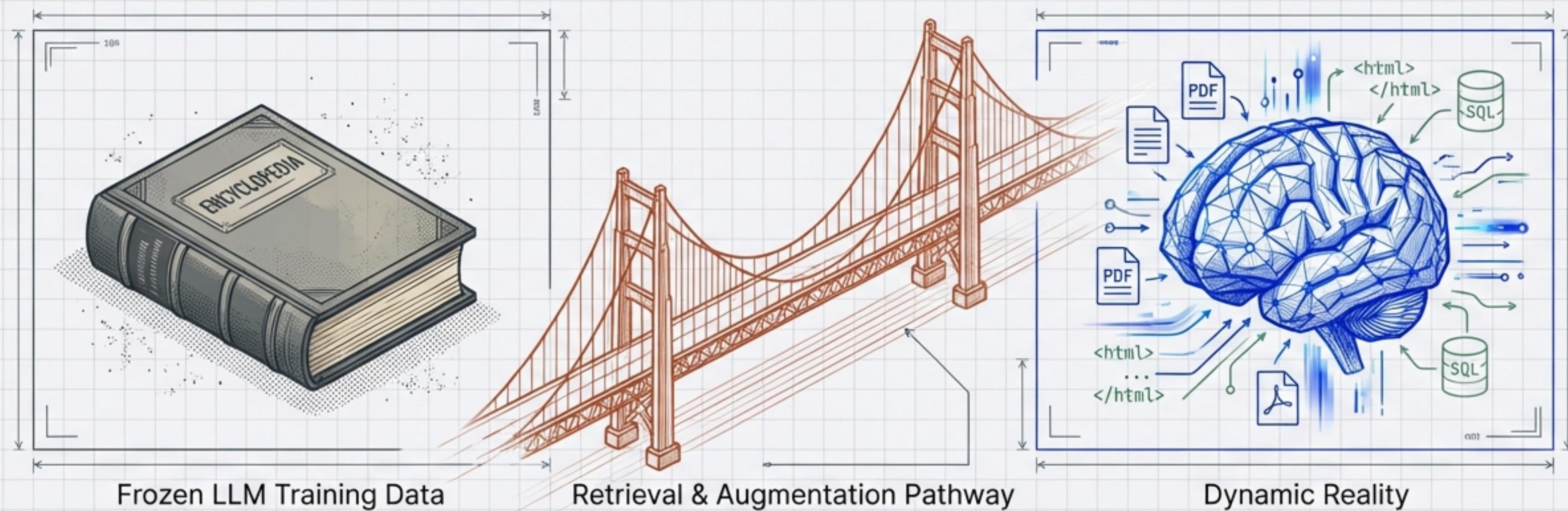


RAG: Bridging Static Knowledge and Dynamic Intelligence

A comprehensive guide to Retrieval Augmented Generation—Fundamentals, Pipelines, and Advanced Techniques



Source Material: Krish Naik, KodeKloud, freeCodeCamp

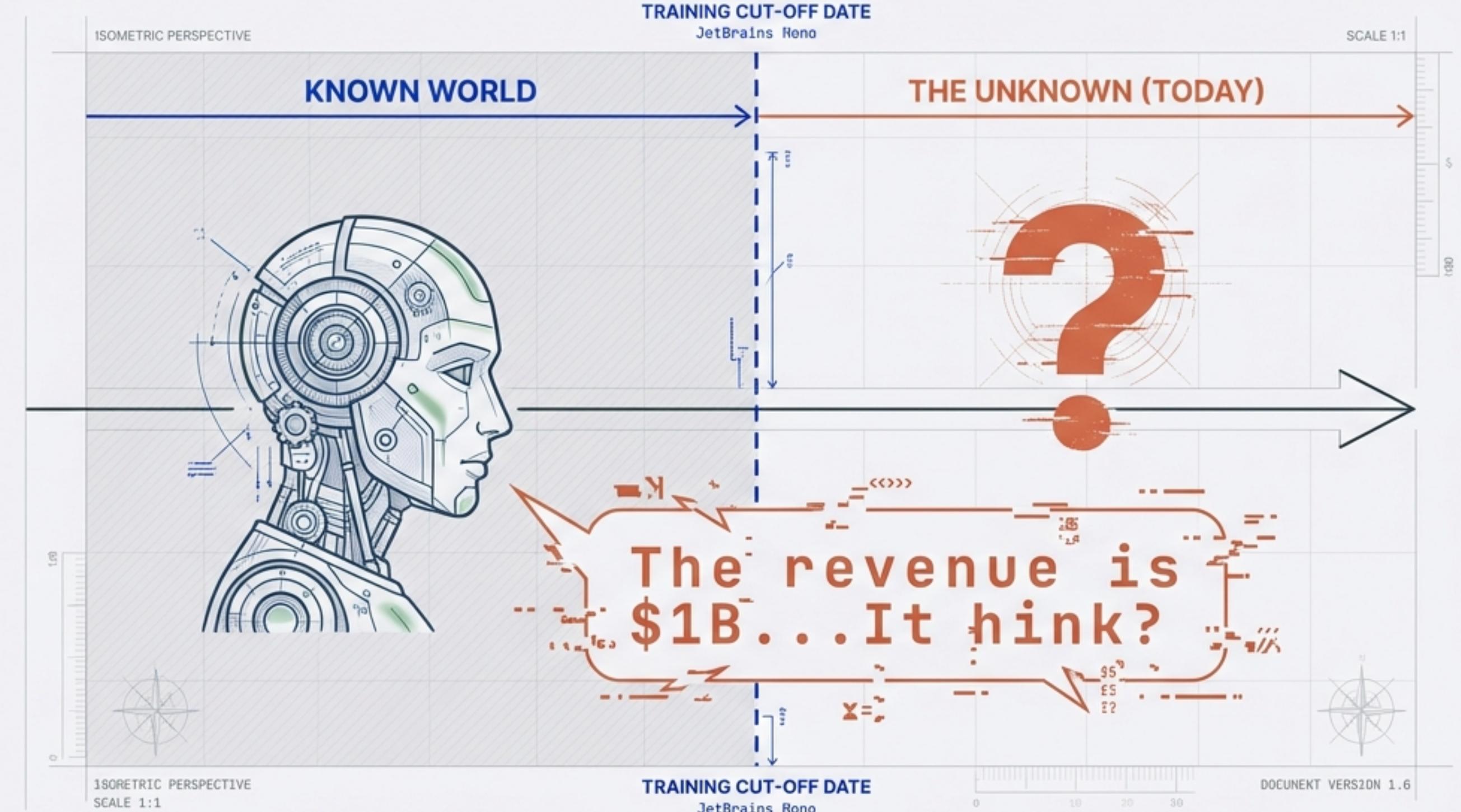


THE FROZEN BRAIN PROBLEM

WHY LLMS HALLUCINATE

LLMs are snapshots of the internet at a specific moment in time. They lack awareness of:

- Recent Events (Post-training history)
- Private Data (Company policies, financial records)
- User-Specific Context



RAG is the Open-Book Exam for AI

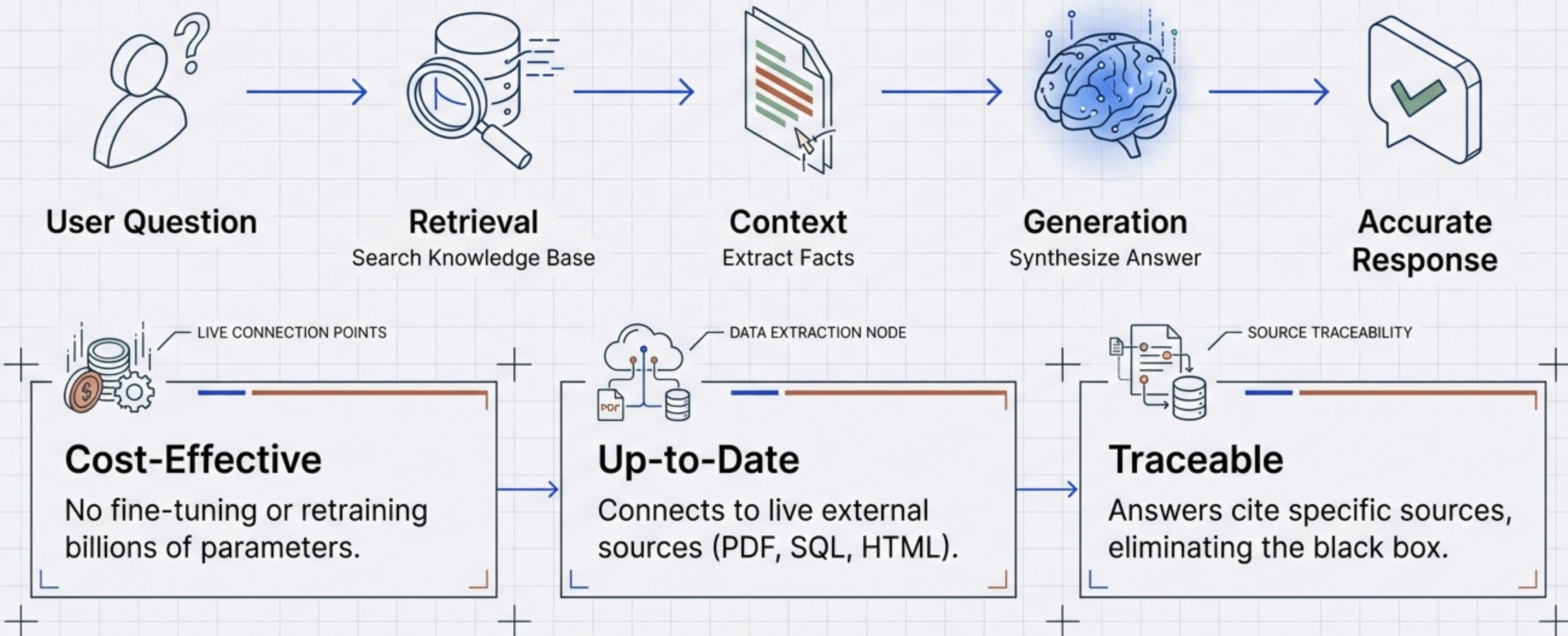


TECHNICAL BLUEPRINT / SWISS EDITORIAL

Helvetica Now Display

SCALE 1:1

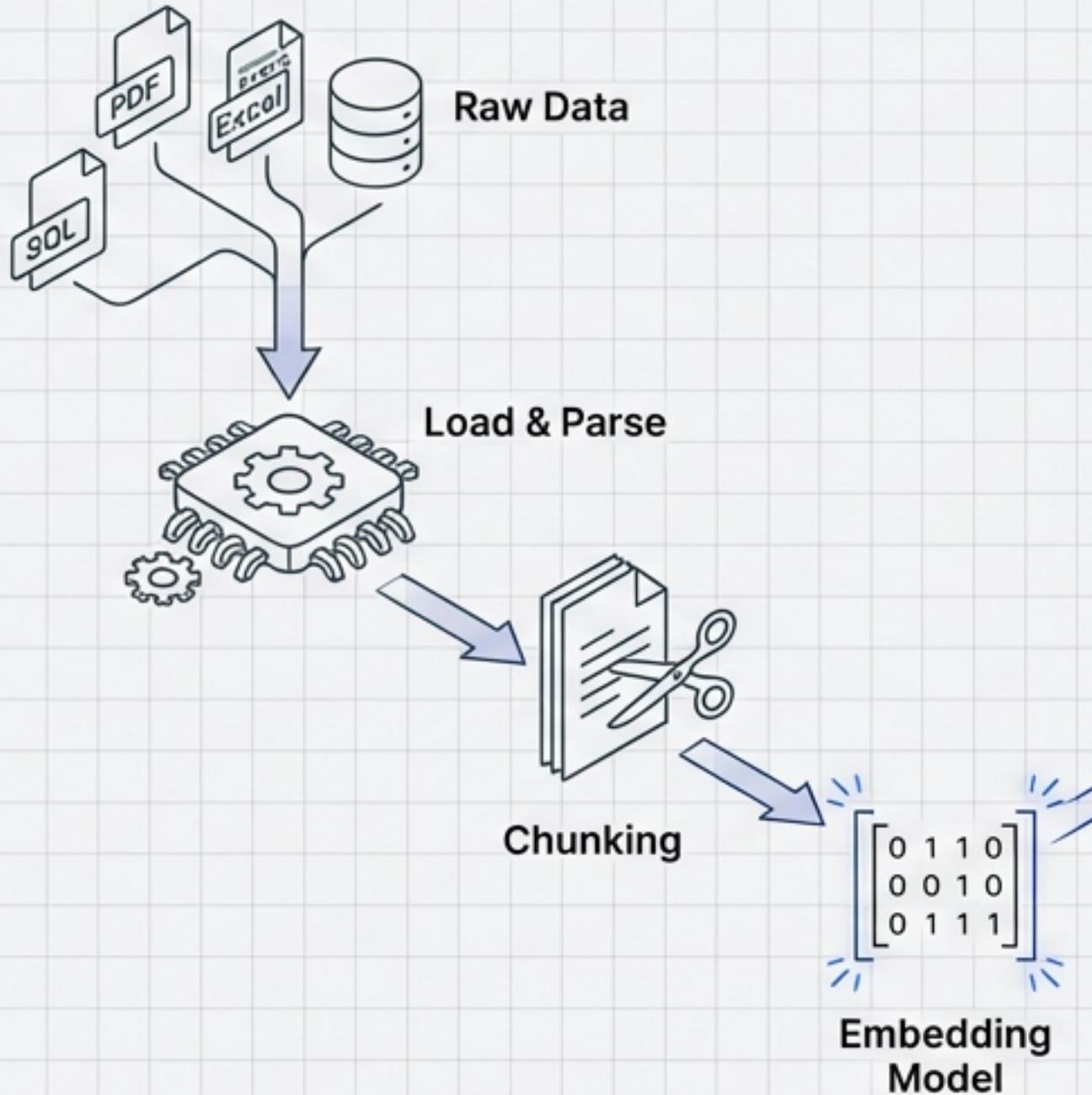
DOCUMENT VERSION 1.7



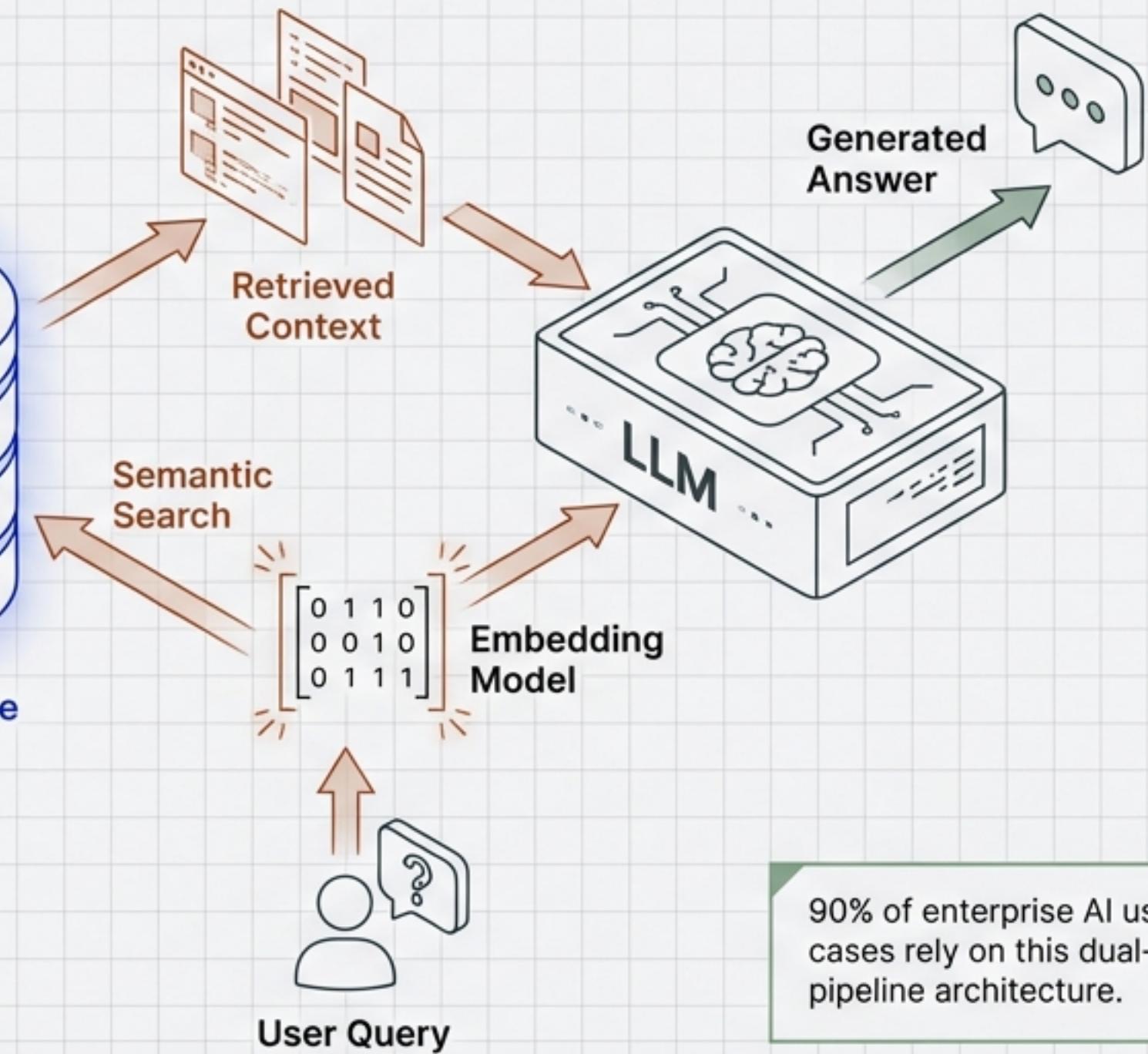


The Anatomy of a RAG System

Pipeline A: Ingestion (Offline/Async)



Pipeline B: Retrieval (Runtime)



90% of enterprise AI use cases rely on this dual-pipeline architecture.

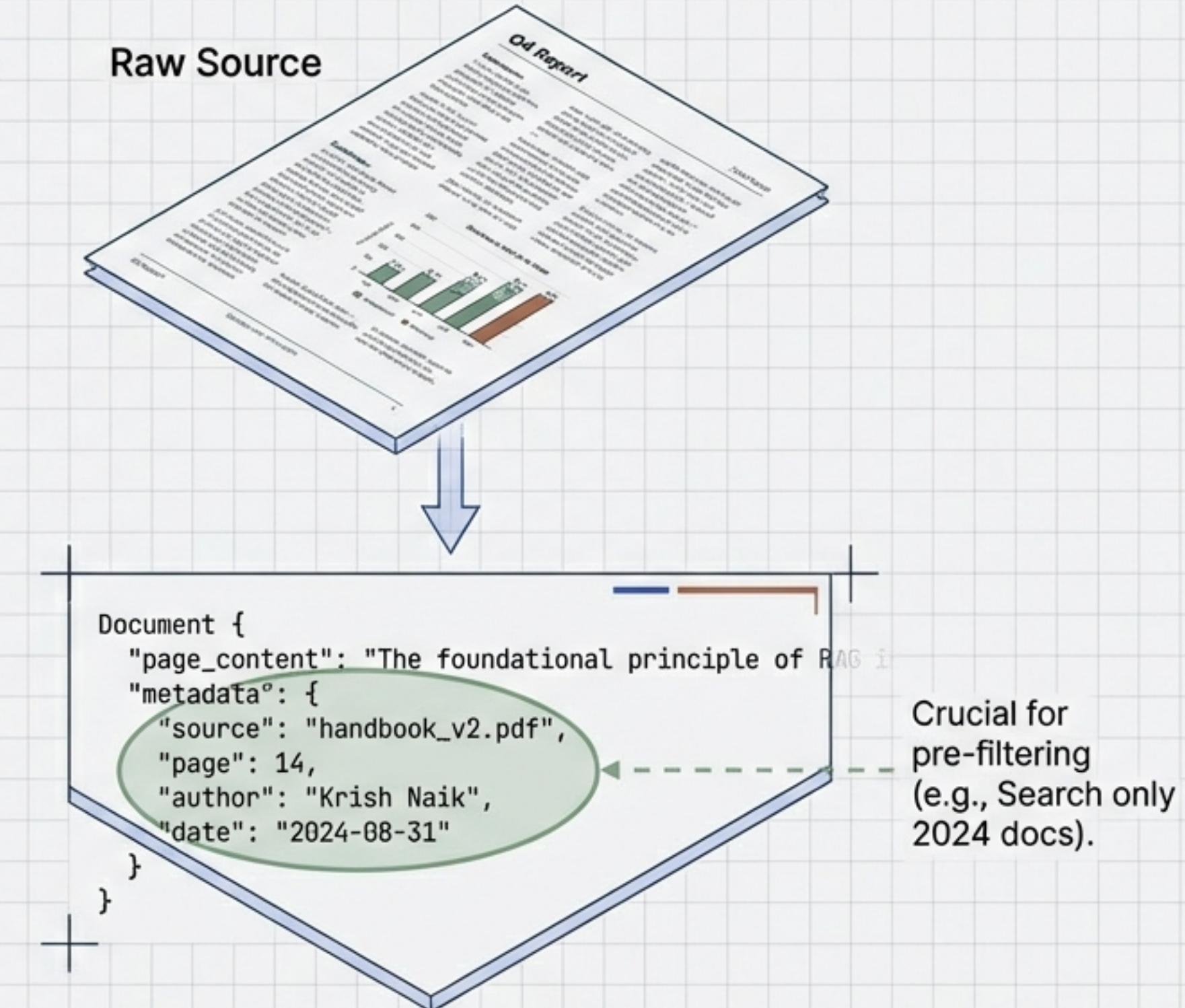


Pipeline A: Data Ingestion & Parsing

The foundation of RAG is extracting usable text from unstructured sources.

“Garbage in, garbage out.”

Tools: LangChain Document Loaders, PyPDF, PyMuPDF.

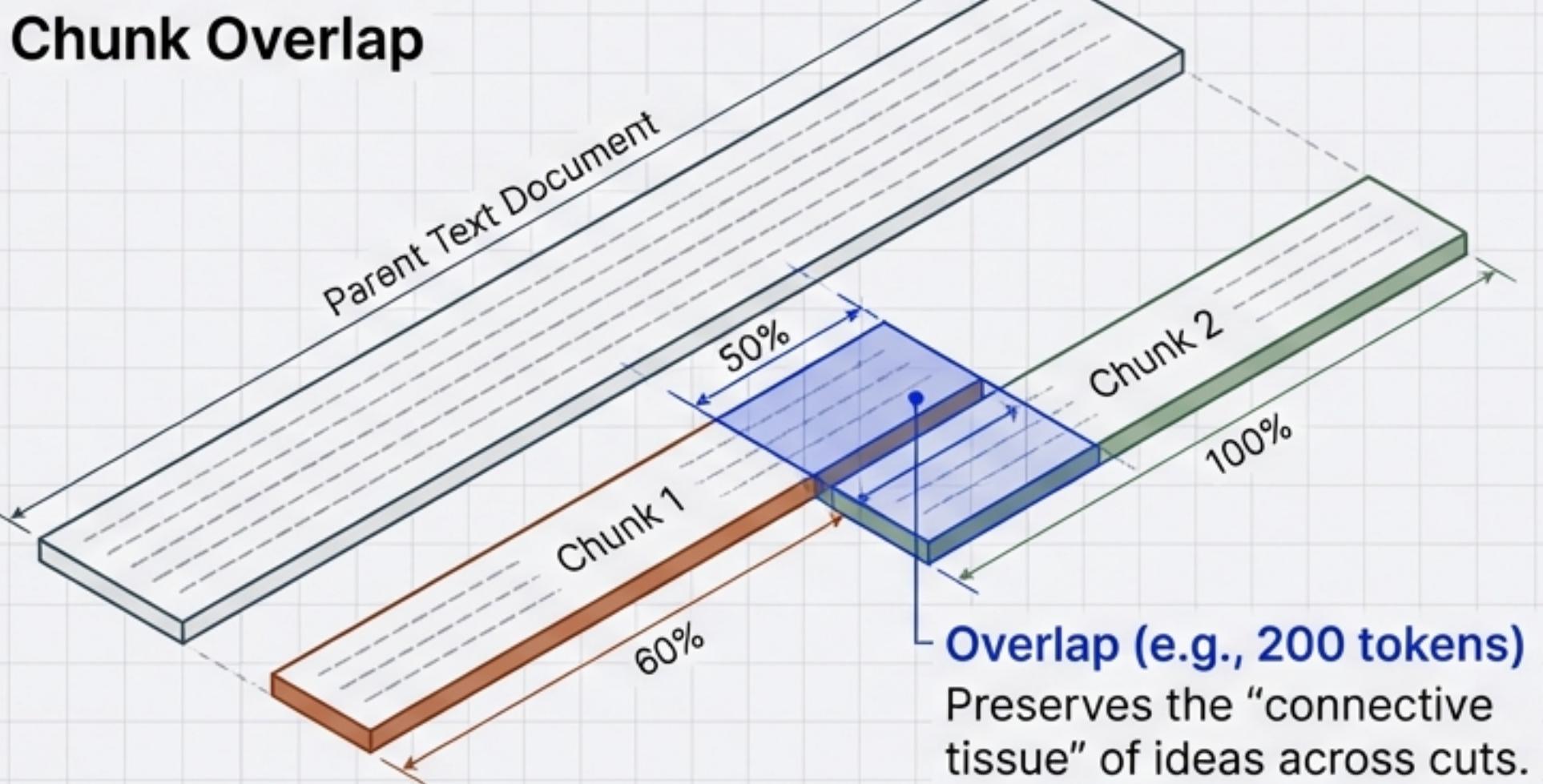


The Art of Chunking

Balancing Context and Precision

TECHNICAL BLUEPRINT / SWISS EDITORIAL
Helvetica Now Display
SCALE 1:1 DOCUMENT VERSION 2.0

LLMs have a fixed “Context Window.” We split text to fit relevant information into the prompt without losing meaning.



Strategy:
RecursiveCharacterTextSplitter

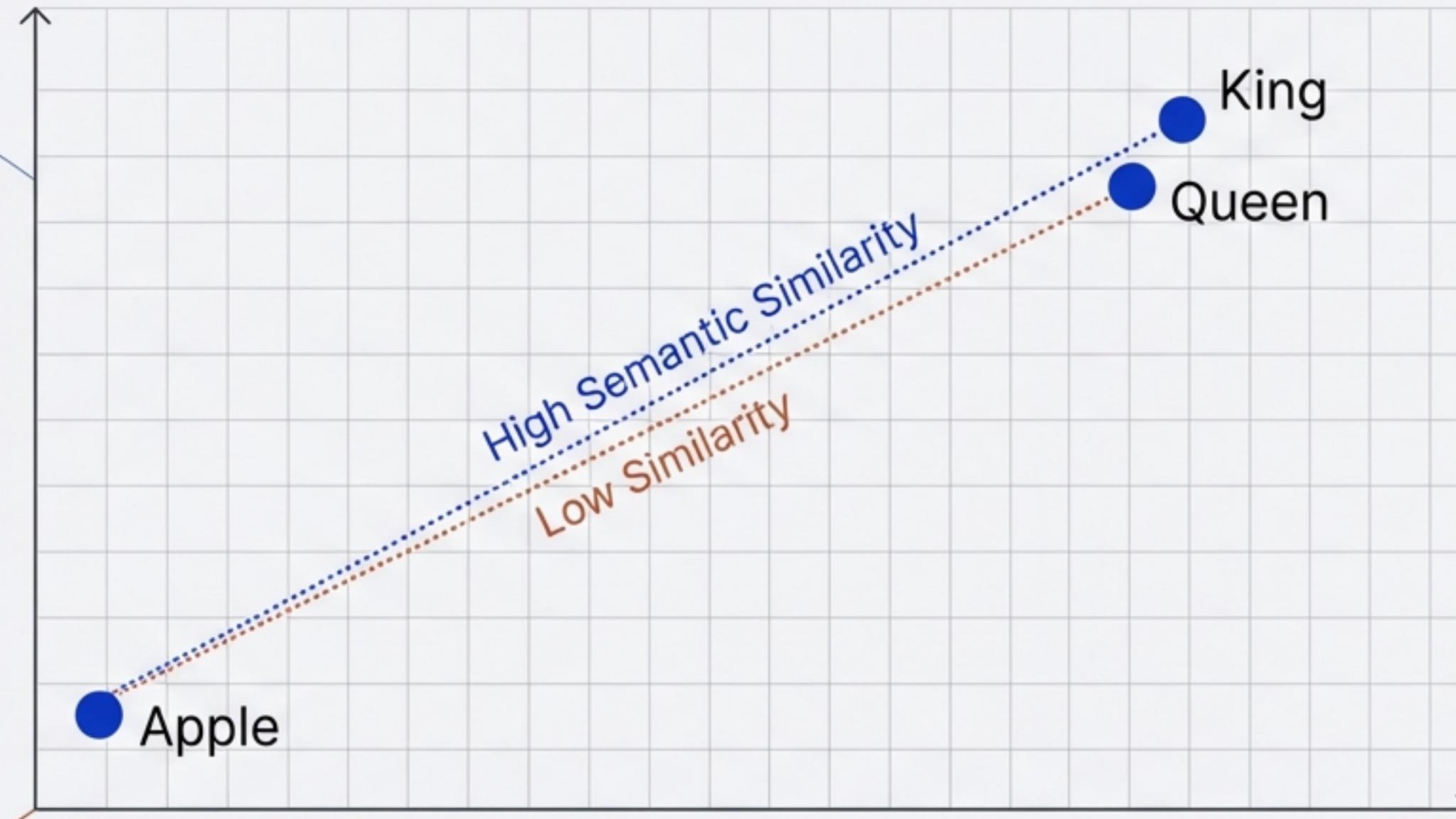
```
chunk_size = 1000  
chunk_overlap = 200
```

Embeddings: Converting Language to Math

Helvetica Now Display

SCALE
1:1

DOCUMENT VERSION
3.0



The computer doesn't read words.
It calculates distance between numbers.

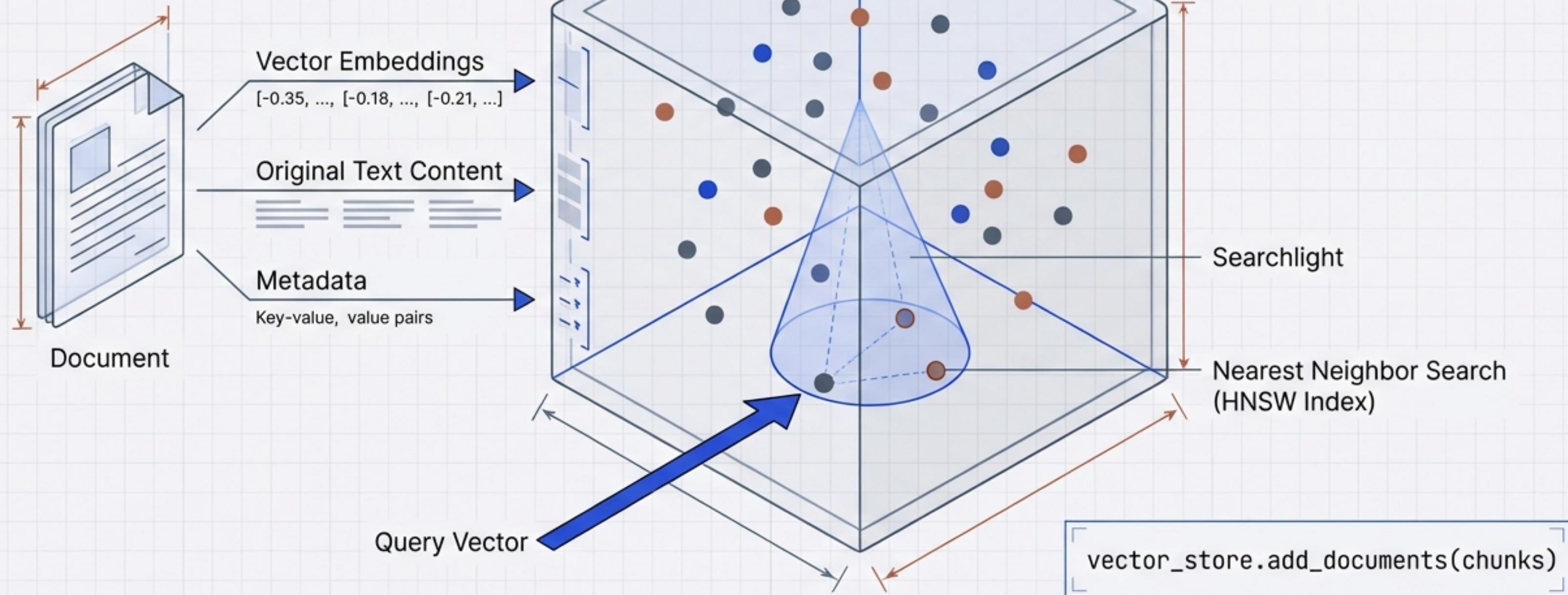
Input: 'Dogs allowed'

Process: Embedding Model
(all-MiniLM-L6-v2)

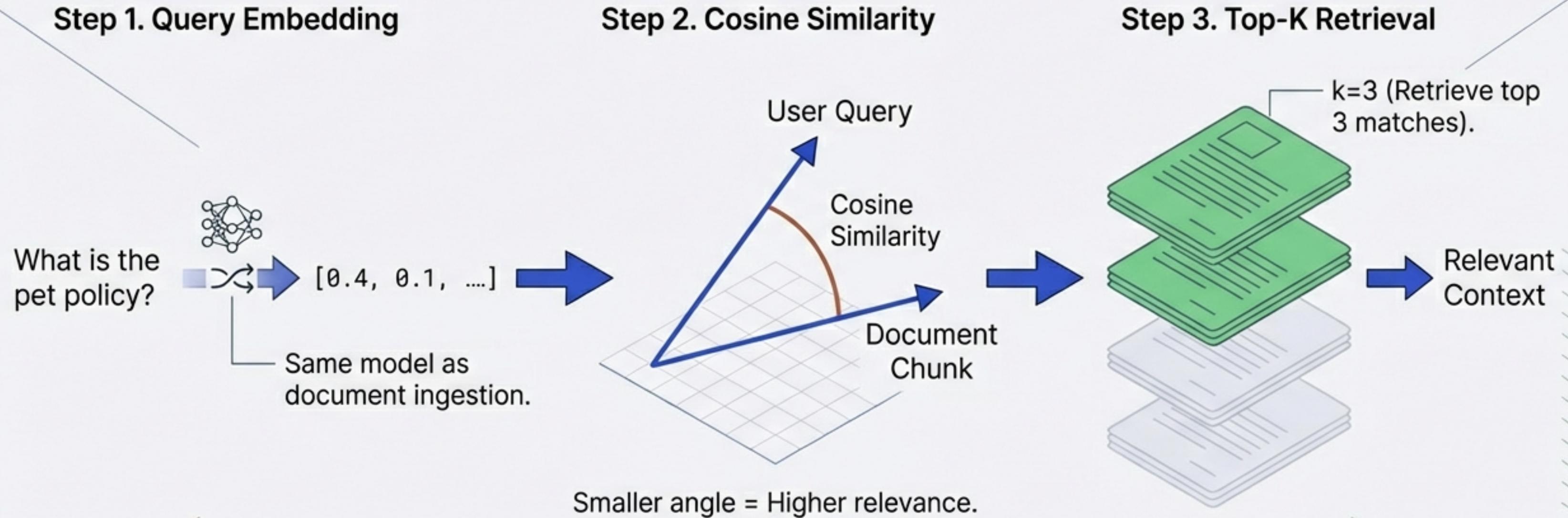
Output: Vector
[0.002,
 -0.41,
 0.98,
 0.55,
 -0.12,
 -0.12,
 ...]
}

The Vector Database

The Brain's Long-Term Memory



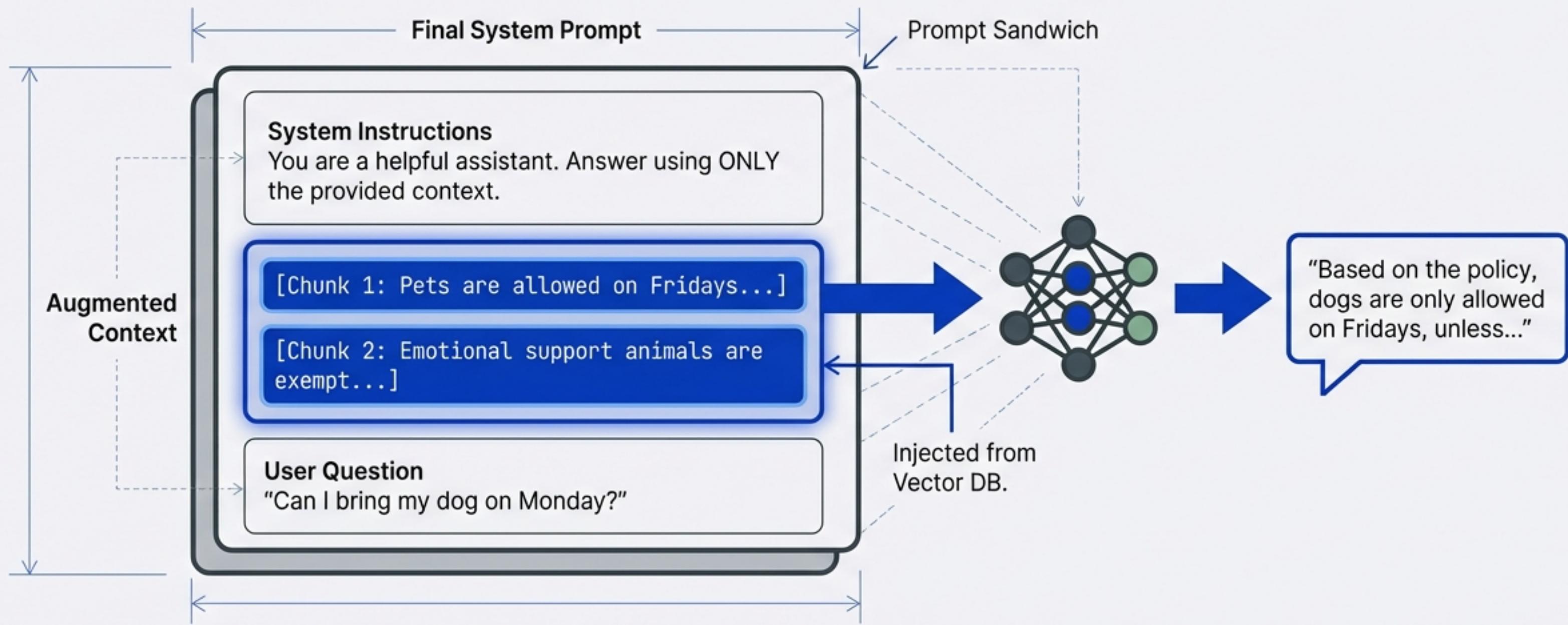
Pipeline B: Retrieval & Semantic Search





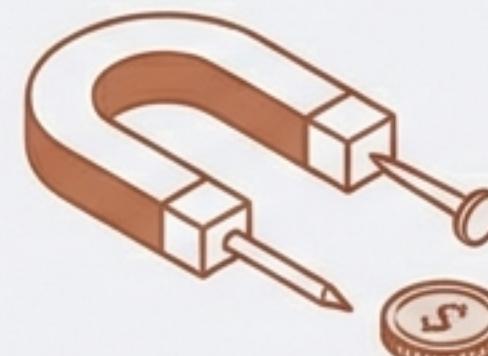
Augmentation & Generation

Synthesizing the Answer





The Reality Check: Why Naive RAG Fails



Retrieval Mismatch

Semantic search finds keywords but misses intent. Queries like "Tell me about it" fail completely.



Missing Context

Documents discussing related concepts (e.g., "Climate" vs "Polar Bears") might be stored too far apart.



Lost in the Middle

LLMs tend to ignore context buried in the middle of a long retrieved list.



Integration Noise

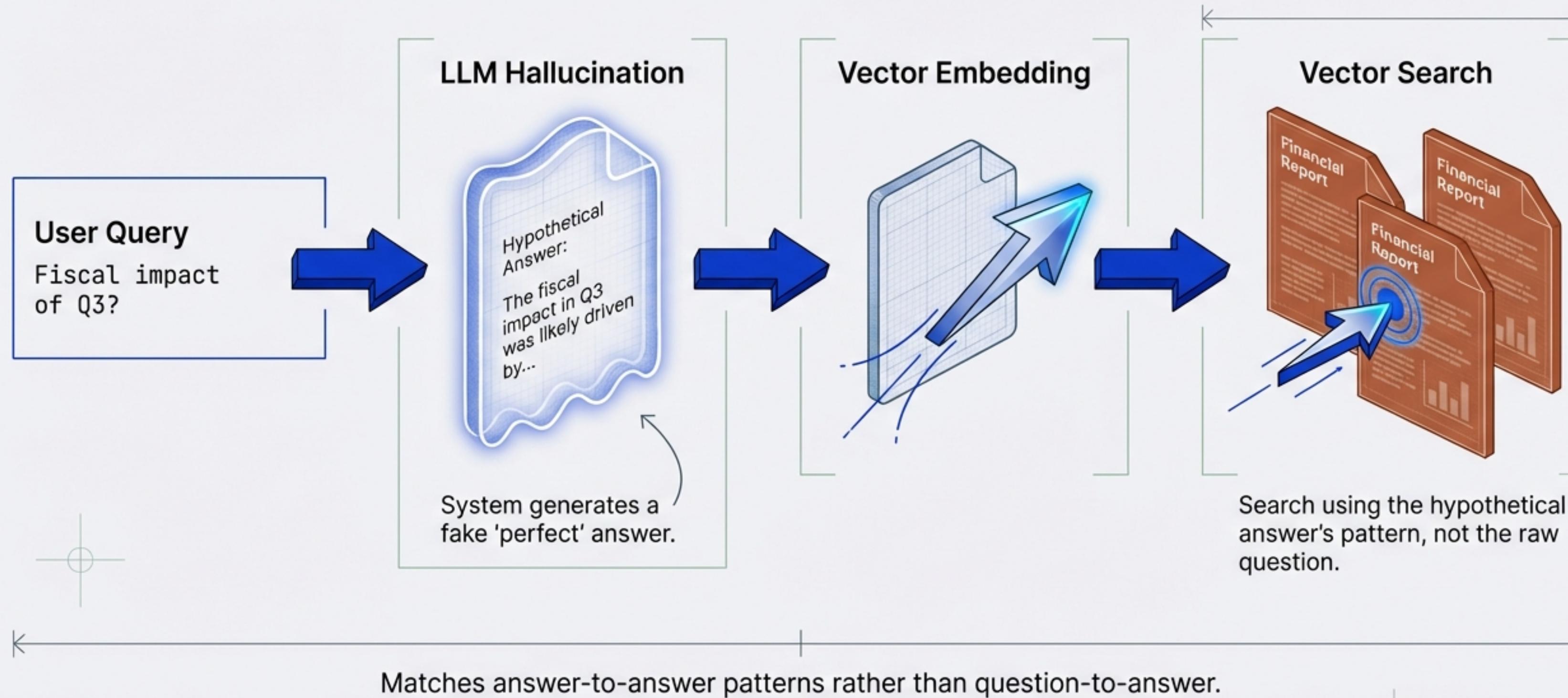
Retrieving irrelevant chunks confuses the LLM, leading to hallucinations.

Solution: Advanced RAG Techniques.

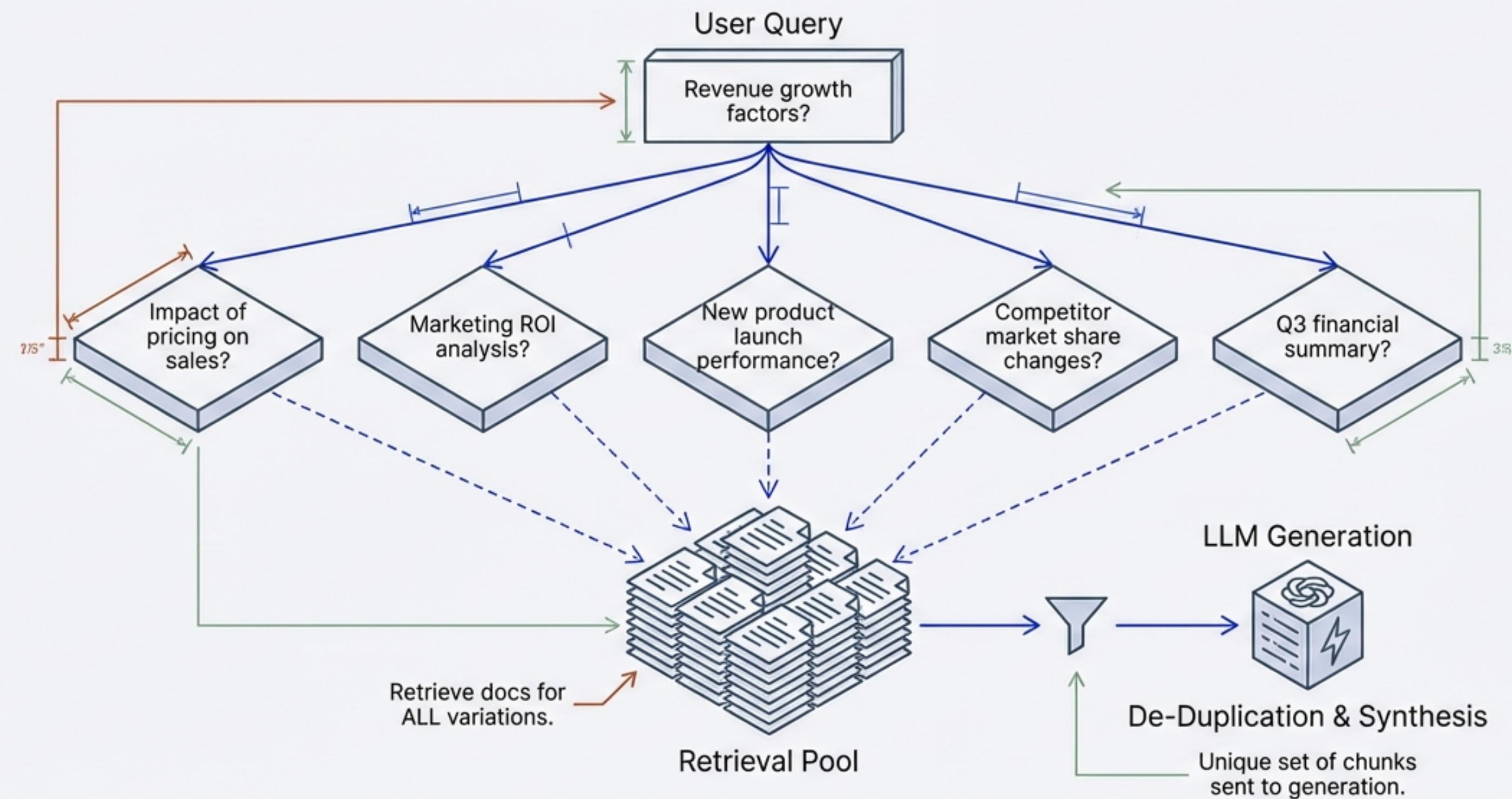


Advanced Technique: HyDE

Hypothetical Document Embeddings



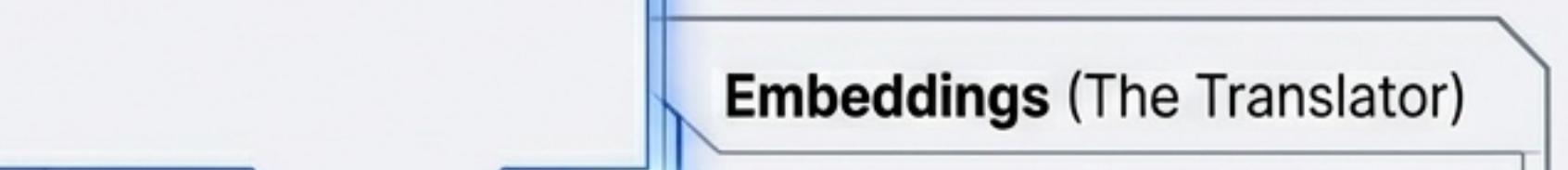
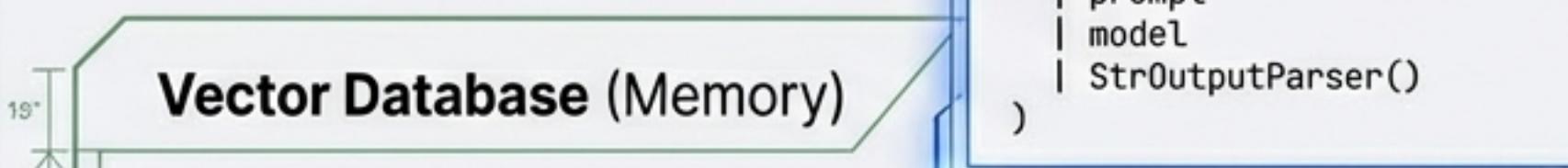
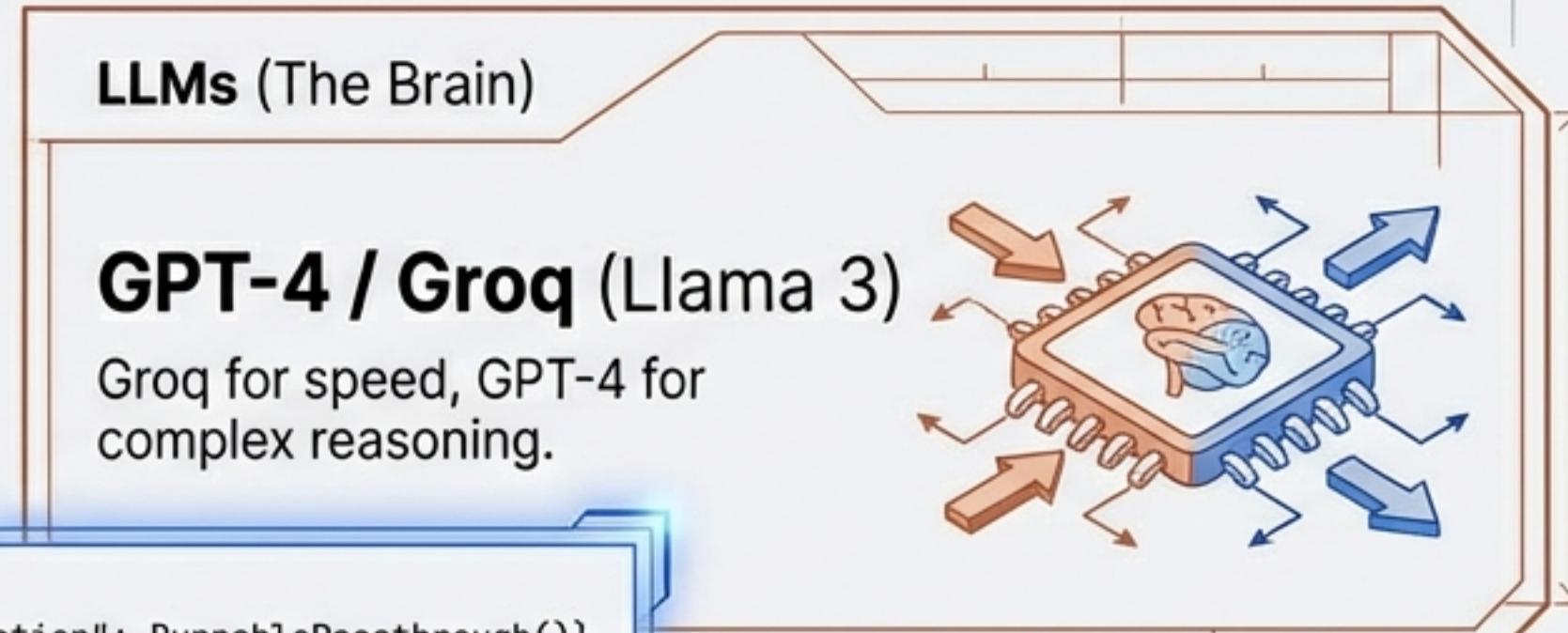
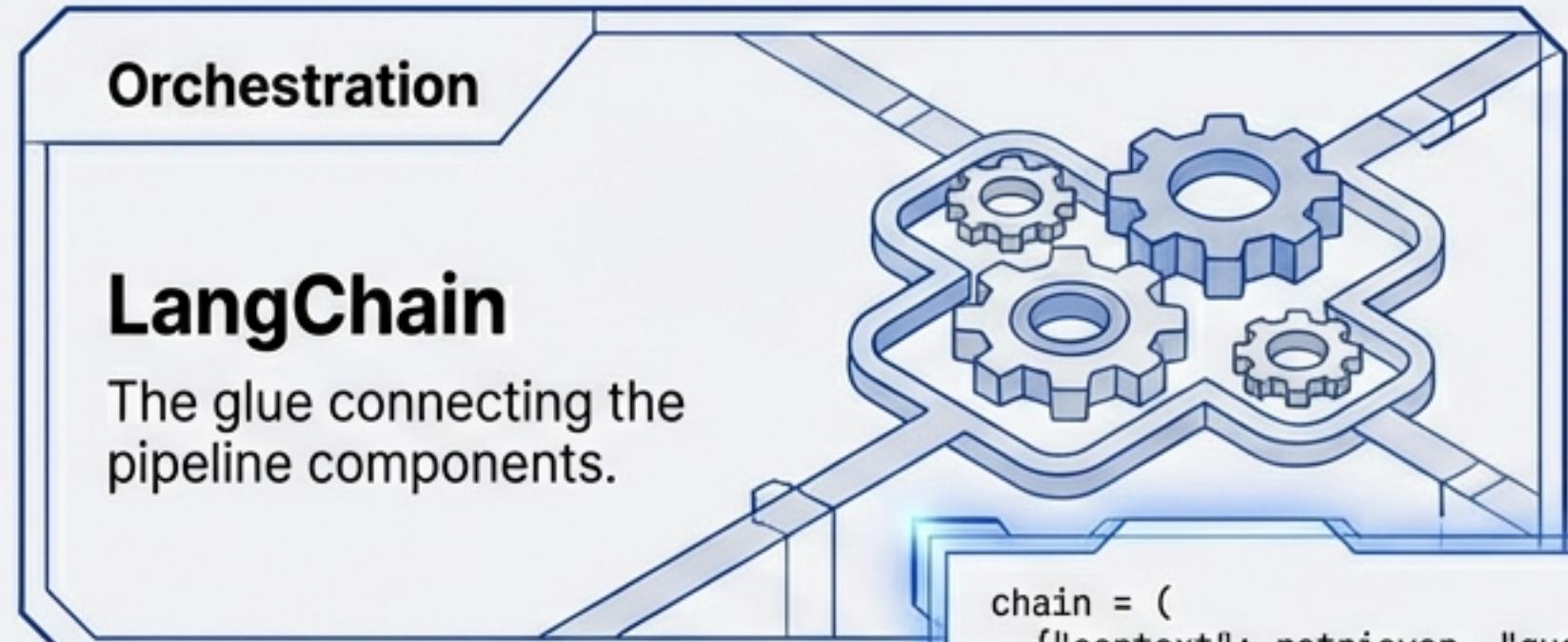
Advanced Technique: Multi-Query Expansion



Significantly increases recall by addressing multiple potential interpretations of the user's intent.

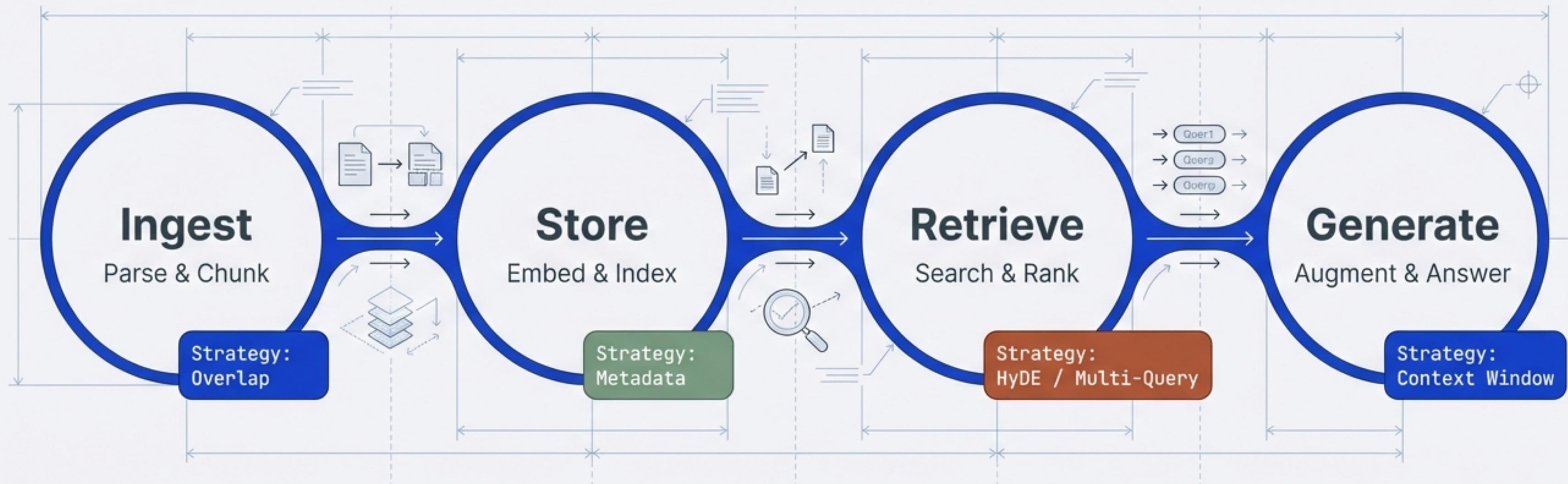


The RAG Technology Stack



A modular and scalable architecture for building advanced Retrieval-Augmented Generation systems.

From Naive Implementation to Advanced Intelligence



RAG is not just a search engine. It is a dynamic architecture that enables AI to reason over private data without the cost of fine-tuning.