# From Vectors to Agents: Managing RAG in an Agentic World

Rajiv Shah
Chief Evangelist, Contextual AI
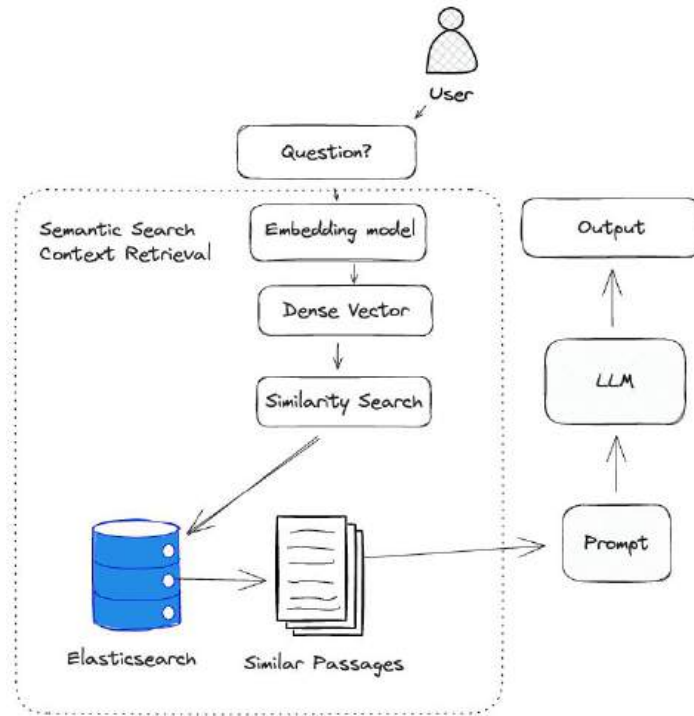rajiv.shah@contextual.ai

https://github.com/rajshah4/LLM-Evaluation

@rajistics

# Building RAG is Easy

# Building RAG is Easy

```python
docs = TextLoader("docs/sample.txt").load()
chunks = RecursiveCharacterTextSplitter(chunk_size=800).split_documents(docs)
vdb = FAISS.from_documents(chunks, OpenAIEmbeddings())
retriever = vdb.as_retriever(search_kwargs={"k": 4})
llm = ChatOpenAI(model="gpt-4o-mini", temperature=0)

prompt = ChatPromptTemplate.from_messages([
    ("system", "Answer only using the context below."),
    ("human", "Q: {question}\n\nContext:\n{context}\n\nA:")
])

rag_chain = (
    {"context": retriever | (lambda d: "\n\n".join(x.page_content for x in
d)), "question": RunnablePassthrough()}
    | prompt | llm | StrOutputParser()
)

print(rag_chain.invoke("What warranty terms are mentioned?"))
```

# RAG Reality Check

**95%**

of Gen AI projects fail to reach Production

### Accuracy
**<70%**

fails beyond simple extraction

### Latency
**>45s**

queries are too slow

### Scaling
**>1,000**

fails with more documents

### Cost
**100x**

complex queries more token use
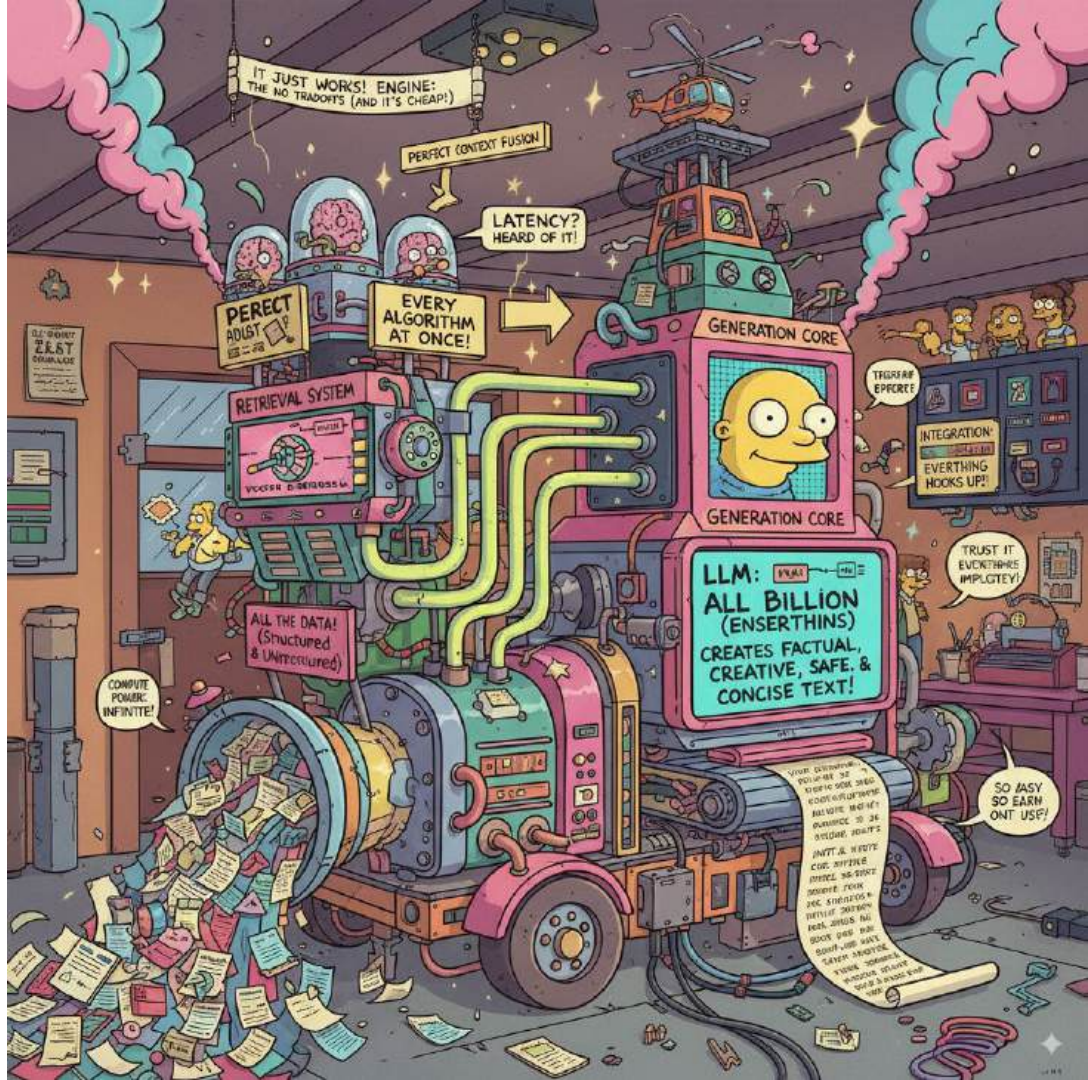
### Compliance
**0%**

access control over documents

https://www.zeta-alpha.com/post/why-genai-pilots-fail-common-challenges-with-enterprise-rag
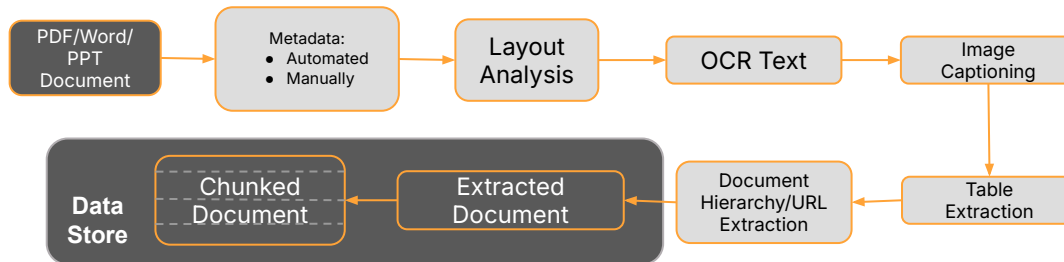
# Maybe try a different RAG?

Basic RAG, Reliable RAG, HyDE (Hypothetical Document Embedding), HyPE (Hypothetical Prompt Embedding), Contextual Chunk Headers, Semantic Chunking, Contextual Compression, Document Augmentation, Fusion Retrieval, Reranking, Multi-faceted Filtering, Hierarchical Indices, Ensemble Retrieval, Dartboard Retrieval, Multi-modal RAG with Captioning, Retrieval with Feedback Loop, Adaptive Retrieval, Iterative Retrieval, DeepEval, GroUSE, Explainable Retrieval, Graph RAG with LangChain, Microsoft GraphRAG, RAPTOR, Self-RAG, Corrective RAG (CRAG), Sophisticated Controllable Agent, Vision-RAG, Cache-Augmented Generation (CAG), Agentic RAG, Retrieval-Augmented Fine-Tuning (RAFT), Self-Reflective RAG, RAG Fusion, Temporal Augmented Retrieval (TAR), Plan-then-RAG (PlanRAG), GraphRAG, FLARE, Contextual Retrieval, GNN-RAG
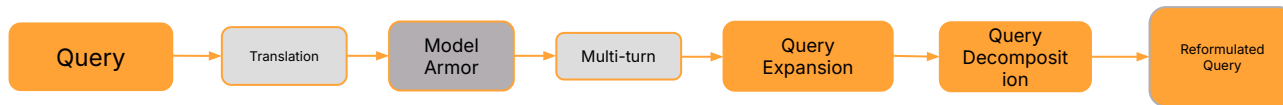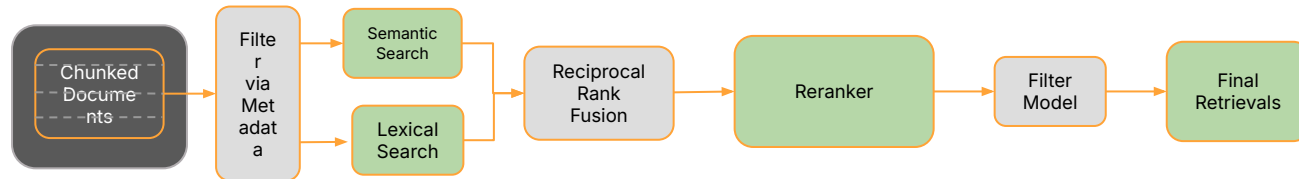
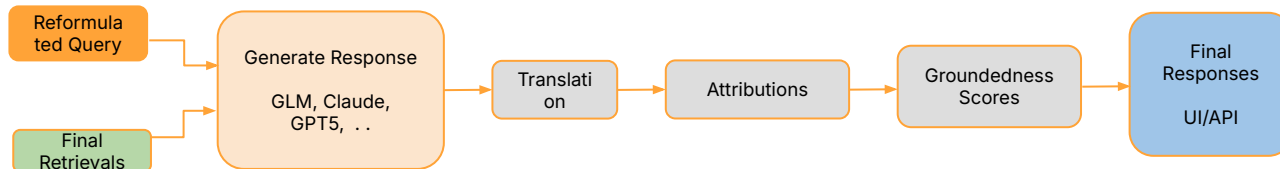# Ultimate RAG Solution



@rajistics

RAG as a system

**1. Parsing**

PDF/Word/PPT Document → Metadata: • Automated • Manually → Layout Analysis → OCR Text → Image Captioning

Data Store — Chunked Document ← Extracted Document ← Document Hierarchy/URL Extraction ← Table Extraction

**2. Querying**

Query → Translation → Model Armor → Multi-turn → Query Expansion → Query Decomposition → Reformulated Query

**3. Retrieving**

Chunked Documents → Filter via Metadata → Semantic Search / Lexical Search → Reciprocal Rank Fusion → Reranker → Filter Model → Final Retrievals

**4. Generation**

Reformulated Query / Final Retrievals → Generate Response GLM, Claude, GPT5, . . → Translation → Attributions → Groundedness Scores → Final Responses UI/API

@rajistics

# Designing a RAG Solution

Problem Complexity
(instead of accuracy)

RAG
Tradeoffs

Latency

Cost

Practical:
Cost of a mistake

# RAG Considerations

- Extraction
- Latency
- Amount of Queries
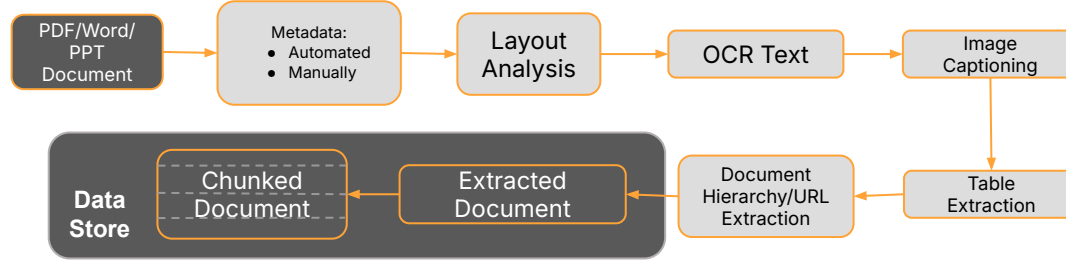- Multilingual
- Domain difficulty
- Data Quality

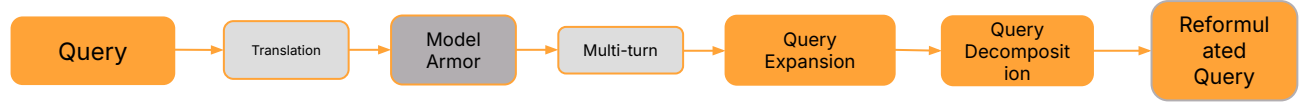| Generation → | 1. Simple Fact | 2. Summarization | 3. Multi-Source Synthesis | 4. Deep Reasoning/Analysis |
|---|---|---|---|---|
| Retrieval 1: Single-hop | Basic factual Q&A | Short doc summary | Summarize from 2–3 texts | Single-hop but deep reasoning |
| Retrieval 2: Multi-hop | Factual, but requires combining 2 steps to retrieve | Summaries that rely on multi-step retrieval | Synthesize multi-doc, multi-hop context | Multi-hop with multi-step logic in generation |
| Retrieval 3: Cross-domain | Straight pass-through, but from different data sources | Summaries that span multiple domains (e.g., news + scientific articles) | Cross-domain synthesis (e.g., financial + technical) | Complex reasoning across domain boundaries |
| Retrieval 4: Ambiguous / advanced | Passing through uncertain context or ambiguous queries | Summaries that handle contradictory / ambiguous sources | Complex bridging across ambiguous queries + multi-sources | Highest difficulty: multi-hop + cross-domain + advanced reasoning |

# Consider Query Complexity

**Simple Keyword**
1. What is Tesla's total revenue in 2023?

**Semantic variation**
2. How much bank did Tesla make last year from its operations?

**Multi-hop**
3. Compare Tesla's revenue growth in 2023 with Rivian's net loss in the same year.

**Cross-document**
4. Summarize how EV companies described supply chain issues in their 2023 filings.

**Out of corpus**
5. In Rivian's 10-K, they mention compliance with the Clean Air Act. What specific obligations does this impose on them?

**Agentic scenario**
6. If I were evaluating Rivian's environmental liabilities, how do the obligations under the Clean Air Act and California's Zero Emission Vehicle mandate intersect with the risks they disclosed in their last two annual reports?"
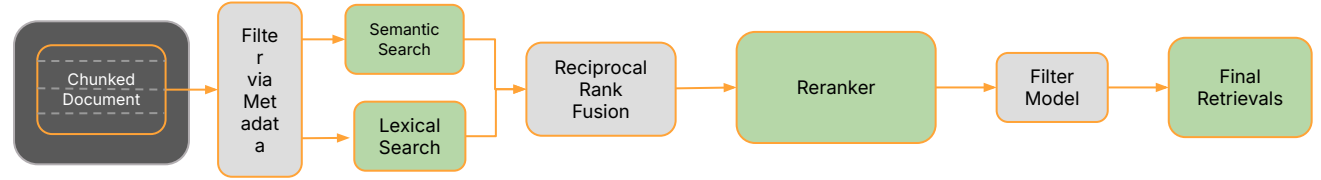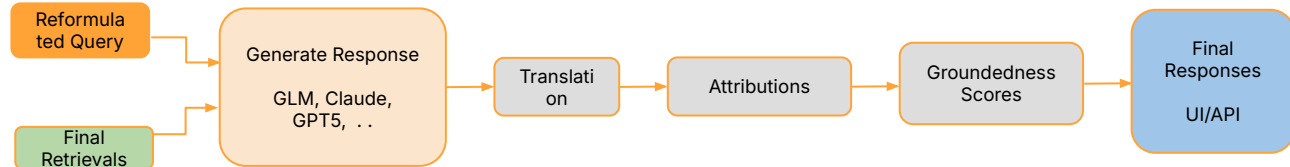
@rajistics

**1. Parsing**

PDF/Word/PPT Document → Metadata: • Automated • Manually → Layout Analysis → OCR Text → Image Captioning

Data Store — Chunked Document ← Extracted Document ← Document Hierarchy/URL Extraction ← Table Extraction

**2. Querying**

Query → Translation → Model Armor → Multi-turn → Query Expansion → Query Decomposition → Reformulated Query

**Retrieval**

**3. Retrieving**

Chunked Document → Filter via Metadata → Semantic Search / Lexical Search → Reciprocal Rank Fusion → Reranker → Filter Model → Final Retrievals

**4. Generation**

Reformulated Query / Final Retrievals → Generate Response GLM, Claude, GPT5, . . → Translation → Attributions → Groundedness Scores → Final Responses UI/API

@rajistics

# Retrieval Approaches

**BM25**
Keyword-based retrieval

**Language Models**
Semantic meaning with embeddings

**Agentic Search**
Dynamic
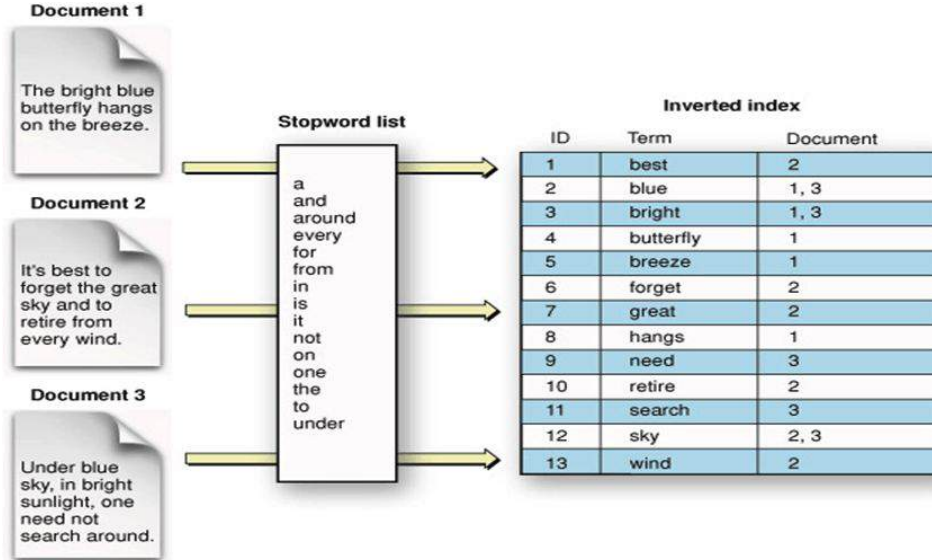using LLM Reasoning

@rajistics

# Building RAG is Easy

```python
docs = TextLoader("docs/sample.txt").load()
chunks = RecursiveCharacterTextSplitter(chunk_size=800).split_documents(docs)
vdb = FAISS.from_documents(chunks, OpenAIEmbeddings())
retriever = vdb.as_retriever(search_kwargs={"k": 4})
llm = ChatOpenAI(model="gpt-4o-mini", temperature=0)

prompt = ChatPromptTemplate.from_messages([
    ("system", "Answer only using the context below."),
    ("human", "Q: {question}\n\nContext:\n{context}\n\nA:")
])

rag_chain = (
    {"context": retriever | (lambda d: "\n\n".join(x.page_content for x in
d)), "question": RunnablePassthrough()}
    | prompt | llm | StrOutputParser()
)

print(rag_chain.invoke("What warranty terms are mentioned?"))
```

# BM25



Probabilistic lexical ranking function

# BM25 Performance

- Keyword precision

- Efficient at scale

- Battle-tested

| N_docs | Linear (s) | Inverted Index | BM25 (s) |
|--------|-----------|----------------|----------|
| 1000 | 3.468 | 0.005 | 0.028 |
| 3000 | 10.188 | 0.014 | 0.097 |
| 6000 | 20.608 | 0.025 | 0.24 |
| 9000 | 30.092 | 0.061 | 0.36 |

## BM25 Failure Cases

Lexical, probabilistic matching can mis-rank when meaning diverges from exact word overlap.

### Synonym Gap (Vocabulary Mismatch)

Query: "physician salary cap policy"

Doc A (relevant): "Doctor compensation limits for hospitals..."  ●

Doc B (distractor): "Company salary cap policy for managers..."  ✕

**BM25** Overlaps on "salary", "cap", "policy" → often ranks Doc B above Doc A.

Why: No shared token for physician↔doctor.

### Aliases & Abbreviations

Query: "International Business Machines layoffs 2024"

Doc A (relevant): "IBM announced workforce reductions in 2024..."  ●

Doc B (distractor): "International business trends show slower layoffs..."  ✕

**BM25** Matches literal words → Doc B can outrank Doc A.
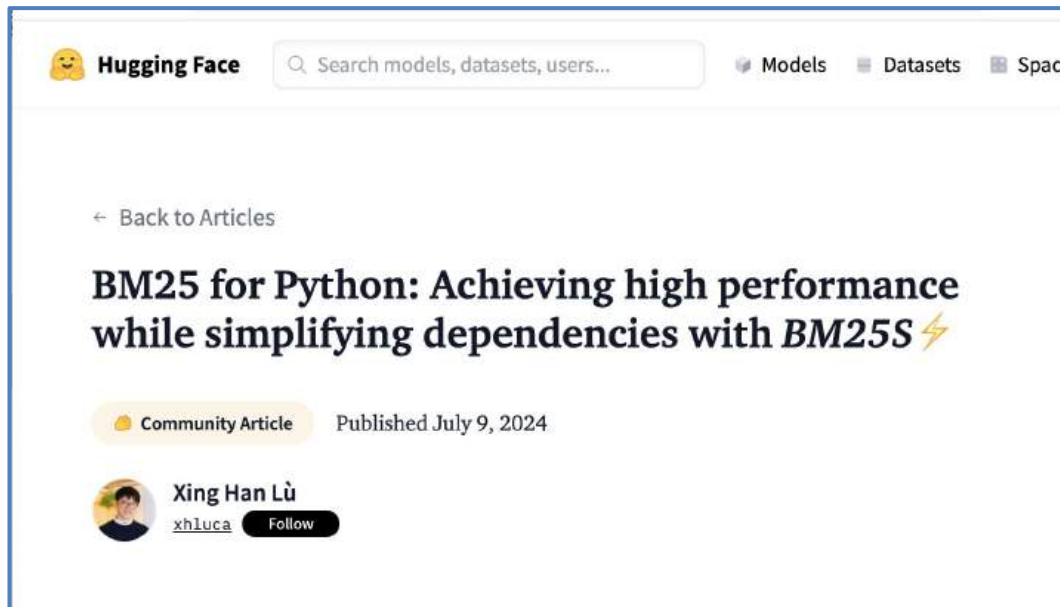
Why: Misses IBM ≡ International Business Machines.

Takeaway:

BM25 is a strong baseline

If you have keyword-heavy queries and need sub-second response →
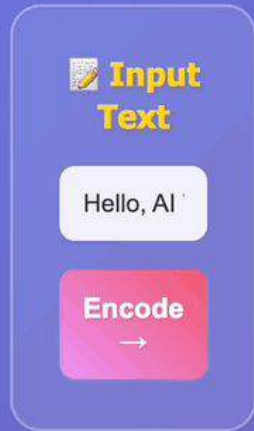
BM25 might be sufficient

@rajistics

# Hands on: BM25s

- Fast lexical search implementing BM25 in Python using Numpy, Numba and Scipy



https://github.com/xhluca/bm25s

# Enter Language Models



@rajistics

# Embeddings Visualized

# Semantic search is widely used





@rajistics

# Which language model?



Inference Speed vs NDCG@10 Scores on the NanoBEIR Benchmark

https://huggingface.co/blog/static-embeddings

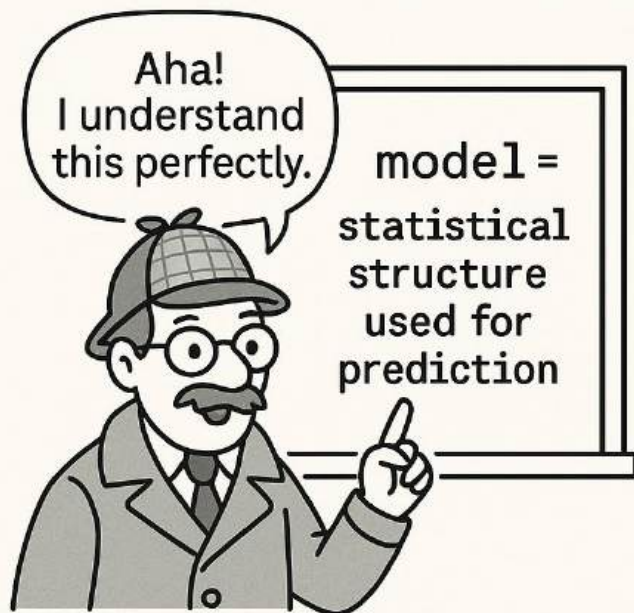# **Static Embeddings**



Model2Vec



- Uncontextualized
  - less accuracy
- Fast
- Lightweight CPU

Older versions: FastText, Word2Vec, Glove

https://github.com/MinishLab/model2vec
https://projector.tensorflow.org/

@rajistics

# Many more models!



Inference Speed vs NDCG@10 Scores on the NanoBEIR Benchmark

# MTEB/RTEB

- 300 Models
- 100+ Tasks
- 1000+ Languages

https://huggingface.co/spaces/mteb/leaderboard
https://huggingface.co/blog/rteb
Hands On BEIR: https://colab.research.google.com/drive/1HfutiEhHMJLXiWGT8pcipxT5L2TpYEdt?authuser=1

# Selecting a embedding model



- Accuracy
- Latency
- Compute (CPU/GPU)

# Selecting a embedding model

Other considerations:
- Model Size
- Architecture (CPU/GPU/Quantization)
- Embedding Dimension (128 to 8960)
- Training Data (Multilingual, Domain)
  - Fine Tuning

# Matryoshka Embedding Models



It's so sunny outside!

Matryoshka Embedding model

dim=1024 ⋯

Truncating

dim=512
dim=256
dim=128
dim=64

https://huggingface.co/blog/matryoshka

@rajistics

# Sentence Transformer

- Designed for Sentence-Level Meaning
- Semantic Search Ready
- Better Performance on Retrieval
- Efficiency

Elmo -> BERT -> DistilBERT
Sentence Transformers

https://sbert.net/

@rajistics

# Cross Encoder / Reranker



https://www.mongodb.com/resources/basics/artificial-intelligence/reranking-models

# Cross Encoder / Reranker



https://www.mongodb.com/resources/basics/artificial-intelligence/reranking-models

@rajistics

# Cross Encoder / Reranker



Boosting accuracy by adding reranker in Llama 3.1 70B powered RAG

https://developer.nvidia.com/blog/how-using-a-reranking-microservice-can-improve-accuracy-and-costs-of-information-retrieval/

@rajistics

# Cross Encoder / Reranker



Execution Flow (Approximate)

0s   1s   2s   3s   4s   5s   6s   7s   8s

| 1s | 0.99s | 1.19s | 2.33s | 2.44s |

**1. Query Setup** (~0.0s – 1.0s)
Initial query processing and reformulation
Sequential
**1s**
12.3%

**2. Query Expansion** (~1.0s – 2.0s)
Query expansion and inference
Sequential
**0.99s**
12.1%

**3. Embedding Generation** (~2.0s – 2.0s)
Converting query to vector embeddings
Sequential
**0.06s**
0.7%

**4. Search & Retrieval** (~2.0s – 3.2s)
Index search + parallel document retrieval
Parallel operations
**1.19s**
14.6%

**5. Reranking** (~3.2s – 5.6s)
Multiple reranking passes for relevance
Overlapping phases
**2.33s**
28.6%

**6. Response Generation** (~5.6s – 8.0s)
Generating and streaming the final response
Includes streaming
**2.44s**
29.9%

**7. Final Processing** (~8.0s – 8.2s)
Attribution, groundedness, and final steps
Final cleanup
**0.14s**
1.7%

# Hands On: Retriever & Reranker



Retrieve & Re-Rank Demo over Simple Wikipedia

This examples demonstrates the Retrieve & Re-Rank Setup and allows to search over Simple Wikipedia.

You can input a query or a question. The script then uses semantic search to find relevant passages in Simple English W smaller and fits better in RAM).

For semantic search, we use `SentenceTransformer('multi-qa-MiniLM-L6-cos-v1')` and retrieve 32 potentia answer the input query.

Next, we use a more powerful CrossEncoder (`cross_encoder = CrossEncoder('cross-encoder/ms-marco-M` that scores the query and all retrieved passages for their relevancy. The cross-encoder further boost the performance, e search over a corpus for which the bi-encoder was not trained for.

https://colab.research.google.com/drive/1lRr0J5fumRBP-RmTm5kD9lMd9nuOlhml?authuser=1#scrollTo=HETfTO4P-otj

# Instruction Following Reranker



The first instruction-following reranker from Contextual AI

@rajistics

**Current Reranker Instruction:**
Default ranking

**#1 Consumer Guide Review** — Score: 0.94
Dec 15, 2024 — Product Review — Professional
The BlendMaster 3000 earned our top safety rating with no reported issues during extensive testing.

**#2 HomeGoods Safety Alert** — Score: 0.87
Feb 25, 2025 — Safety Notice — Official
RECALL: BlendMaster 3000 models with serial numbers starting with BM3-25 have faulty wiring that can cause fires.

**#3 BlendMaster Support Forums** — Score: 0.73
Jan 30, 2025 — User Report — Community
Some users report overheating in the base after 30+ minutes of continuous use.

@rajistics

# Combine Multiple Retrievers

Technical Documents:
- BM25
- BGE (Gemma-2)
- E5 Mistral (7B)
- Voyager-large-2



Best Nuggets (oracle retrieval): 0.719

Best Answer (oracle retrieval): 0.683

Recall@50

BM25: 0.200
BM25 + Rerank: 0.244
BGE (Gemma-2): 0.419
BGE + Rerank: 0.471
E5 Mistral (7B): 0.359
E5 + Rerank: 0.426
Voyage-large-2: 0.458
Voyage + Rerank: 0.511
Fusion (4 models): 0.505
Fusion + Rerank: 0.545

FreshStack: https://arxiv.org/pdf/2504.13128

# Cascading Rerankers in Kaggle



https://www.youtube.com/watch?v=Bnn2m4S22T4

@rajistics

# Best practices



@rajistics

# Families of Embedding Models

A quick taxonomy to orient choices (speed, accuracy, and how they're used).

Faster                                                                                    More accurate (at match quality)

Model2Vec (Static)    DistilBERT (First-gen)    Sentence-Tfrs (Bi-encoder)    OpenAI Embeddings    BGE / Instructor (LLM-based)    Cross-Encoders (Rerankers)

## Model2Vec (Static Embeddings)

- Token-level vectors; fixed vocabulary
- Very fast; tiny footprint
- Weak on synonyms/negation
- Okay for simple similarity

## DistilBERT (First generation)

- CLS/mean pooling for sentences
- General-purpose; light footprint
- Better semantics than static vectors
- Baseline for small/medium use

## Sentence Transformers (Sentence-based)

- Siamese training for pair similarity
- Great recall for search & clustering
- Index-friendly; fast at query time
- Solid default for retrieval

## Cross-Encoders | Reranker

- Token-level interaction per pair
- Highest accuracy for top-k rerank
- Slow O(k) — not full-corpus search
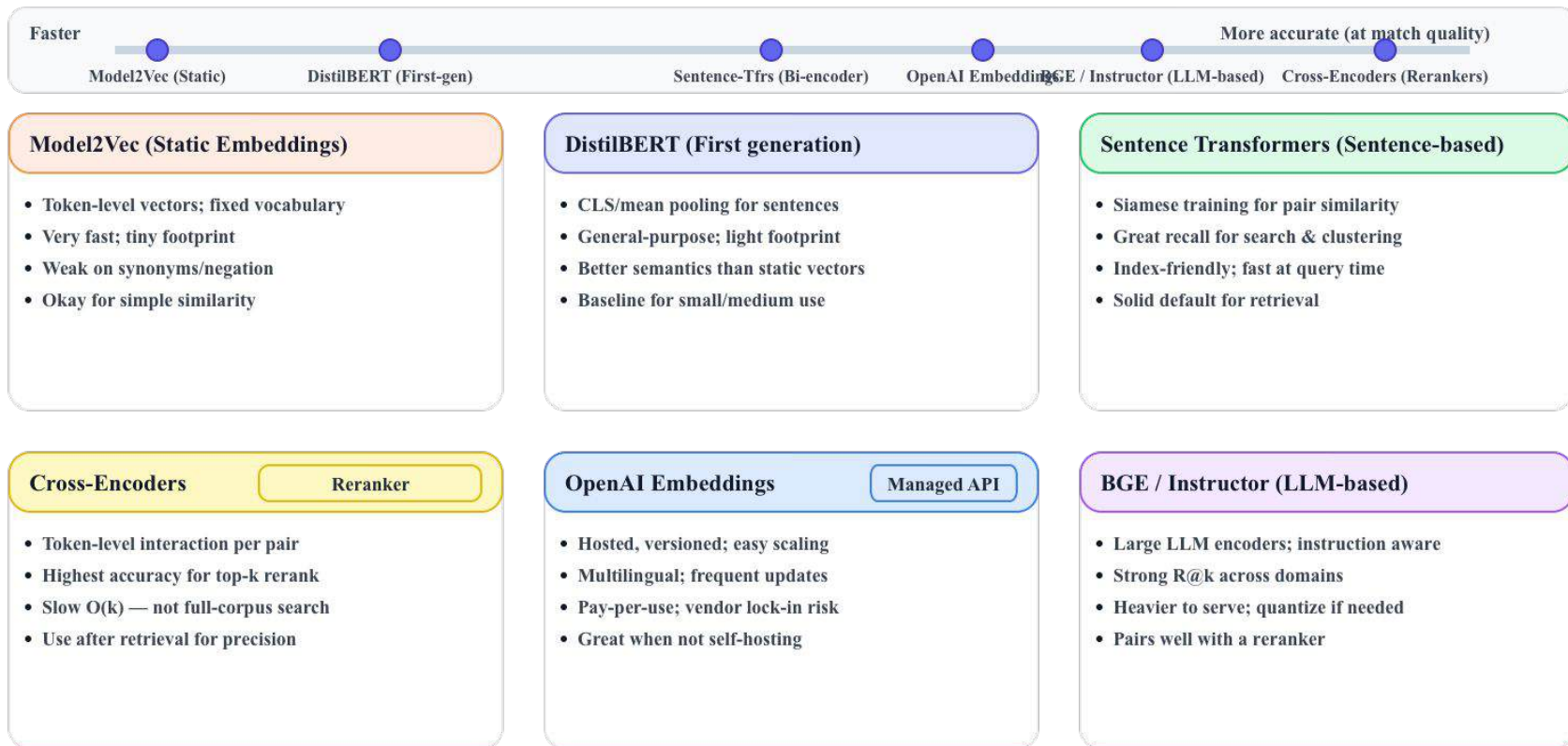- Use after retrieval for precision

## OpenAI Embeddings | Managed API

- Hosted, versioned; easy scaling
- Multilingual; frequent updates
- Pay-per-use; vendor lock-in risk
- Great when not self-hosting

## BGE / Instructor (LLM-based)

- Large LLM encoders; instruction aware
- Strong R@k across domains
- Heavier to serve; quantize if needed
- Pairs well with a reranker

@rajistics

# Lots of New Models



John Hopkins
University



IBM



Google

@rajistics

# Other retrieval methods:

- SPLADE for sparse
- ColBERT Late Interaction
- GraphRAG
- Many RAG flavors

# **Operational Concerns:**

## Computing Embeddings

## Storing Embeddings

**Nearest Neighbor Search**
Finding similar vectors in 10M embeddings

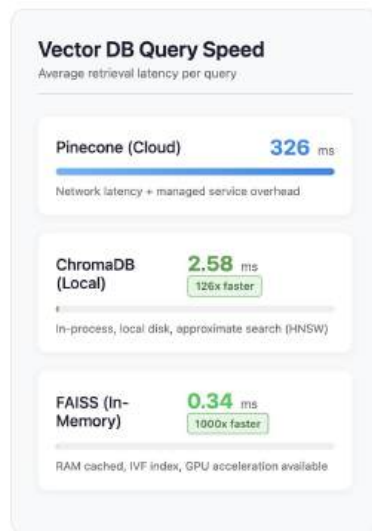Sklearn Brute Force     **25** seconds

Pure Python, O(n) linear scan, no optimization

FAISS IndexFlatL2     **1.7** seconds   [18x faster]

C++ optimized, SIMD instructions, CPU parallelization

⚠ This is for exact search without any indexing – just raw compute optimization

**Vector DB Query Speed**
Average retrieval latency per query

Pinecone (Cloud)     **326** ms

Network latency + managed service overhead

ChromaDB (Local)     **2.58** ms   [126x faster]

In-process, local disk, approximate search (HNSW)

FAISS (In-Memory)     **0.34** ms   [1000x faster]

RAM cached, IVF index, GPU acceleration available

# Vector Database Options

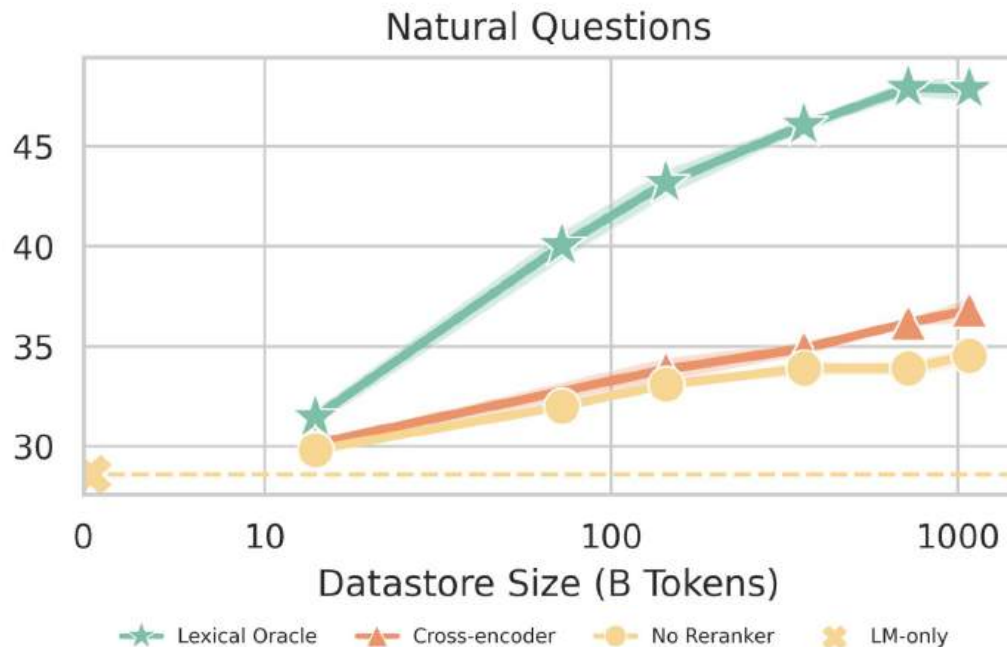## Vector Database Layered Storage Architecture

Storage tier optimization and technical solution configuration based on latency requirements

| Storage Tier | Latency Requirements | Core Application Scenarios | Technical Solutions |
|---|---|---|---|
| Hot Data Layer | < 50ms | Real-time search / Intelligent recommendations / Targeted advertising | Traditional specialized vector databases (e.g., Milvus hot instances) |
| Warm Data Layer | 50-500ms | Standard RAG dialogue / Multi-tenant shared services | S3Vector / Milvus three-tier storage instances |
| Cold Data Layer | > 500ms | Historical data archiving / Offline data analysis | S3+Spark/Daft / Milvus vector data lake |

Note: This architecture demonstrates the latency requirements, applicable scenarios, and corresponding technical solutions for different storage tiers in vector database layered storage. The hot→warm→cold tier design optimizes overall storage costs and performance.

https://zilliz.com/blog/will-amazon-s3-vectors-kill-vector-databases-or-save-them

@rajistics

# Operational Concerns:

As datastores get bigger, you need to work on improving retrieval performance

## Natural Questions



Trillion Token: https://arxiv.org/pdf/2407.12854

# Search Strategy Comparison

Watch how different approaches explore the solution space

## Traditional RAG

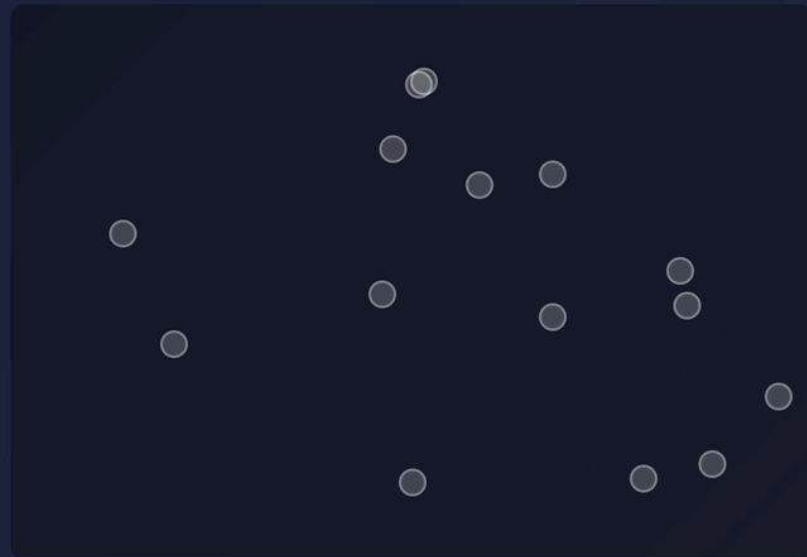Ready to search...

**0**
Queries Made

**0ms**
Time Taken

## Agentic RAG

Ready to explore...

**0**
Queries Made

**0ms**
Time Taken

# Tools use / Reasoning

Use reasoning models to keep using queries until satisfied



Scientific research

```python
scientific_research.py

from agno.agent import Agent
from agno.models.openai import OpenAIChat

task = (
    "Read the following abstract of a scientific paper and provide a critical evalu
    "results, conclusions, and any potential biases or flaws:\n\n"
    "Abstract: This study examines the effect of a new teaching method on student p
    "A sample of 30 students was selected from a single school and taught using the
    "The results showed a 15% increase in test scores compared to the previous seme
    "The study concludes that the new teaching method is effective in improving math
)
reasoning_agent = Agent(
    model=OpenAIChat(id="gpt-5-mini-2024-08-06"), reasoning=True, markdown=True
)
reasoning_agent.print_response(task, stream=True, show_full_reasoning=True)
```
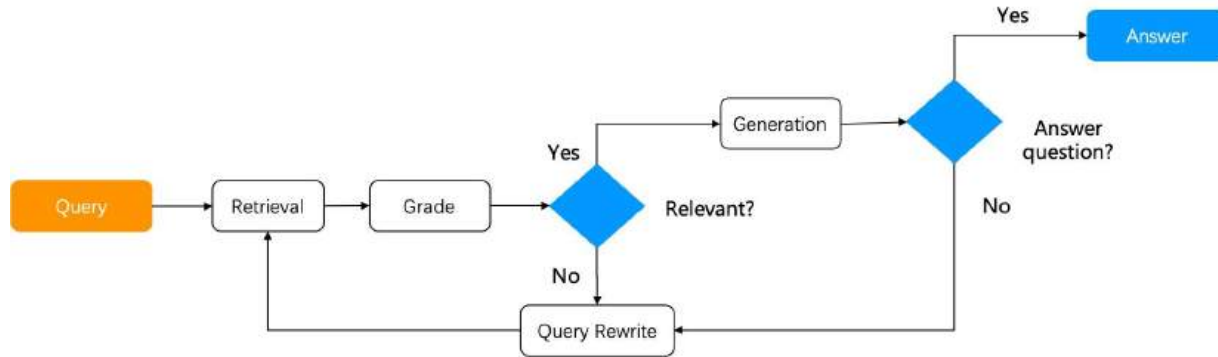
https://docs.agno.com/concepts/reasoning/reasoning-agents

@rajistics

# Agentic RAG



https://ragflow.io/blog/agentic-rag-definition-and-low-code-implementation

@rajistics

# Tools use / Reasoning

**Message**

How to make the published changes draft?

**Reasoning step 1**

**Plan searches and ambiguity**
Action: make a tool call to run multiple context

**Reasoning step 2**

**Expand searches to other features and limitation**
Action: make more contextual queries in parallel

**Reasoning step 3**

**Assess gathered info and spot issues**
We gathered info covering blog revert to draft, draft limitation, store product visibility, email events result included an incorrect step referen able to revert to draft after publish—need to va not supported and workarounds like noindex and p published states. Confirm that restoring site hi like Editor X/Studio Editor similar. We should r for pages draft, events, blog, site history limi feature called "make published changes draft" in

**Reasoning step 4**

**Synthesize answer structure**
Action: Prepare final answer with bullets and mi question at end.

**Tool Calls**

- think(title=Plan searches and ambiguity, thought=Th Blog post: unpublish/move to draft; (2) Site-wide cha (Content Manager) items: move published item back to revert published, only duplicate; (5) Wix Stores prod doesn't have a page draft toggle, but there's "Hide f "Unpublish a page" in some contexts; Editor allows "H published—there is a "Show in menu" vs "Hide"; Also " including limitations and workarounds. I'll run multi contextual queries in parallel, confidence=0.7)
- get_contextual_results(query=Wix revert published c
- get_contextual_results(query=Wix blog move publishe
- get_contextual_results(query=Wix Editor undo publis history)
- get_contextual_results(query=Wix CMS items change s Manager status draft published)
- get_contextual_results(query=Wix page draft feature limitations cannot save page as draft)
- think(title=Expand searches to other features and l unpublish a site, unpublish a page, Wix Stores produc regular pages don't have a draft state. Also search f Site History. Then compile concise answer asking whic more contextual queries in parallel, confidence=0.67)
- get_contextual_results(query=Wix unpublish site how
- get_contextual_results(query=Wix unpublish a page d unpublish page 2024)
- get_contextual_results(query=Wix Stores unpublish p stock or hidden)
- get_contextual_results(query=Wix Email Marketing re Wix)
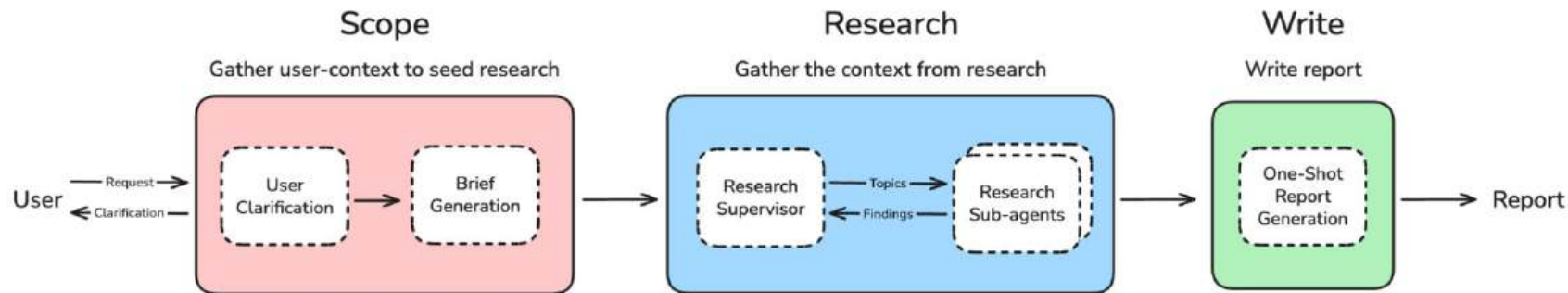- get_contextual_results(query=Wix Events unpublish event draft event status unp

**Response (176.8s)**

It depends on what you published. Here are the common cases in Wix and what y

- Blog posts
  - You can revert a published post to Draft.
  - Steps: Dashboard > Blog > open the post > More/Expand (…) > Revert to D
  - Note: If you restore content via post history, SEO settings/categories/ updated manually.
- Entire site changes
  - Use Site History to restore a previous saved/published version.
  - Steps: Editor or Dashboard > Site History > pick a version > Restore.
  - Important limits: Restoring does NOT roll back Wix app content (e.g., B items). Those apps manage their own data.
- Regular pages
  - Not supported: You cannot "unpublish" just one page or make a page a dr
  - Workarounds:
    - Hide from menu (page remains live if someone has the link).
    - Restrict access (password-protect or Members Only).
    - Noindex the page in SEO settings to keep it out of search.
    - Keep working changes unpublished by not clicking Publish (changes sa
- CMS (Content Manager) collection items
  - Use item visibility instead of "draft."
  - Steps: Dashboard > CMS > open collection > enable Control item visibili to Visible or Hidden; you can also Schedule visibility.
- Wix Events
  - After publishing an event, you cannot revert it to Draft.
  - Workaround: Control which events display on pages (e.g., use "Select Ma published event doesn't show).
- Wix Stores products
  - Hide products instead of "drafting" them, or schedule visibility.
  - Steps: Dashboard > Store Products > hover product > Hide; or use Store
- Email Marketing campaigns

🔬 **Open Deep Research**

https://github.com/langchain-ai/open_deep_research

@rajistics

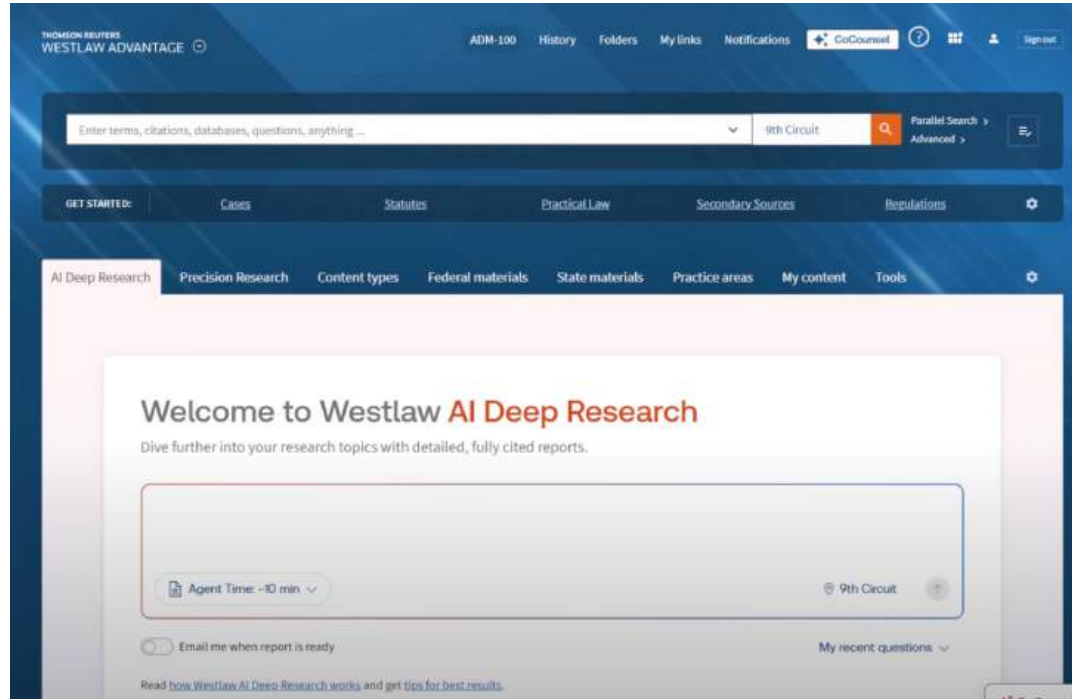# DeepResearch Bench

100 PhD-level research tasks

| Rank | model | overall | comp. | insig |
|------|-------|---------|-------|-------|
| 1 | 🚀 gemini-2.5-pro-deepresearch | 49.71 | 49.51 | 49.45 |
| 2 | 🚀 openai-deepresearch | 46.45 | 46.46 | 43.73 |
| 3 | 🚀 claude-research | 45 | 45.34 | 42.79 |
| 4 | 🚀 kimi-researcher | 44.64 | 44.96 | 41.97 |
| 5 | 🚀 doubao-deepresearch | 44.34 | 44.84 | 40.56 |
| 6 | 🚀 langchain-open-deep-research | 43.44 | 42.97 | 39.17 |
| 7 | nvidia-aiq-research-assistant | 40.52 | 37.98 | 38.39 |

| Deep Research Bench Submission | c0a160b | openai:gpt-4.1-nano | openai:gpt-4.1 | openai:gpt-4.1 | $87.83 | 207,005,549 |
|---|---|---|---|---|---|---|

https://huggingface.co/spaces/Ayanami0730/DeepResearch-Leaderboard

# Westlaw AI Deep Research



https://www.youtube.com/watch?v=tvpH36uT7hw

# Agentic RAG



Self RAG: https://arxiv.org/pdf/2310.11511

@rajistics

# Agentic RAG

https://www.reddit.com/r/LangChain/comments/1njmb1r/i_taught_my_retrievalaugmented_generation_system/

# Agentic RAG



Best-Information2493 OP · 1d ago

You're absolutely right about the inefficiency!

**Non-relevant docs:** The system usually tries to rewrite and retrieve again, but it should have better fallbacks like using pure LLM knowledge or graceful exit.

**Resource waste:** Going through the full pipeline just to restart is brutal. Better approaches would be:

- Early stopping at each step,

- Circuit breakers to prevent endless loops,

- Caching intermediate results

The paper prioritizes accuracy over efficiency real production systems definitely need smarter resource management.
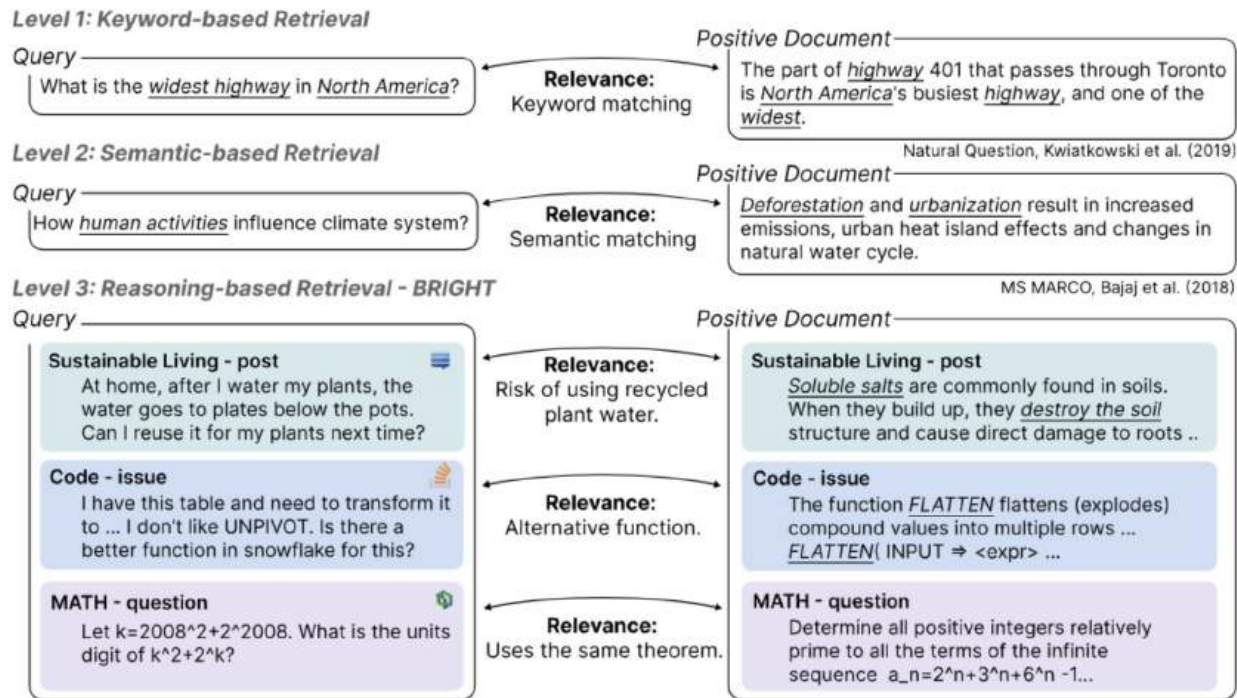
⬆ 2 ⬇    💬 Reply    🏅 Award    ↗ Share    …

https://www.reddit.com/r/LangChain/comments/1njmb1r/i_taught_my_retrievalaugmented_generation_system/
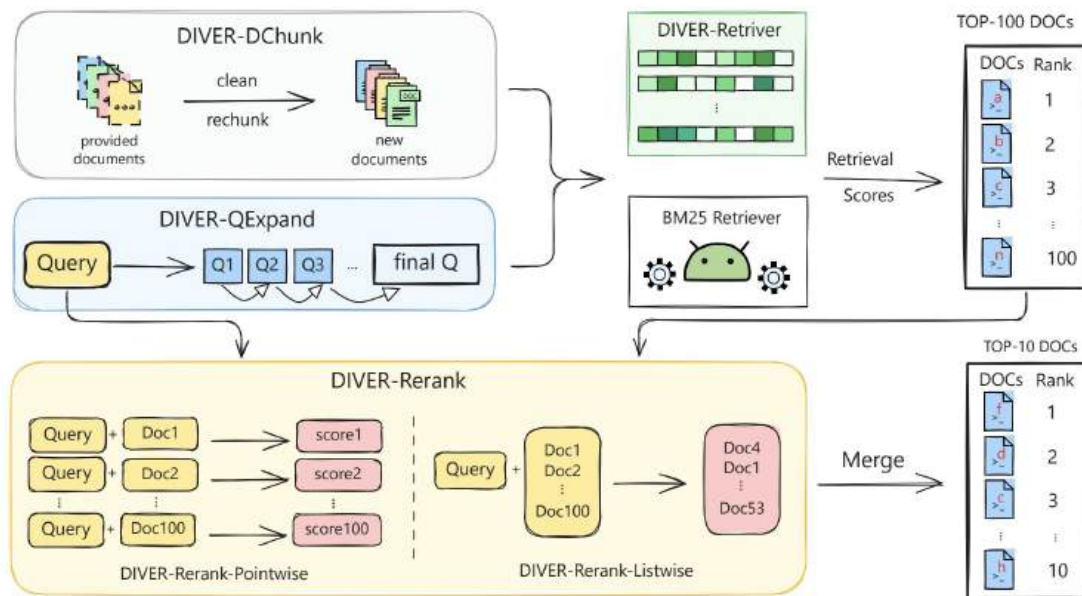
# Research: BRIGHT

Analyzing retrieval reasoning



**Level 1: Keyword-based Retrieval**
Query
What is the *widest highway* in *North America*?

**Relevance:** Keyword matching

Positive Document
The part of *highway* 401 that passes through Toronto is *North America*'s busiest *highway*, and one of the *widest*.

Natural Question, Kwiatkowski et al. (2019)

**Level 2: Semantic-based Retrieval**
Query
How *human activities* influence climate system?

**Relevance:** Semantic matching

Positive Document
*Deforestation* and *urbanization* result in increased emissions, urban heat island effects and changes in natural water cycle.

MS MARCO, Bajaj et al. (2018)

**Level 3: Reasoning-based Retrieval - BRIGHT**
Query

**Sustainable Living - post**
At home, after I water my plants, the water goes to plates below the pots. Can I reuse it for my plants next time?

**Relevance:** Risk of using recycled plant water.

Positive Document

**Sustainable Living - post**
*Soluble salts* are commonly found in soils. When they build up, they *destroy the soil* structure and cause direct damage to roots ..

**Code - issue**
I have this table and need to transform it to ... I don't like UNPIVOT. Is there a better function in snowflake for this?

**Relevance:** Alternative function.

**Code - issue**
The function *FLATTEN* flattens (explodes) compound values into multiple rows ... *FLATTEN*( INPUT ⇒ <expr> ...

**MATH - question**
Let $k=2008^2+2^{2008}$. What is the units digit of $k^2+2^k$?

**Relevance:** Uses the same theorem.

**MATH - question**
Determine all positive integers relatively prime to all the terms of the infinite sequence $a_n=2^n+3^n+6^n -1$...

BRIGHT: https://arxiv.org/pdf/2407.12883

# BRIGHT #1: DIVER

Reasoning-intensive

Information Retrieval



https://arxiv.org/pdf/2508.07995

# BRIGHT #1: DIVER

Reasoning-
intensive

Information
Retrieval

Table 1: Prompts used in DIVER-QExpand for query expansion. Braces {} denote placeholders.

| Prompt Stage | LLM Instruction |
| --- | --- |
| First Round | Given a query and the provided passages (most of which may be incorrect or irrelevant), identify helpful information from the passages and use it to write a correct answering passage. Use your own knowledge, not just the example passages! Query: {query} Possible helpful passages: {top-k retrieved documents} |
| Subsequent Rounds | Given a query, the provided passages (most of which may be incorrect or irrelevant), and the previous round's answer, identify helpful information from the passages and refine the prior answer. Ensure the output directly addresses the original query. Use your own knowledge, not just the example passages! Query: {query} Possible helpful passages: {top-k retrieved documents} Prior generated answer: {last-round expansion} |

https://arxiv.org/pdf/2508.07995

# Agentic RAG on WixQA

Pick:
- Accuracy
- Latency

(6s versus 50s)

# Rethink your Assumptions

Querying with
LLM using
BM25



BRIGHT: https://arxiv.org/pdf/2407.12883

# Agentic RAG with BM25

# Agentic RAG for Code Search

- **Claude Code (Lexical / Iterative Search)**

- Keeps searching (like grep) until it finds or rules out a function/dependency

https://x.com/pashmerepat/status/1926717705660375463
https://www.tigerdata.com/blog/why-cursor-is-about-to-ditch-vector-search-and-you-should-too#reading-the-cursor-tea-leaves

@rajistics



Jacky Liang · You
📐 always learning // ai @ tigerdata && founder answerhq.co // pinecone,...
6d · 🌐

I have a bold prediction.

Cursor is going to rip out their entire vector search implementation, and replace it with pure lexical (a smarter-sound way of saying keyword) search akin to Anthropic Claude Code's implementation

Claude Code specifically uses grep, find, and other exact file/text search commands.

Note that Cursor already uses lexical search tool calling in its Cursor Agent product, but it is nowhere near as good as Claude Code's.

If Cursor does rip out their entire vector search implementation, this is a major loss of a large customer for turbopuffer, which powers Cursor's code vector search infrastructure.

The facts:

- the incredible Claude Code uses ONLY lexical search (no vector search) for context discovery, which is leaps and bounds better than Cursor's
- Boris Cherny and Catherine Wu, the chief architects of Claude Code (and

# Combine Retrieval Approaches

Response
Guardrail:
2 Tier System of
Text + LLM



Output Guardrail System

https://careersatdoordash.com/blog/large-language-modules-based-dasher-support-automation/
UAR: https://arxiv.org/html/2406.12534v1

@rajistics

# Hands on: Agentic RAG



> ## ∨ Agentic RAG with Hugging Face smolagents vs Vanilla RAG
>
> Author: @MariaKhalusova
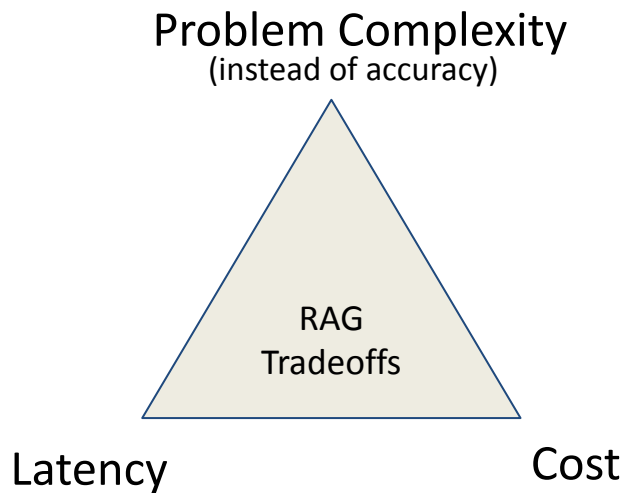>
> Last updated: Jan 9th, 2025
>
> ∨ What you'll learn:
>
> 1. Parsing PDF documents from S3 into DataStax AstraDB with Unstructured Platform
> 2. Building Vanilla RAG in pure Python without using specialized frameworks
> 3. Differences between Vanilla RAG and Agentic RAG
> 4. Creating Agentic RAG with Hugging Face `smolagents` library
> 5. Whether Agentic RAG can produce better answers (spoiler: it can!)
>
> In Vanilla RAG, your system uses the user's question to perform a single retrieval step and get a batch of documents that are r
> relevant to the query. These documents are then passed on to the LLM to generate an answer grounded in the context of thos
>
> However, this approach has limitations. If the results of the retrieval are inadequate (either irrelevant or incomplete), this will h
> negative impact on generation. There are many different methods one can employ to improve the retrieval quality, such as cho
> embedding model, switching to a different retrieval method (e.g., BM25, or hybrid, metadata filtering, etc.), increasing the num
> documents, and adding a reranker. However, there may still be situations where a single retrieval step, or retrieving based on t
> "as is," may not produce optimal results.

Smolagents: https://colab.research.google.com/drive/1hG3dPgd8wjrO9wSD0K0Feo7EY1iXqrEN
Page Index: https://github.com/VectifyAI/PageIndex

# Solutions for a RAG Solution

Problem Complexity
(instead of accuracy)

RAG
Tradeoffs

Latency                    Cost

- High cost of mistakes + budget →
  Rerankers
- Need <5s latency → BM25 +
  Static Embeddings
- Complex multi-hop queries →
  Agentic RAG

# Retriever Checklist

- Keyword / BM25
- Semantic Search / Embedding Model
- Agentic / Reasoning LLM

**BM25**
Keyword-based retrieval
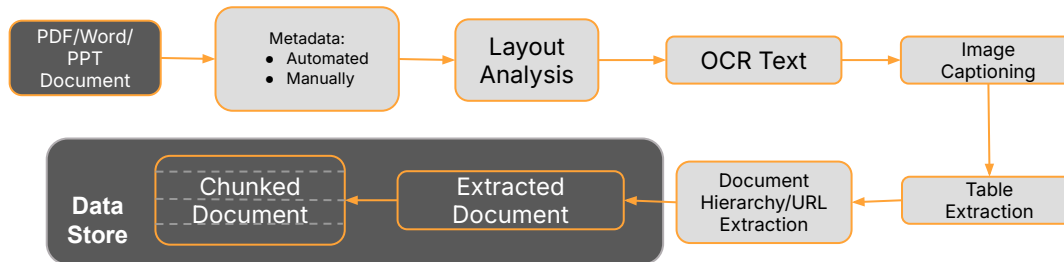
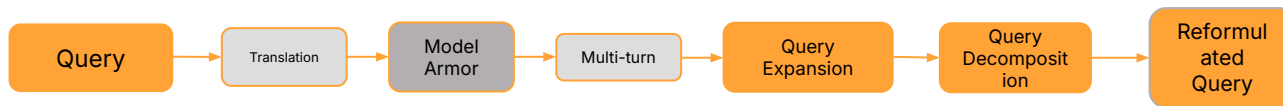**Language Models**
Semantic meaning with embeddings

**Agentic Search**
Dynamic using LLM Reasoning

@rajistics

**1. Parsing**

PDF/Word/PPT Document → Metadata: Automated, Manually → Layout Analysis → OCR Text → Image Captioning → Table Extraction → Document Hierarchy/URL Extraction → Extracted Document → Chunked Document (Data Store)

**2. Querying**

Query → Translation → Model Armor → Multi-turn → Query Expansion → Query Decomposition → Reformulated Query

**Retrieval**

**3. Retrieving**

Chunked Documents → Filter via Metadata → Semantic Search / Lexical Search → Reciprocal Rank Fusion → Instruction Following Reranker → Filter Model → Final Retrievals

**4. Generation**

Reformulated Query, Final Retrievals → Generate Response (GLM, Claude, GPT5, ..) → Translation → Attributions → Groundedness Scores → Final Responses UI/API

@rajistics

# RAG - Generation

Don't want a list of search results
So use a generation model

Greatest area of technical improvement for RAG
in the last few years

# RAG - Generation

Less interesting, because either choose
- Best generation model that fits your cost/latency budget
- Special needs
    - Low hallucination (Contextual GLM)
    - Domain Specific (Fine Tuned Healthcare LLM)
    - Language Specific
- Don't overindex on Context Window size -> Context Rot post



**Choose what's best for you!**
We added more models

Generation Settings
Parameters that affect response generation.

Generation Model
The model to use for generating responses.
Select generate model
✓ Contextual GLM
Gemini 2.5 Pro
Gemini 2.5 Flash
Gemini 2.0 Flash
Gemini 2.0 Flash Lite
Claude Opus 4
Claude Sonnet 4
GPT-5

Microsoft's total revenue for FY2024 was $245,122 million, with an operating income of $109,433 million. 3

**All models include:**
☑ - Inline attributions
☑ - Grounding checks

**Contextual AI GLM**
- Grounded answers

**Anthropic Opus 4**
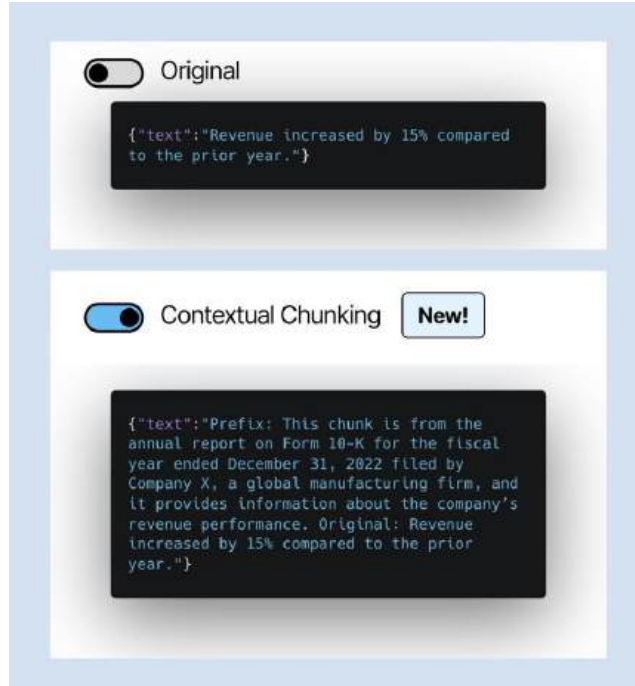- Deep Reasoning

**Google Gemini Pro 2.5**
- Long Form Content

**OpenAI GPT-5**
- Structured outputs & code

@RAJISTICS

@rajistics

# Chunking approaches

@rajistics

# From Vectors to Agents: Managing RAG in an Agentic World

Rajiv Shah
Chief Evangelist, Contextual AI
rajiv.shah@contextual.ai

https://github.com/rajshah4/LLM-Evaluation

@rajistics