# FROM POC TO PRODUCTION

Chaitanya Pathak

Chief Product and Technology Officer –Analyttica Datalabs

INSIGHTS FROM IMPLEMENTING ADVANCED RAG ARCHITECTURES
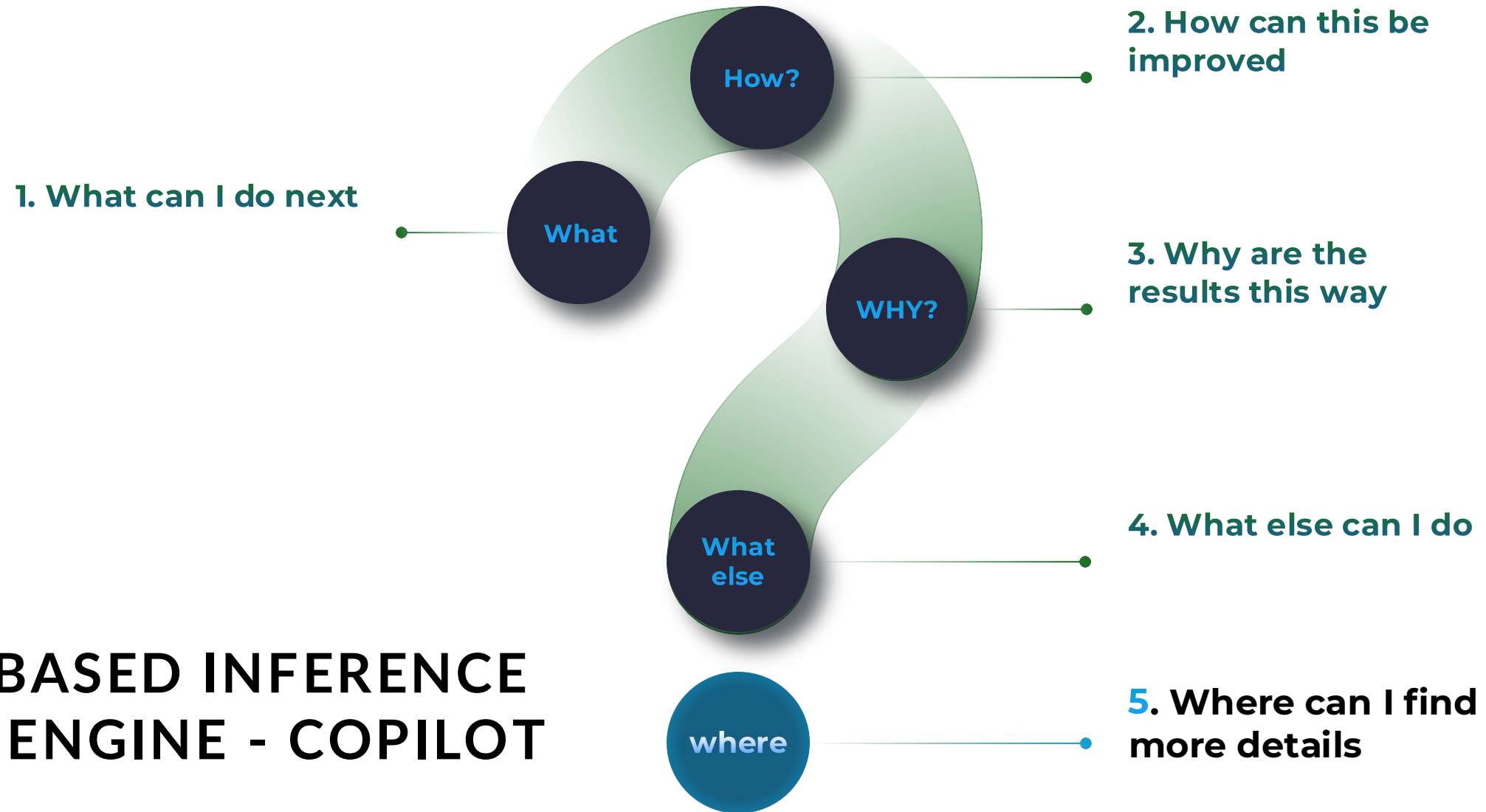
# TALKING POINTS

1. Context

2. Components

3. Prompts and RAG

4. Next Iteration

LEAPS PROGRAMS

2. How can this be improved

1. What can I do next

How?

What

WHY?

3. Why are the results this way

What else

4. What else can I do

RAG BASED INFERENCE ENGINE - COPILOT

where

5. Where can I find more details

# EVALUATION LANDSCAPE

LEAPS PROGRAMS

**Approaches**
o Reference based
o LLM based
o Hybrid

**Evaluation Scope**
o End to end
o Component

**Eval Data**
o Building Gold standard data
o Human vs synthetic

**Metrics Driven Development**

**Generative**
o Faithfulness
o Answer Relevancy

**IR Metrics**
o Rank based (Hit Rate. MRR, NDCG
o Predictive ( F1 score, precision , recall)
o Context Precision , Context Recall

**Design Considerations**
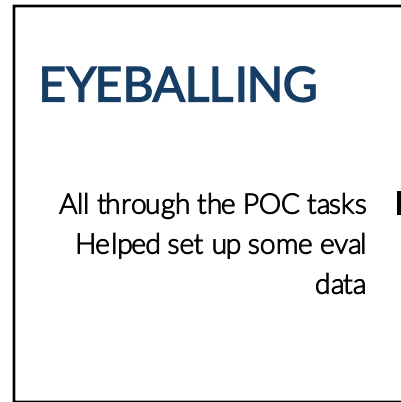o Deterministic vs nondeterministic
o Turns
o Prompt Flows

## EYEBALLING

All through the POC tasks
Helped set up some eval data

**1**

**2**

## SUPERVISED

Labelled Eval data for ground truth
Domain drives the complexity

## LLM AS JUDGE

Instrumentation of Evaluation
Required Sophisticated prompt engineering

**3**

# EVALUATION JOURNEY

META DATA

PARSING

SEMANTIC CHUNKING

SPARSE & DENSE RETRIEVERS

DATA PREPROCESSING

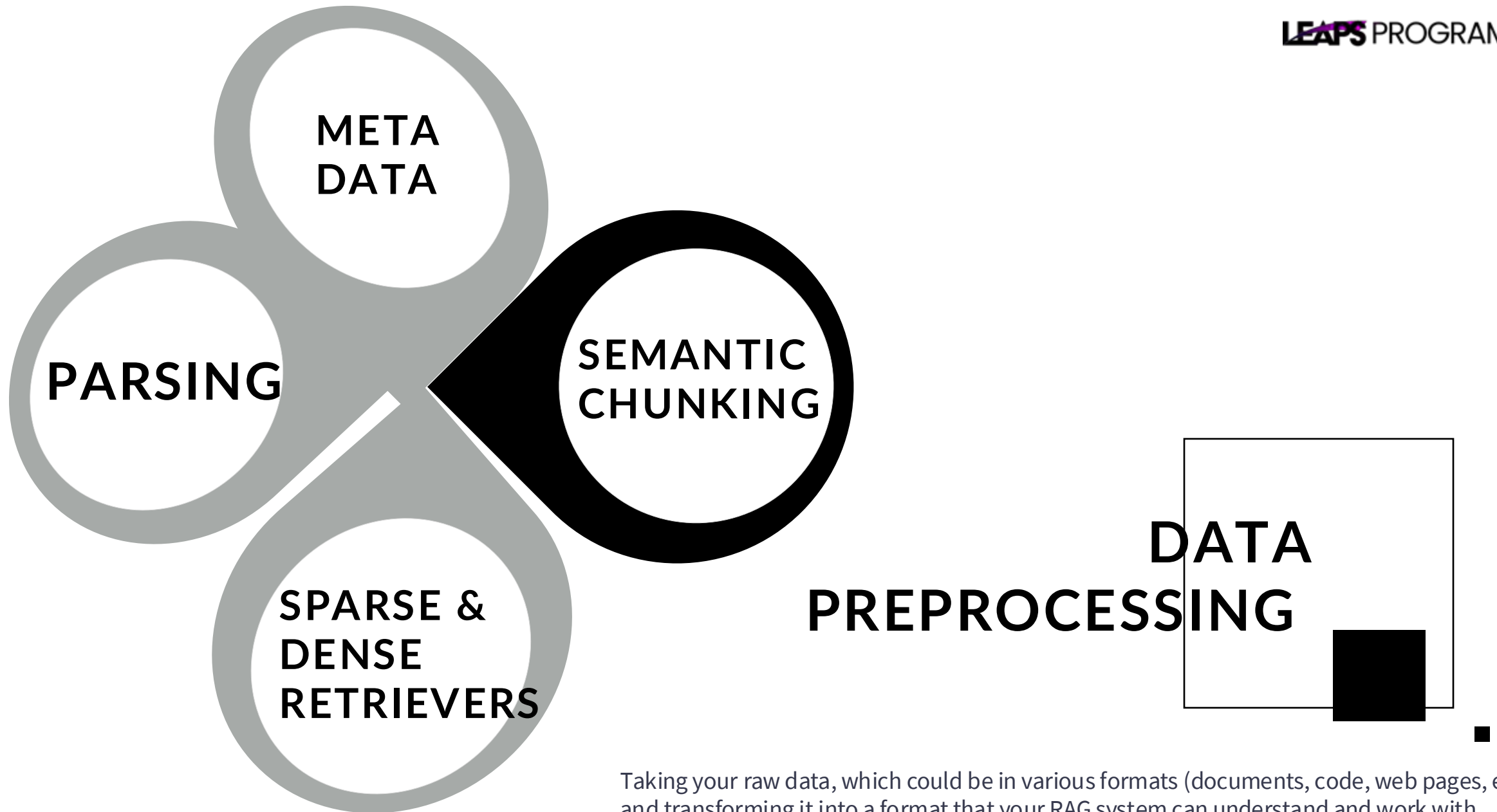Taking your raw data, which could be in various formats (documents, code, web pages, etc.), and transforming it into a format that your RAG system can understand and work with efficiently.

# QUERY ENHANCEMENT – Benefits, Best practices and Tradeoffs

**01**

Intent

Helps s understand intent better

**02**

Breakdown

Breakdown complex queries into sub queries which are more consumable by the search , retrieval and LLM models

**03**

Tradeoffs

Latency and complexity

**04**

HyDE

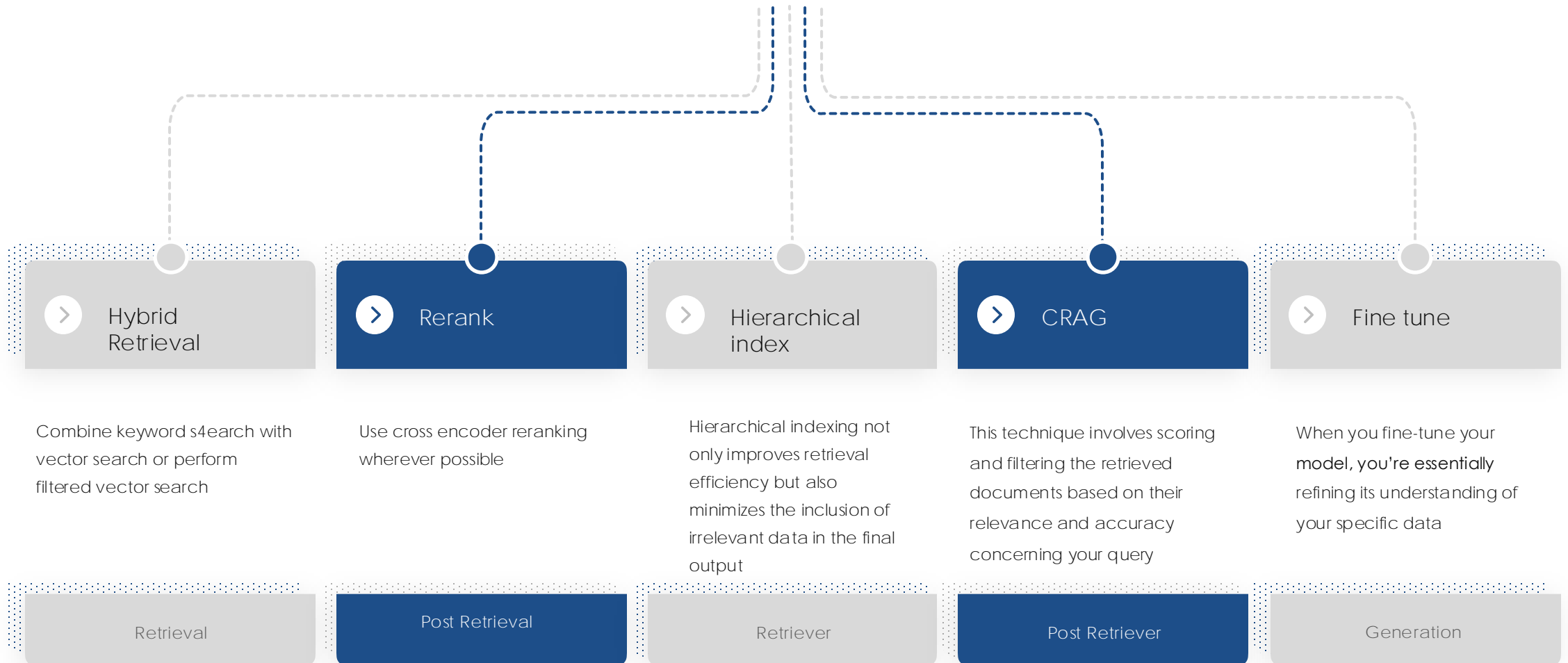Generate documents / context wherever not available

**05**

Routing

This is essential best practice in a multi database / context scenarios
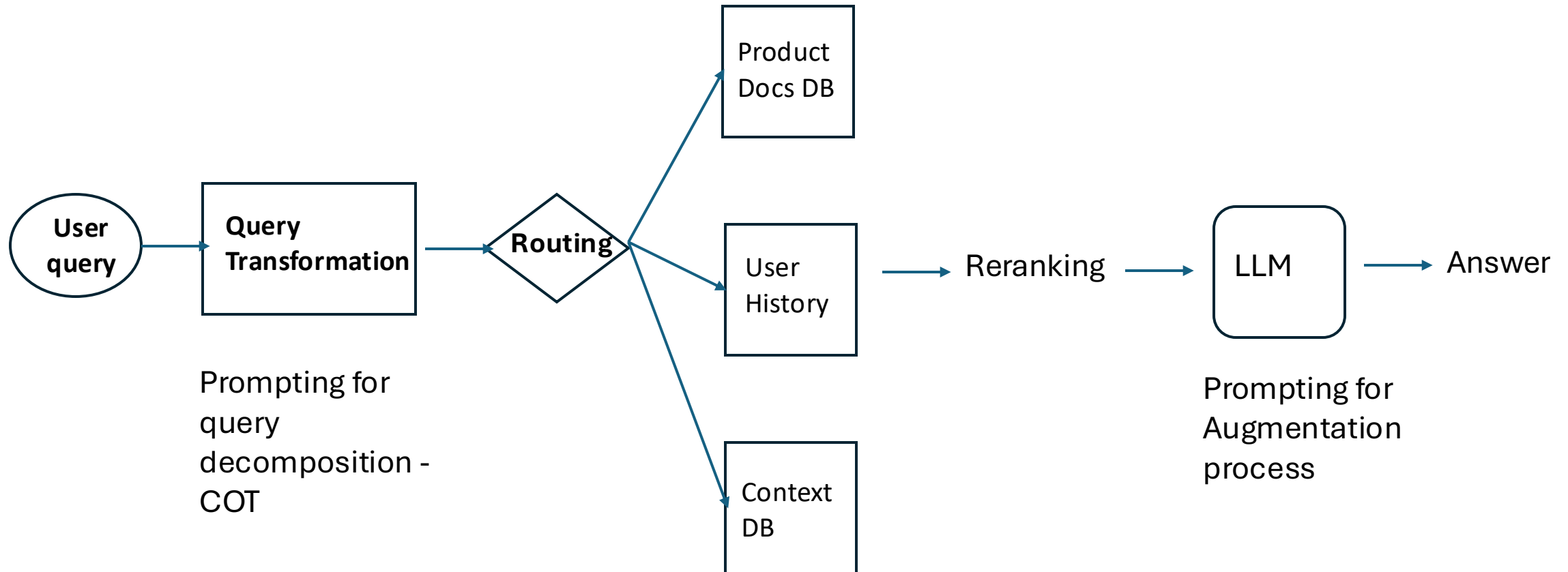
An essential part of Rag systems that improves retrieval accuracy and relevancy of responses

# Retriever and Reranking

## Hybrid Retrieval

Combine keyword s4earch with vector search or perform filtered vector search

Retrieval

## Rerank

Use cross encoder reranking wherever possible

Post Retrieval

## Hierarchical index

Hierarchical indexing not only improves retrieval efficiency but also minimizes the inclusion of irrelevant data in the final output

Retriever

## CRAG

This technique involves scoring and filtering the retrieved documents based on their relevance and accuracy concerning your query

Post Retriever

## Fine tune

When you fine-tune your model, you're essentially refining its understanding of your specific data

Generation

# Prompt engineering and RAG architecture

User query

Query Transformation

Routing

Product Docs DB

User History

Context DB

Reranking

LLM

Answer

Prompting for query decomposition - COT

Prompting for Augmentation process

**Take away** - careful prompt engineering will be employed for component(s) of RAG pipeline that uses LLM

# NEXT STEPS

- Move towards modular RAG architecture

- Full versioning of prompts

- Unit testing – Promptfoo

- Agentic abilities – use of tools

HAPPY TO TAKE QUESTIONS
CHAITANYA.PATHAK@ANALYTTICA.COM