# Simplifying RAG with Databricks Vector Search

Wilson Mok

Mar 14, 2025

Dear Azure

# Bio



## Wilson Mok

Sr. Data architect & Consultant

/wilson-mok

/wilson-mok

@the-analytics-lab
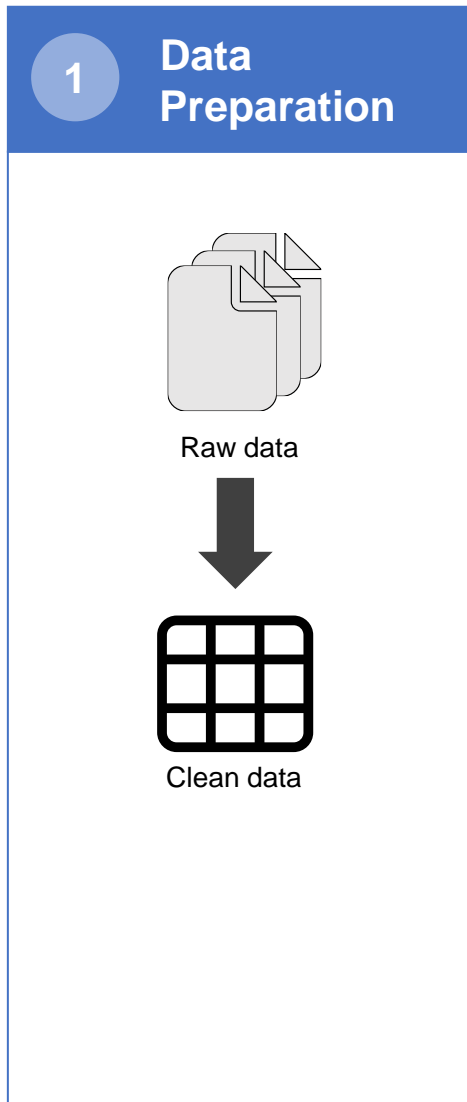
My experiences includes:
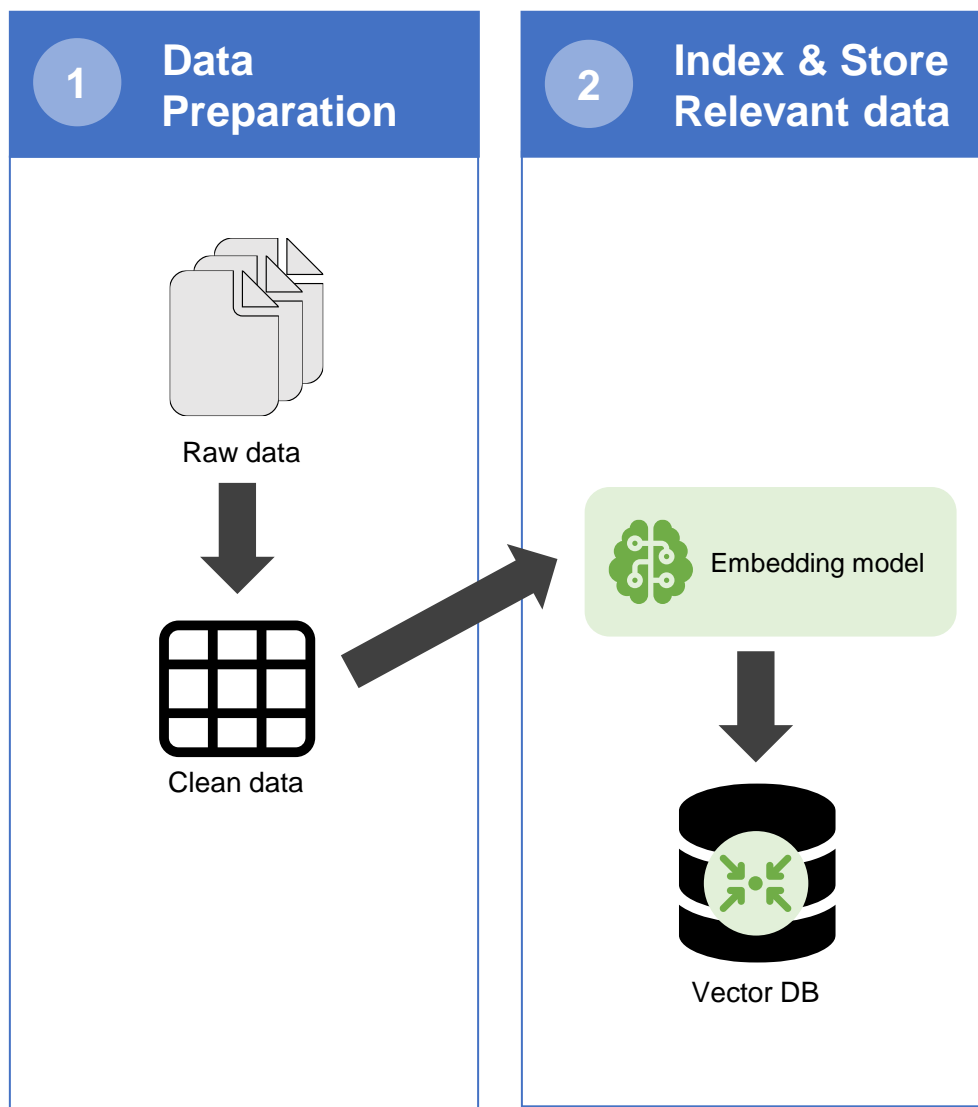- Avanade
- CAE
- Air Canada

# Expected Audience and Presentation Structure

- This session, we will cover:
  - Discuss the complexity of RAG pipelines.
  - Simplifying the RAG pipelines using Databricks Vector Search.
  - Using Unity Catalog's delta table and volume to store and manage our data.
  - Quick Introduction on Databricks Genie.
  - Demo: Implementing RAG using Databricks Unity Catalog and Vector Search.

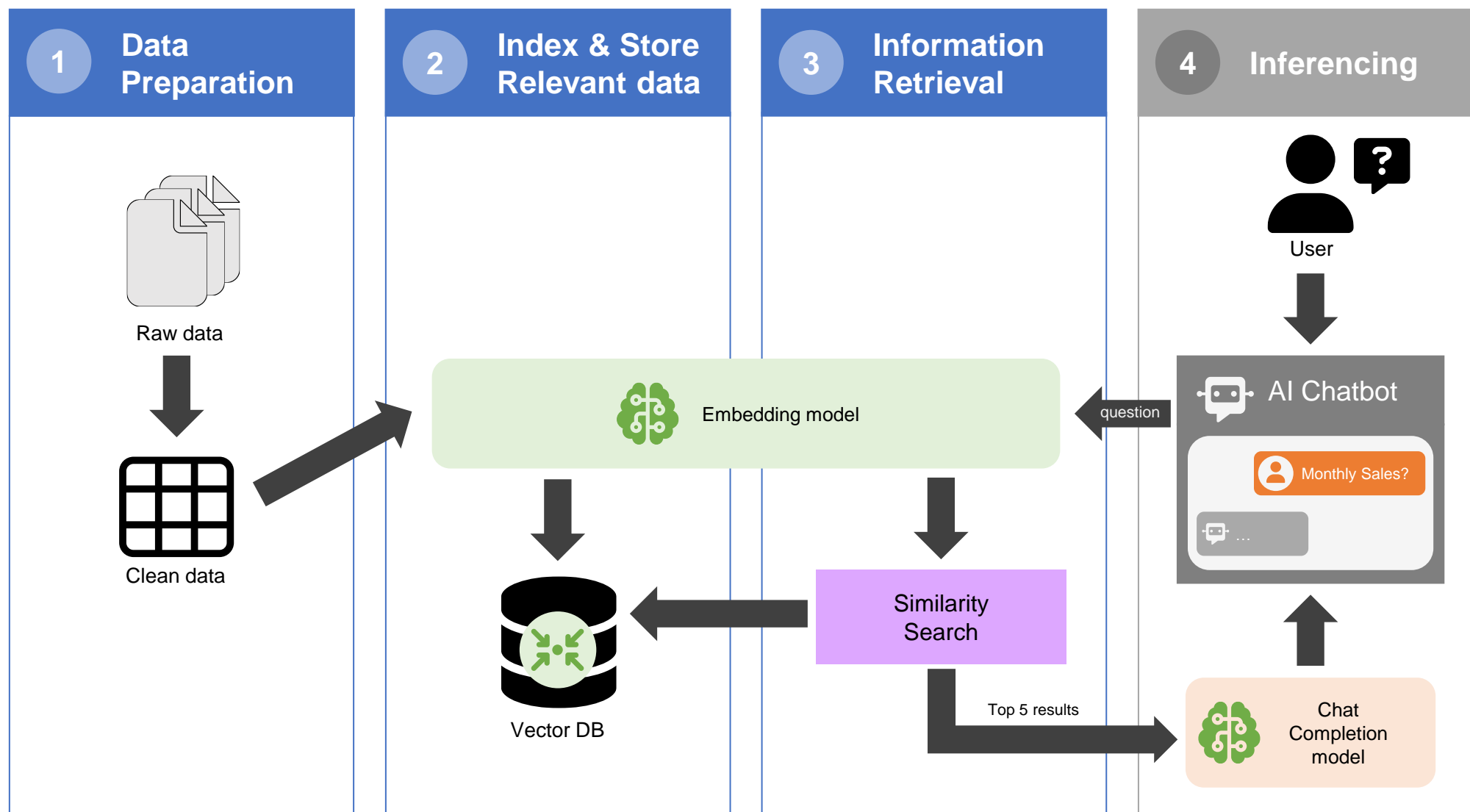- Fundamental knowledge: Gen AI, Databricks and Python would be beneficial.

# There are four key steps in a RAG workflow to create an AI Chatbot.

| 1 | Data Preparation |
|---|---|

Raw data

Clean data

# There are four key steps in a RAG workflow to create an AI Chatbot.



**1** **Data Preparation**

Raw data

Clean data

**2** **Index & Store Relevant data**

Embedding model

Vector DB

# There are four key steps in a RAG workflow to create an AI Chatbot.



**1** **Data Preparation**

Raw data

Clean data

**2** **Index & Store Relevant data**

Embedding model

Vector DB

**3** **Information Retrieval**

Similarity Search

Top 5 results

**4** **Inferencing**

User

question

AI Chatbot

Monthly Sales?
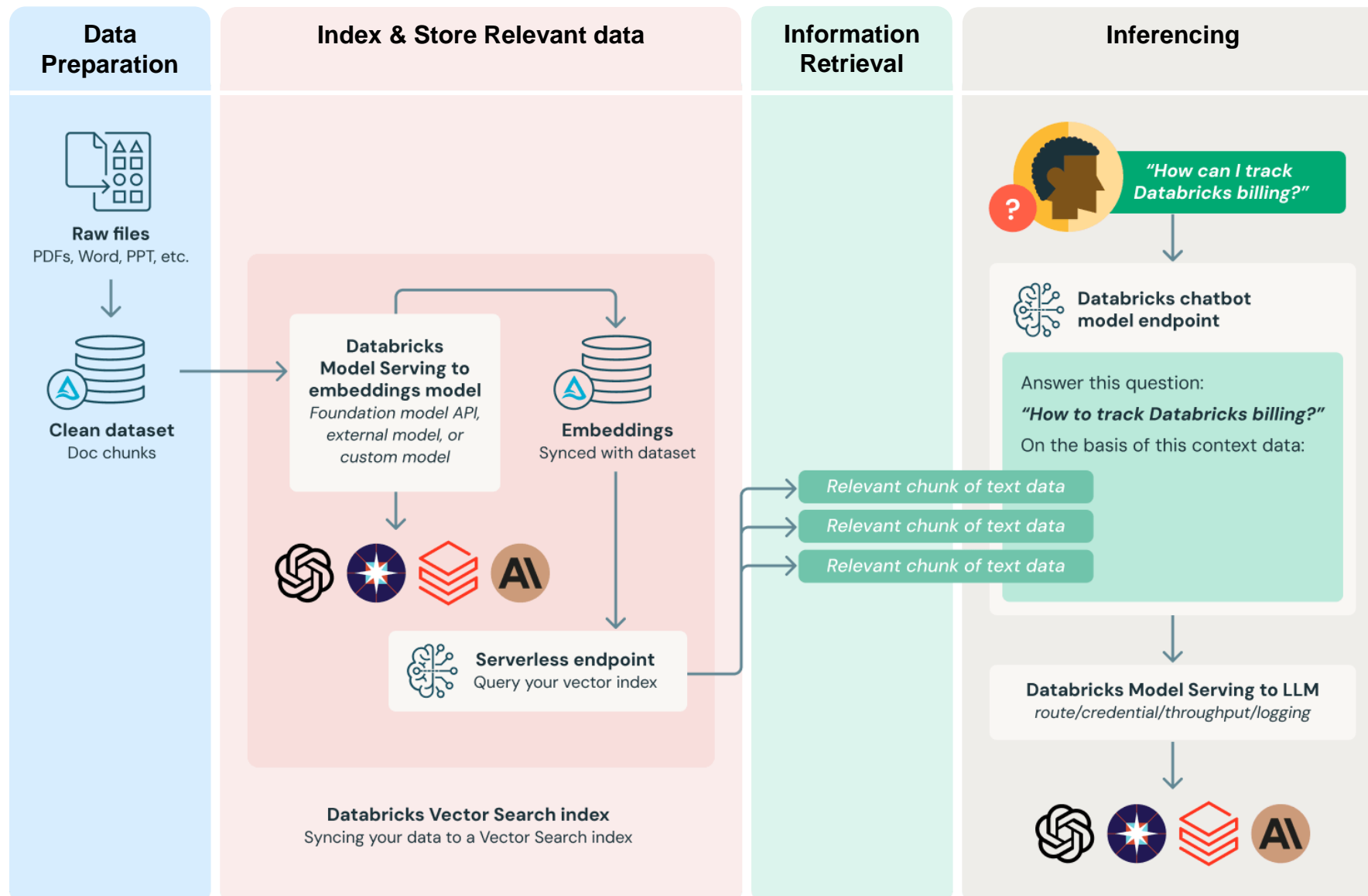
...

Chat Completion model

\* This presentation will not be covering Chat Completion in detail

6

# Each step in the RAG workflow requires a set of specialized tools and technologies…

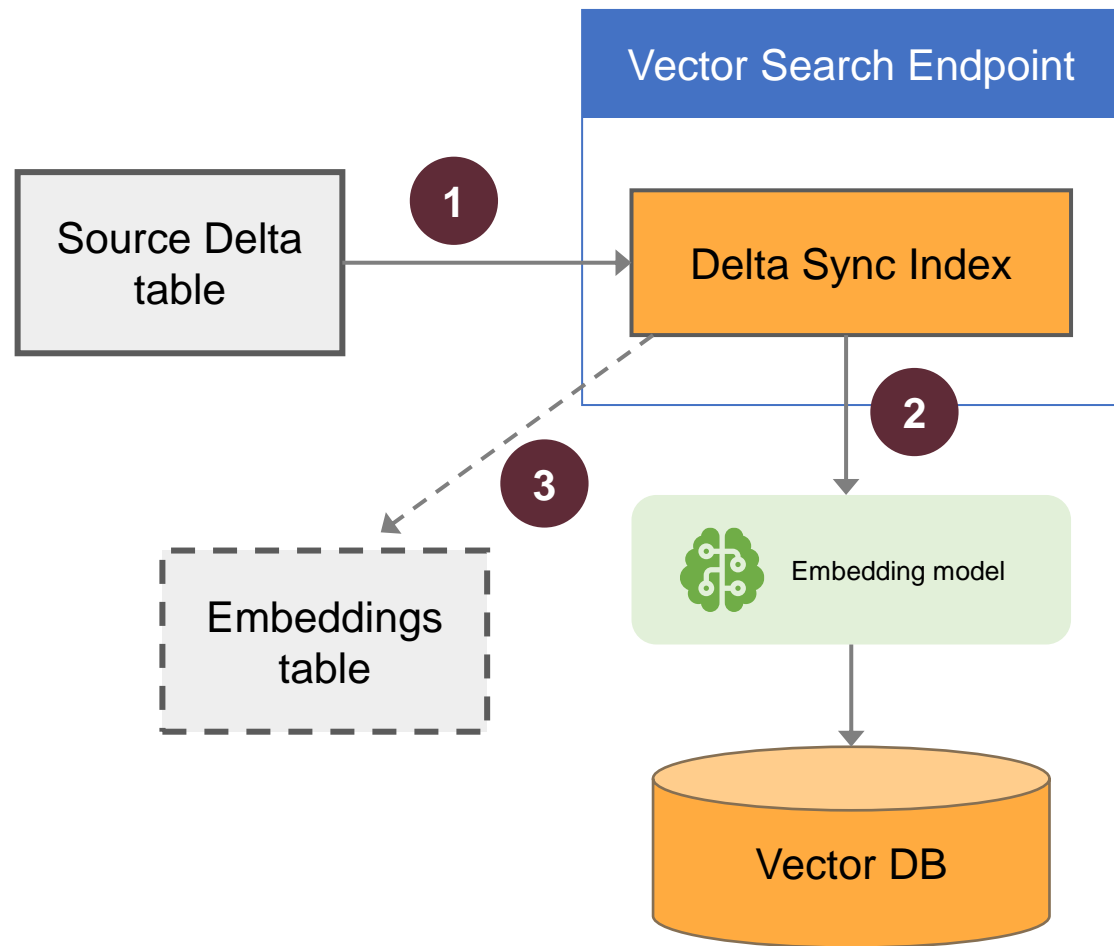| | |
|---|---|
| **Data processing** | • Azure Cognitive Search<br>• Azure AI Vision<br>• Azure AI Speech-to-Text<br>• Spark (Databricks, MS Fabric)<br>• Orchestrator (Azure Data Factory, LangChain) |
| **Vector Database** | • Pinecone<br>• Azure AI Search<br>• Postgres DB (pgvector) |
| **Embedding Model** | • Open AI: text-embedding-3<br>• Google: text-embedding-gecko<br>• Open Source: GTE |
| **LLM Model** | • Open AI: GPT-4o<br>• Google: Gemini<br>• Anthropic: Claude<br>• Open Source: Meta Llama 3.3, DeepSeek |

# Databricks streamlines this workflow by…

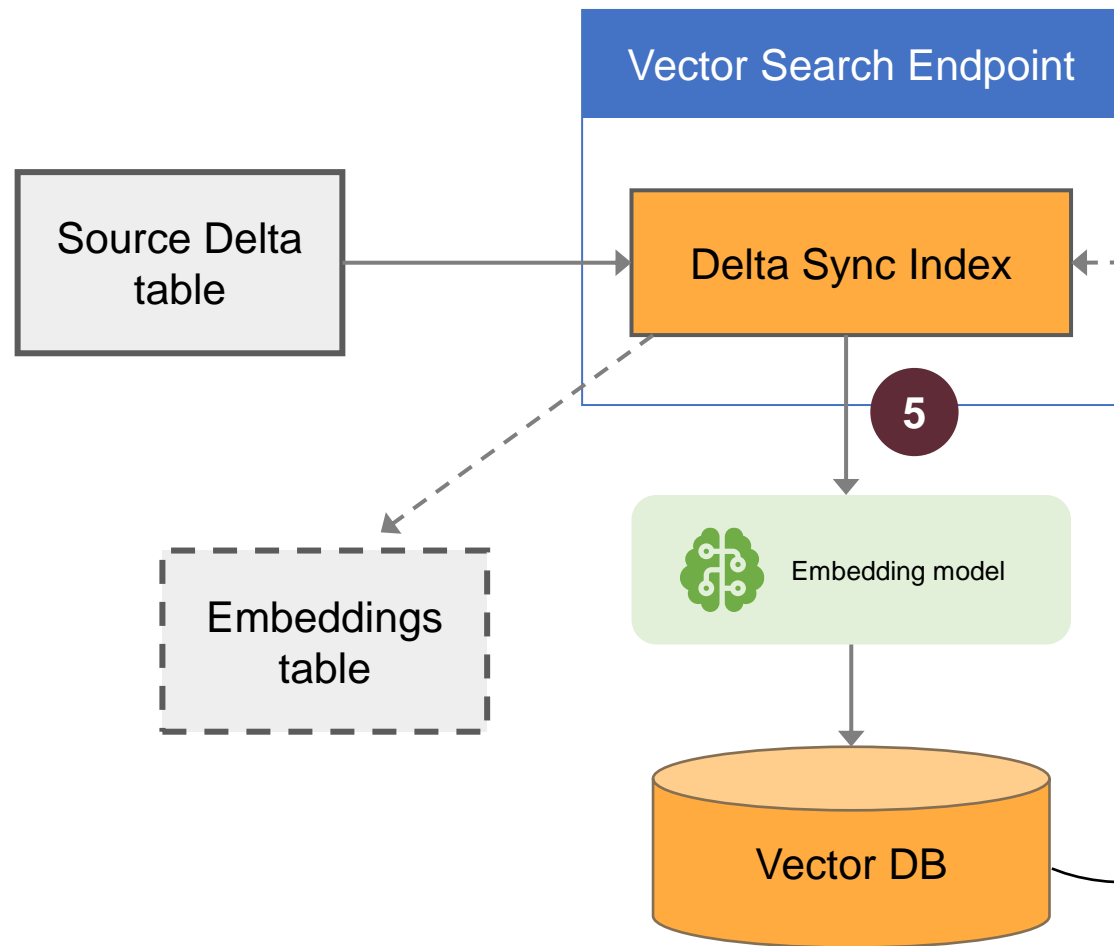Credit: https://www.databricks.com/glossary/retrieval-augmented-generation-rag

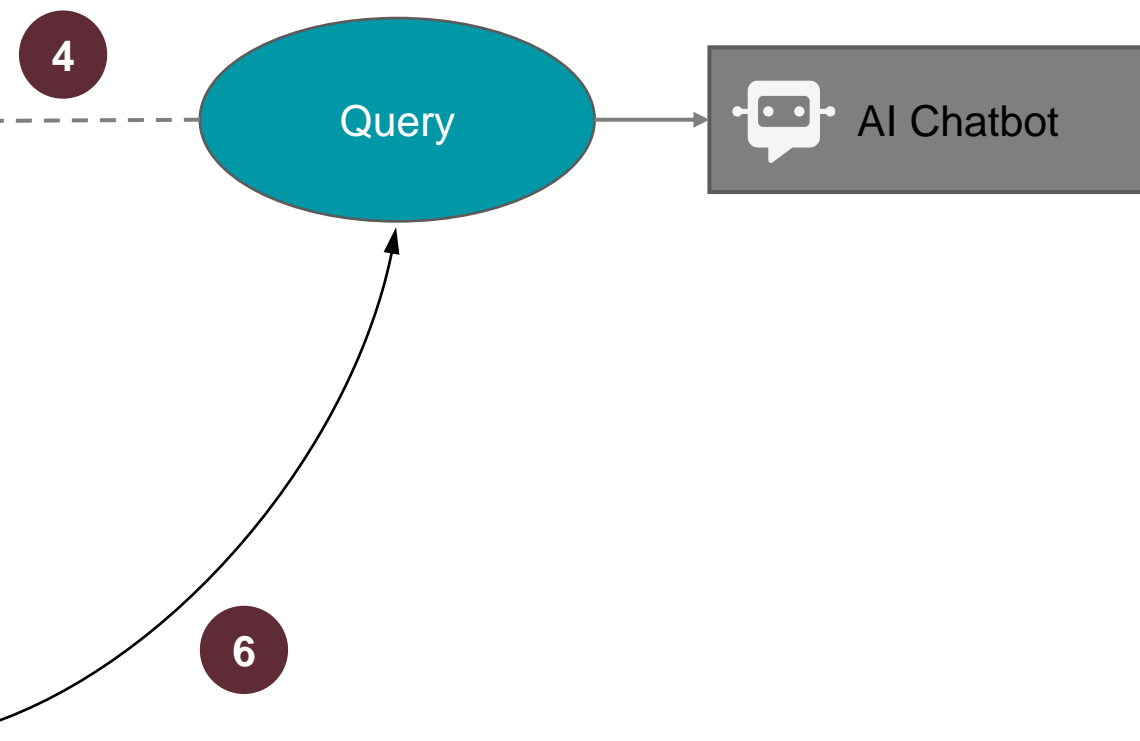# A Closer look into Databricks – Mosaic AI Vector Search…

Index & Store data
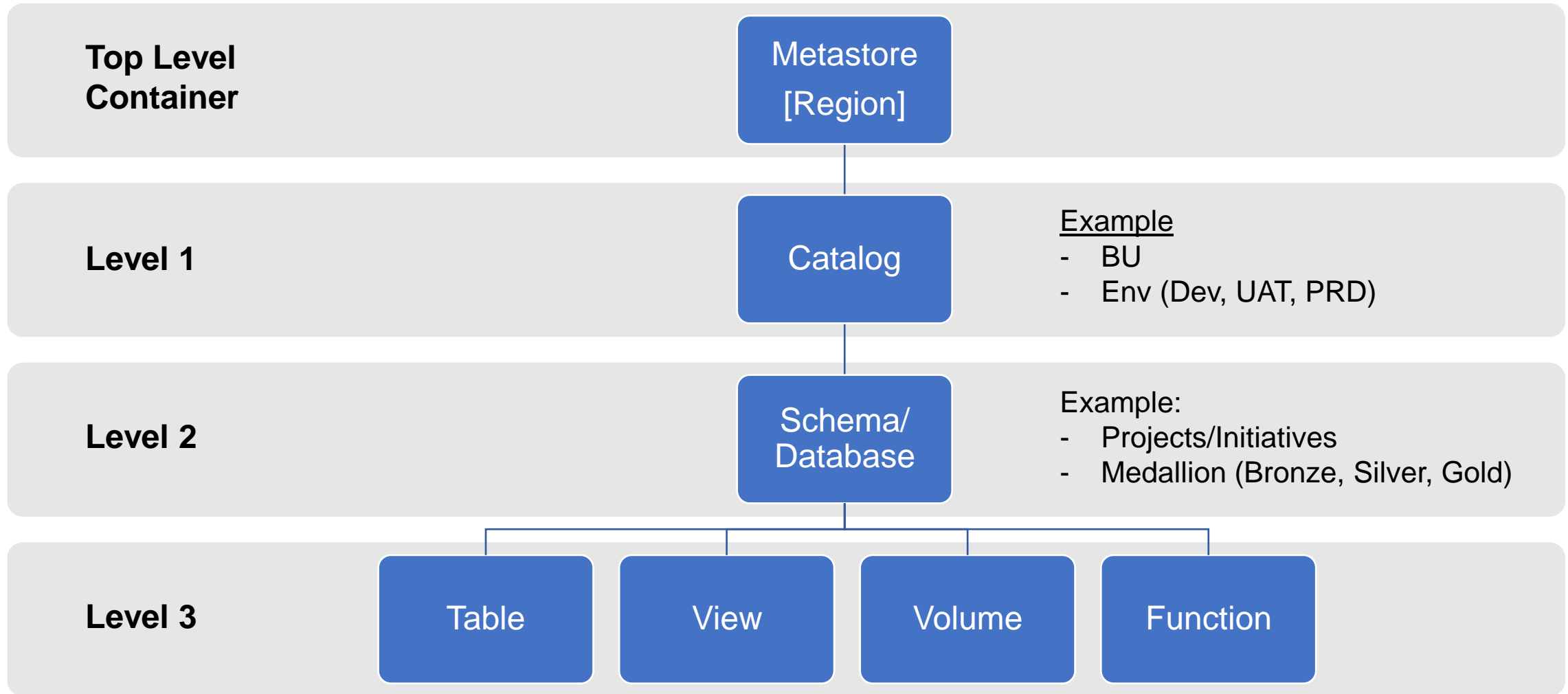
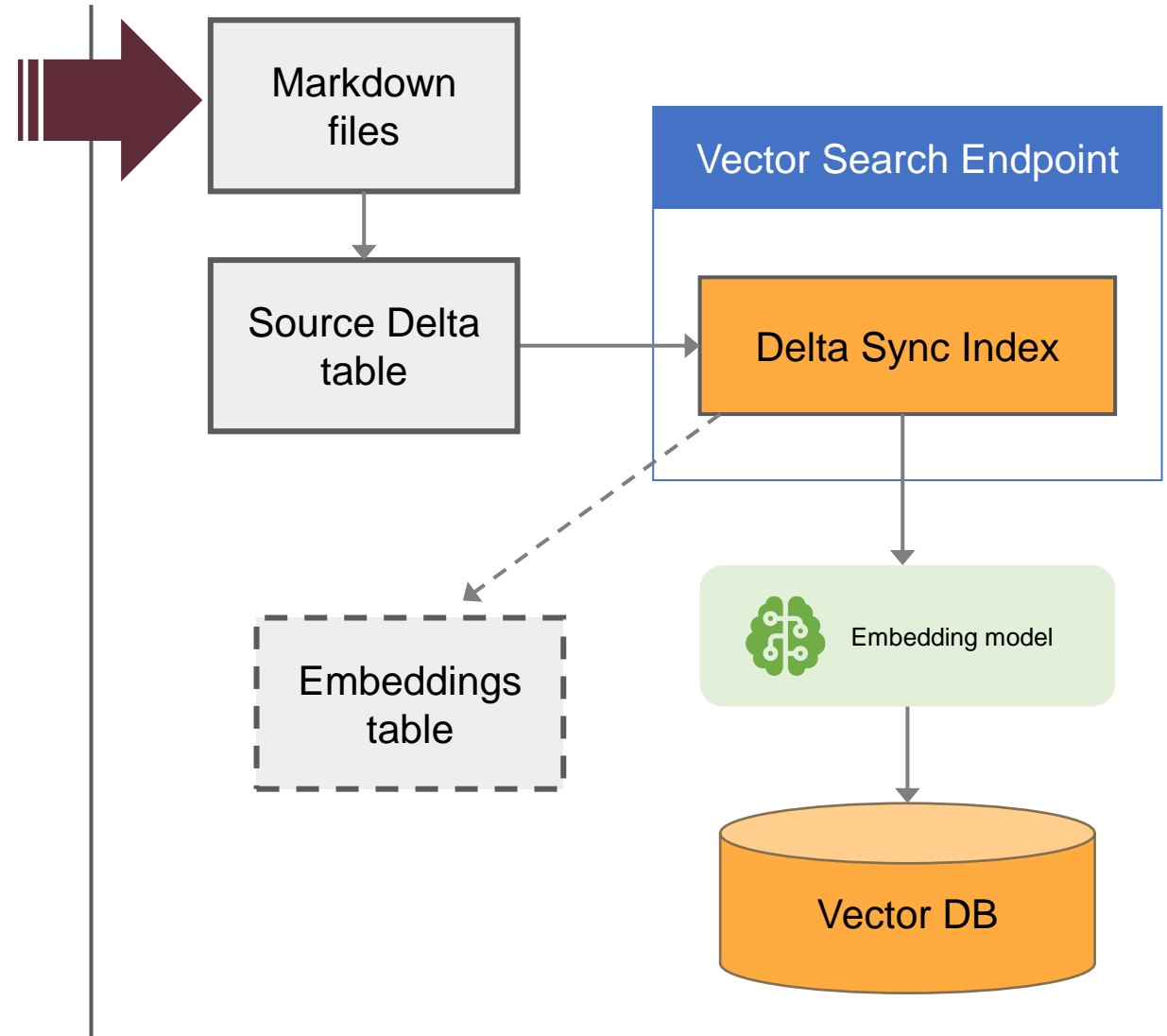# A Closer look into Databricks – Mosaic AI Vector Search…

Index & Store data

Information Retrieval

Vector Search Endpoint

Source Delta table

Delta Sync Index

Embeddings table

Embedding model

Vector DB

Query

AI Chatbot

4

5

6

# Unity Catalog provides a way to govern and manage…

| | | |
|---|---|---|
| **Top Level Container** | Metastore [Region] | |
| **Level 1** | Catalog | Example - BU - Env (Dev, UAT, PRD) |
| **Level 2** | Schema/ Database | Example: - Projects/Initiatives - Medallion (Bronze, Silver, Gold) |
| **Level 3** | Table   View   Volume   Function | |

# Demo 1: Implementing RAG using Databricks Unity Catalog and Vector Search



https://schiiss.github.io/blog/

# Q & A 1

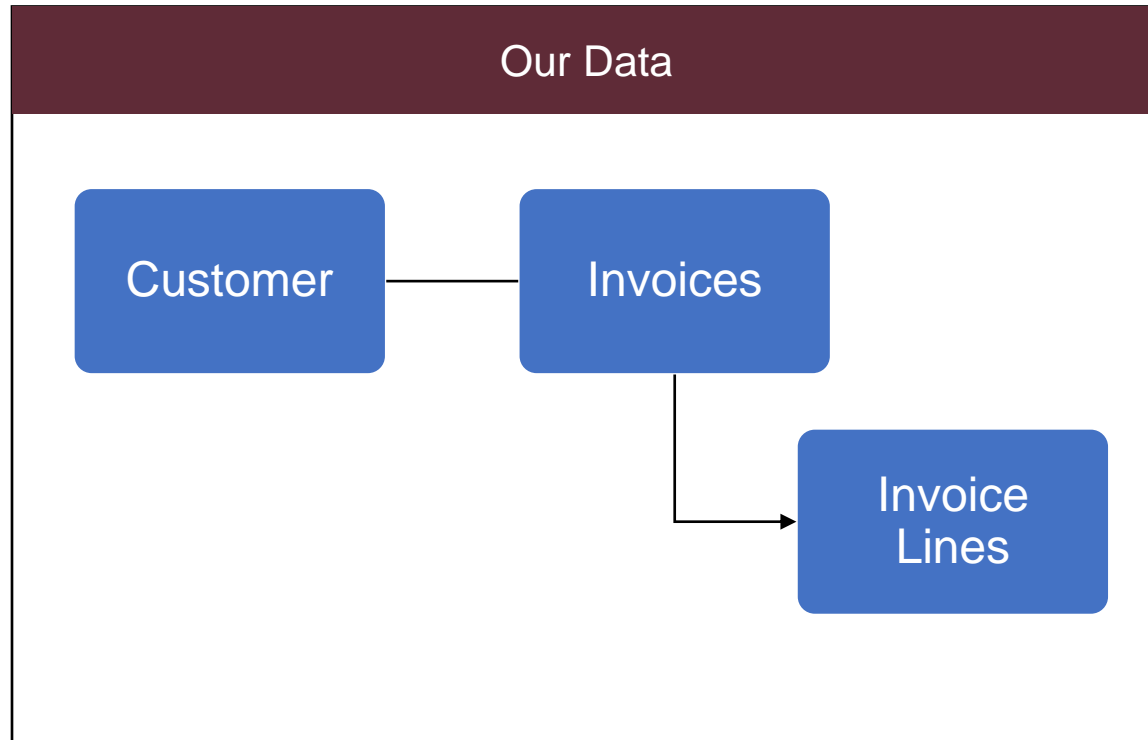Credit: https://www.databricks.com/blog/onboarding-your-new-aibi-genie

## What is an AI/BI Genie?

- Conversational AI assistant powered by Open AI.
- Talk to your data in Natural Language.
- Genie translates them into SQL and retrieve relevant insights.
- Can generate Visualizations like Charts and Graphs.

# Demo: Databricks' AI/BI Genie



## Our Data

Customer → Invoices

Invoices → Invoice Lines

Credit: https://www.databricks.com/blog/onboarding-your-new-aibi-genie

# Q & A 2

# After all this, we have run through…

1. Mosaic AI Vector Search: Starts with a Delta Table, creates a Delta Sync Index to create embedding via Vector Search Endpoint.
2. Unity Catalog: Provides a governed, structured approach to manage your data using three-level hierarchy. We used Delta tables and Volumes in action today.
3. AI/BI Genie: An AI-powered assistant that enables natural language querying, automated insights, and seamless BI integration.

- The code is in my Github repo.
  - https://github.com/wilson-mok/demo/tree/main/2025/DearAzure/1-Databricks-Vector-Search

# Thank you

LinkedIn