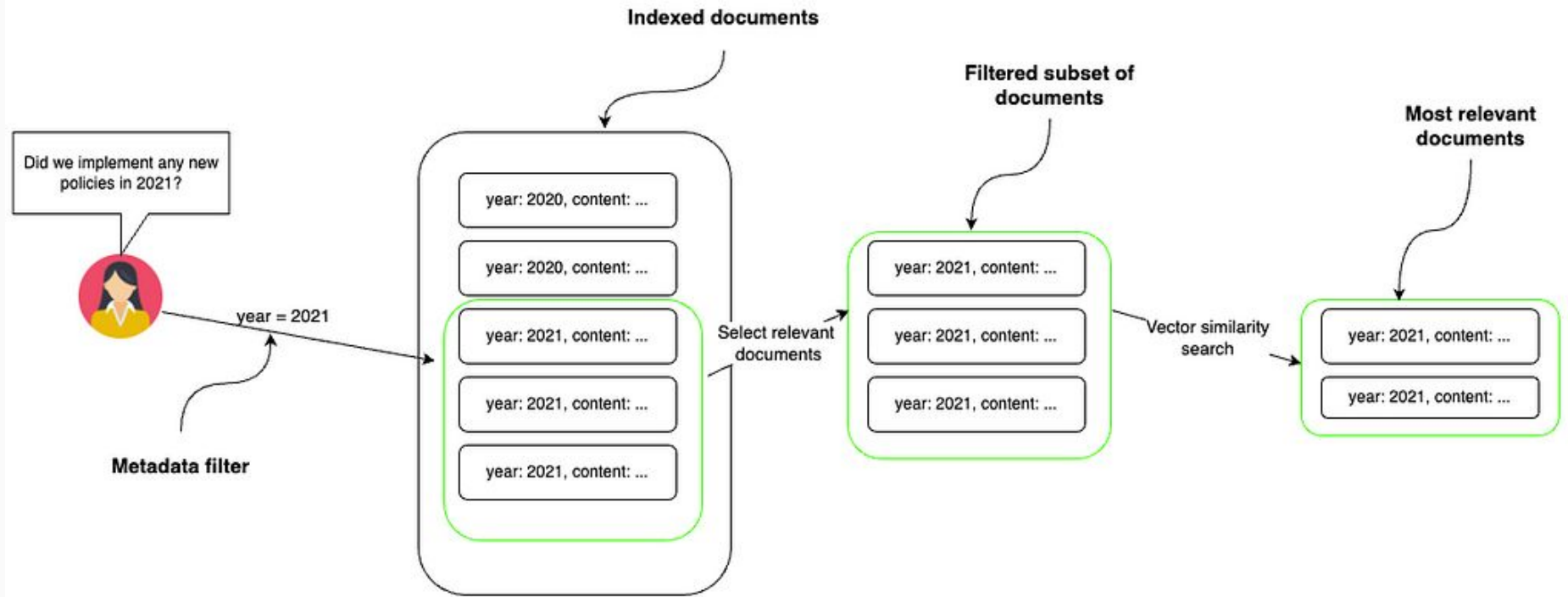- Intent-based & Recursive

- Strategy-based

# Hybrid retrieval

Capitalize on traditional information retrieval techniques.

- Ensemble (rank aggregation / weighted averaging)

- Cascade (filtering layers)

# Keyword + embedding retrieval



Fusion retrieval / hybrid search

*Wei, Huang, and Wang: Retrieval-Augmented Generation for LLM Applications*

# Metadata filtering
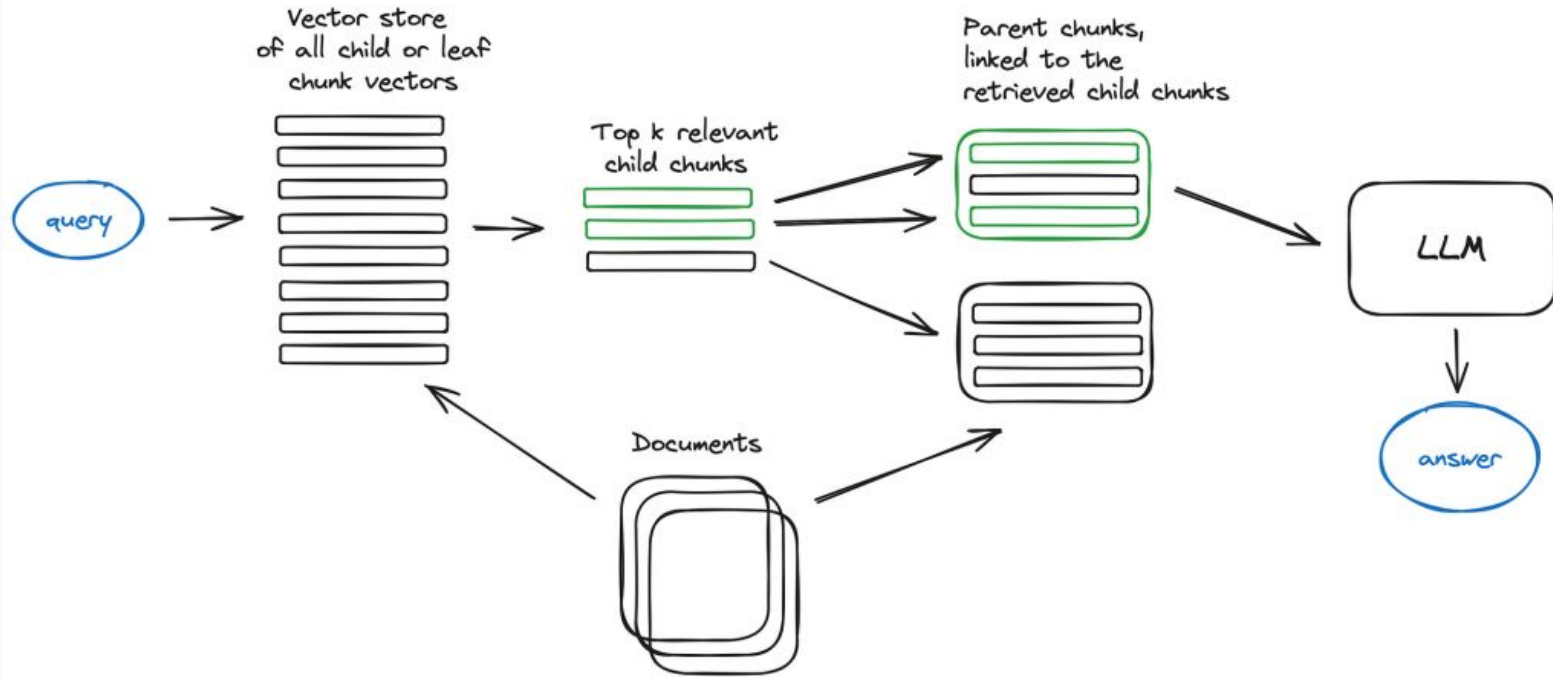
# Small-to-big retrieval

Optimize for retrieval *and* generation by expanding context after initial search.

- Sentence window retrieval
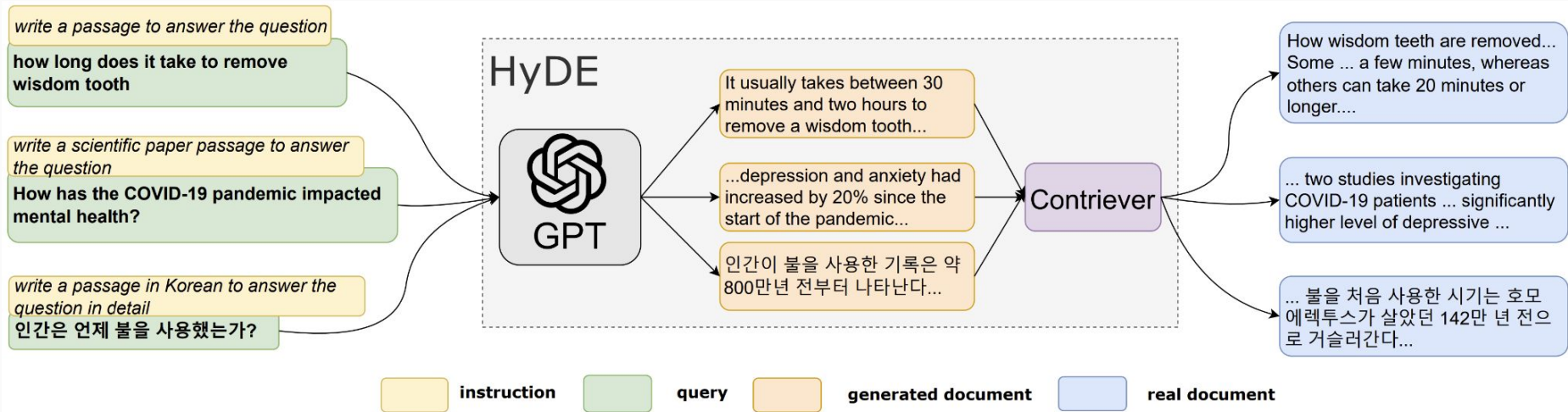
- Parent document crawler

Similarly, retrieval could proceed recursively to build context.

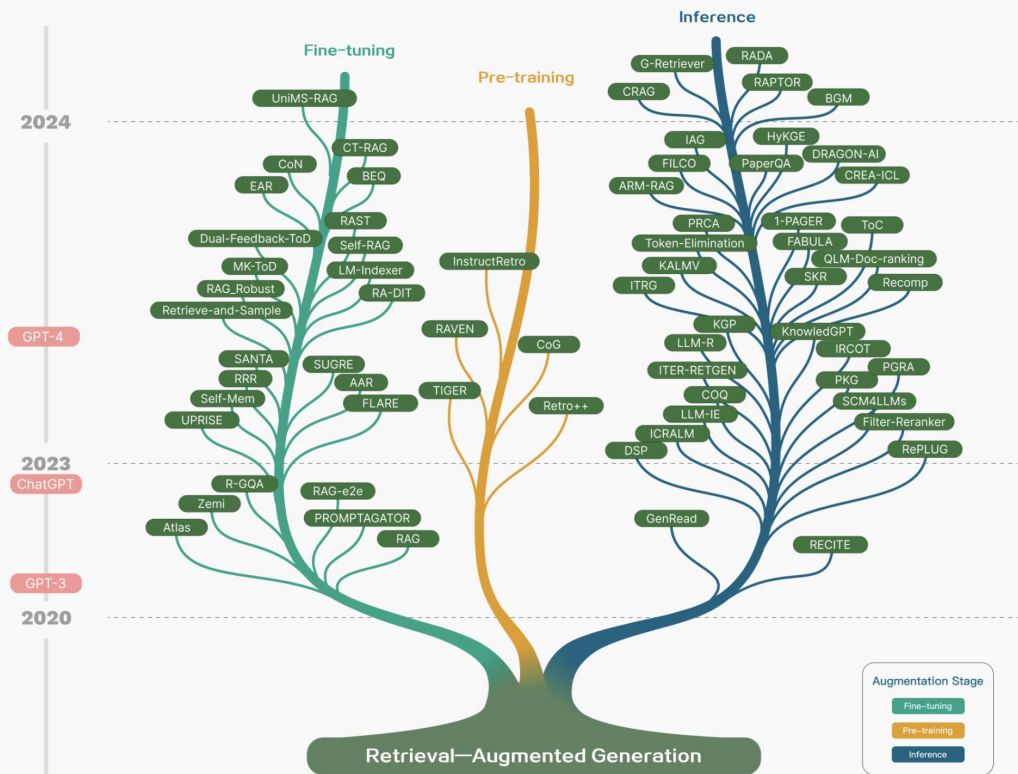# Parent document crawler



Parent-child chunks retrieval
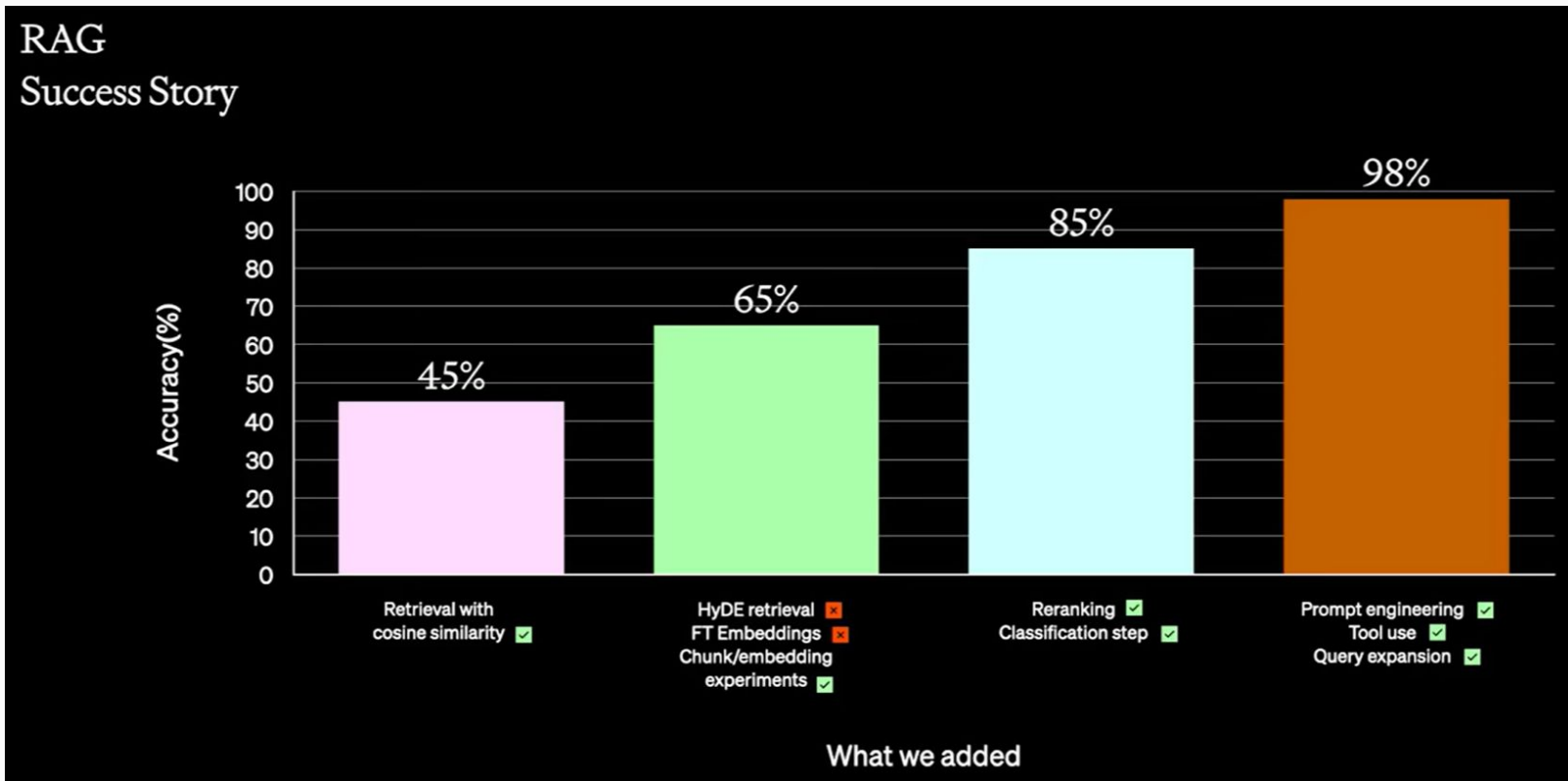
# Hypothetical Document Embeddings (HyDE)



*Precise Zero-Shot Dense Retrieval without Relevance Labels*

# And many more...



_Retrieval-Augmented Generation for Large Language Models: A Survey_

# OpenAI "RAG Success Story"



*A Survey of Techniques for Maximizing LLM Performance*

# RAG evaluation

# Information retrieval metrics

- Recall

- Precision

- Mean reciprocal rank (MRR)

- Normalized discounted cumulative gain (NDCG)

# Building an evaluation dataset

- Benchmarks

- Web data

- Synthetic data

- Proprietary data

# RAG testing pyramid

As with traditional applications, it's beneficial to test each unit (retrieval, augmentation, and generation), their integration, and the end-to-end result.

# RAG evaluation dimensions

- Retrieval relevance

- Groundedness

- Adaptability

- Toxicity

- Efficiency

# Retrieval relevance

- ROUGE (Recall-Oriented Understudy for Gisting Evaluation)

- RAGAs (RAG Assessment)

# Groundedness

- Entity linking accuracy

- Semantic similarity scores

- FactScore

# Generation-specific dimensions

- Noise robustness

- Negative rejection

- Information integration

- Counterfactual robustness

# Further considerations

# RAG vs long context

SOTA models boast long context windows and have dramatically improved their performance - but RAG may still be preferable given:

- Expanded scope

- Increased precision

- Reduced cost

# RAGOps

- Vector DBs

- Frameworks + libraries

- RAG 2.0