

Talk + Demo

# Db2 for LLM Apps

Vector Search, LLM Frameworks, &  
LLM Inferencing via SQL



Shaikh Quader  
Db2 AI Architect



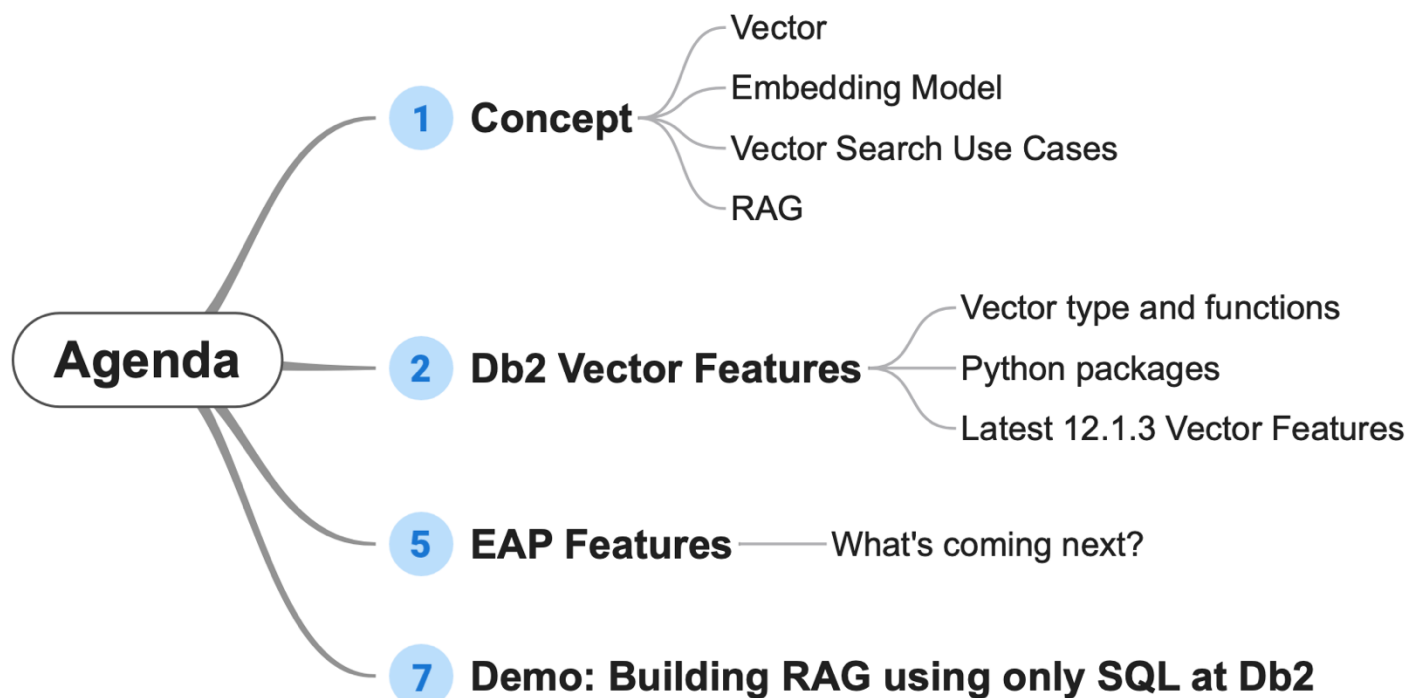
Marhaba Chariwala  
Db2 Developer

Thursday, Nov 13, 2025 @ 1:00 – 2:00 p.m

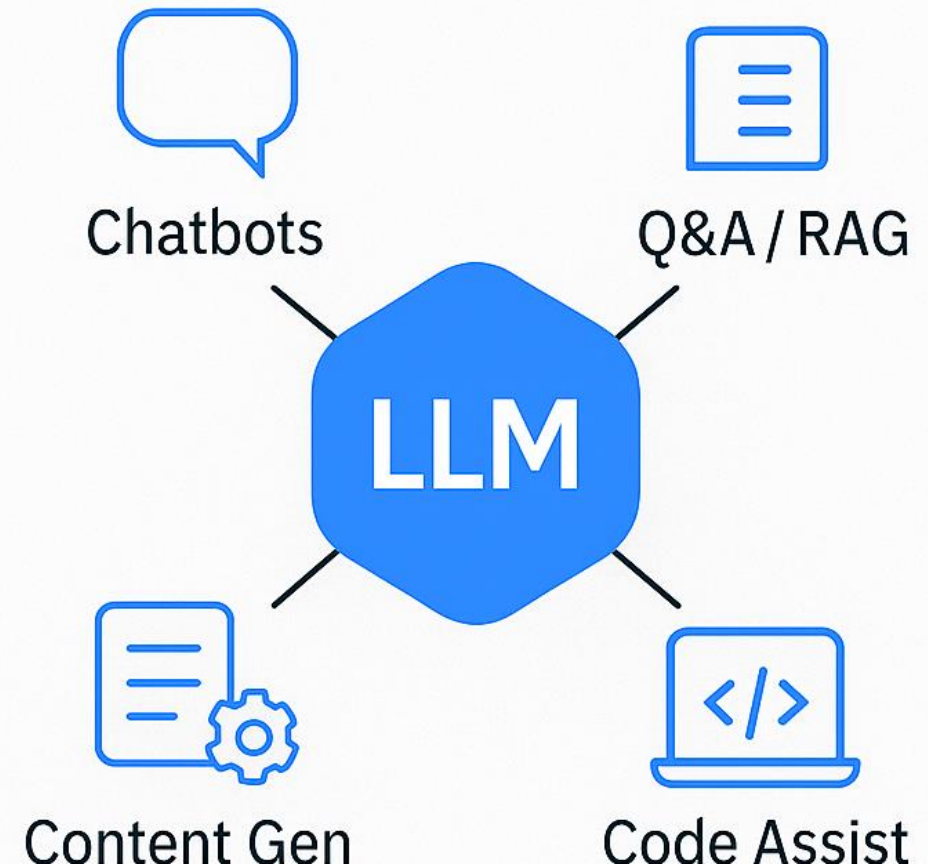
Venue: 8200 Warden Amphitheatre

Online: [ibm.biz/db2talk](https://ibm.biz/db2talk)

## Db2 for LLM Apps Talk @ Nov 13: Building Language Apps with Db2



Apps that use **LLMs** to understand, generate, and interact with natural language.





- A list of numbers, like (1, 2)
- Represents a point in space
- Like map coordinates for cities
- You can measure distance between two vectors



# 3 Common Retrieval Techniques



## Keyword-based

Finds exact word matches such as IDs, emails

E.g., search query: "laptop repair"

Finds: "laptop repair guide", "repair@email" (no ranking)

Misses: "fixing laptops" (fixing ≠ repair), "repairing computers"



## Full-text

Matches word variations, ranks results

e.g., search query: "laptop repair"

Ranked: "Fixing Laptops" > "Computer Tips" > "Support email"

Misses: "Device troubleshooting" (no word overlap)



## Vector-based

Find similar ideas, not just word matches

e.g., search query: "laptop repair"

Finds: "Fixing Laptops", "Device Troubleshooting"

✓ Matches meaning - works with any vocabulary

# Which Two Are More Similar? (What is Similarity, BTW?)





*Transforming real-world objects into numerical representations*



Coffee



Laptop



Cat



**EMBEDDING  
MODEL**



**COFFEE VECTOR**

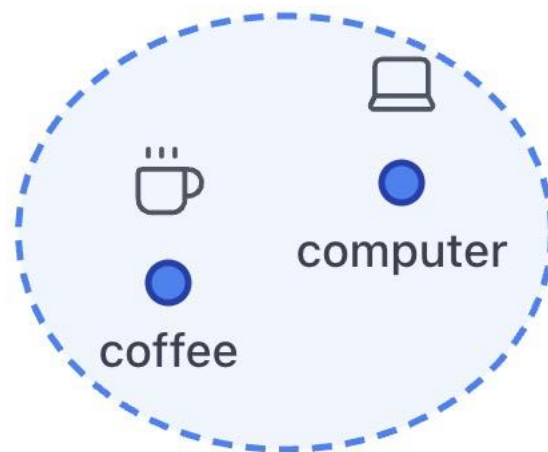
$\begin{bmatrix} 0.32, & -0.17 \end{bmatrix}$

**LAPTOP VECTOR**

$\begin{bmatrix} -0.61, & 0.41 \end{bmatrix}$

**CAT VECTOR**

$\begin{bmatrix} 0.15, & 0.65 \end{bmatrix}$



cat



**So, what? What can Vectors do for you?**



Find similar  
songs you'll like



Shop by  
taking a photo



Movies you  
might enjoy



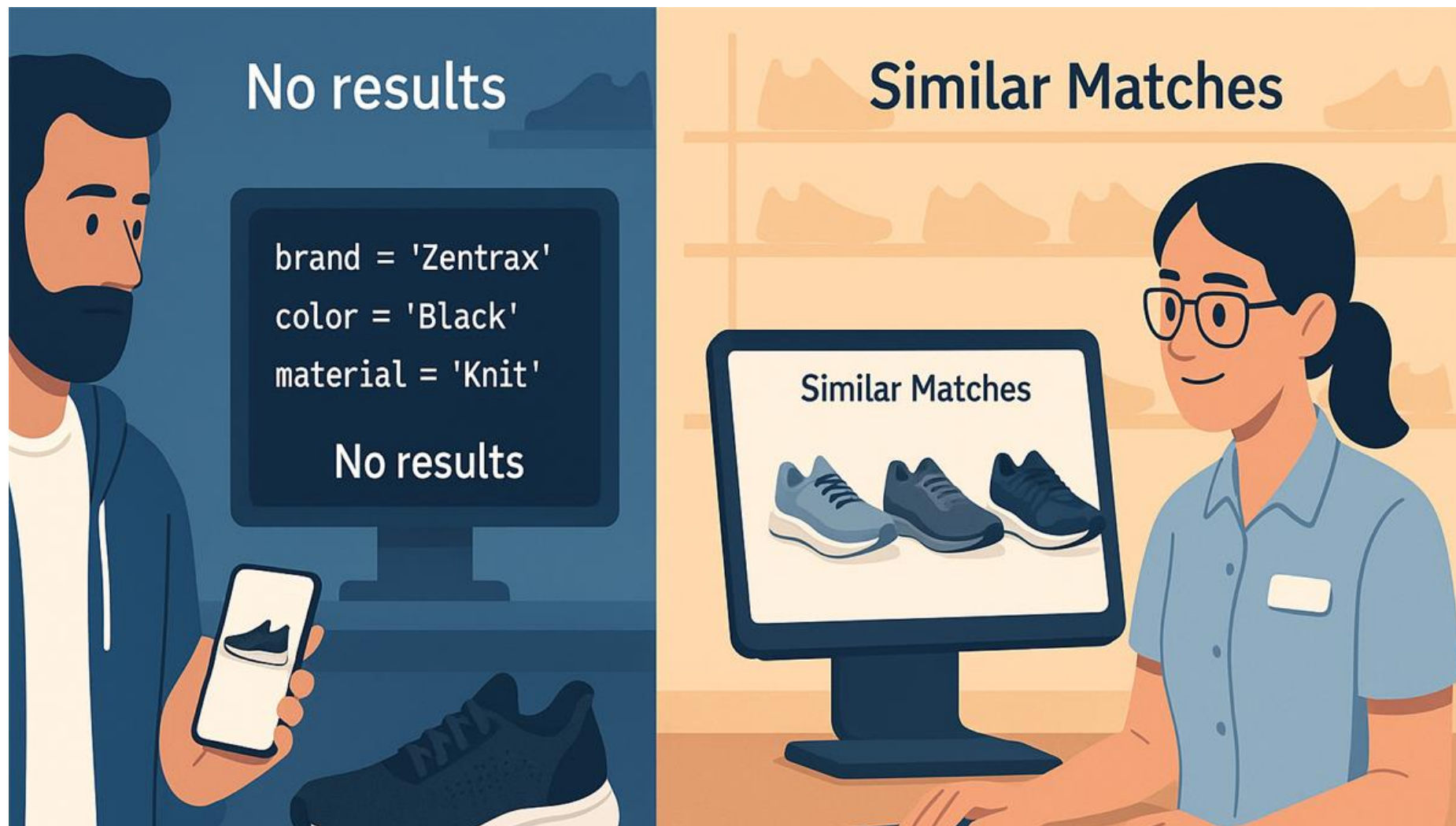
Find your  
photos by faces



Search using  
pictures



Ask questions,  
get answers



## Closed-Book: Vanilla LLM



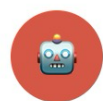
Relies only on internal knowledge  
— may forget or guess

## Open-Book: RAG (Retrieval-Augmented Generation)



Retrieves supporting information  
before answering





## Non-RAG LLM

Traditional Large Language Model

### User Query

Direct input to LLM



### Static Training Data

Fixed knowledge cutoff



### Generation

Response based on training only

Knowledge-Cutoff / Static Knowledge, Hallucinations



## RAG System

Retrieval-Augmented Generation

### User Query

Input converted to vector



### Retrieval

Search external knowledge base



### Augmentation

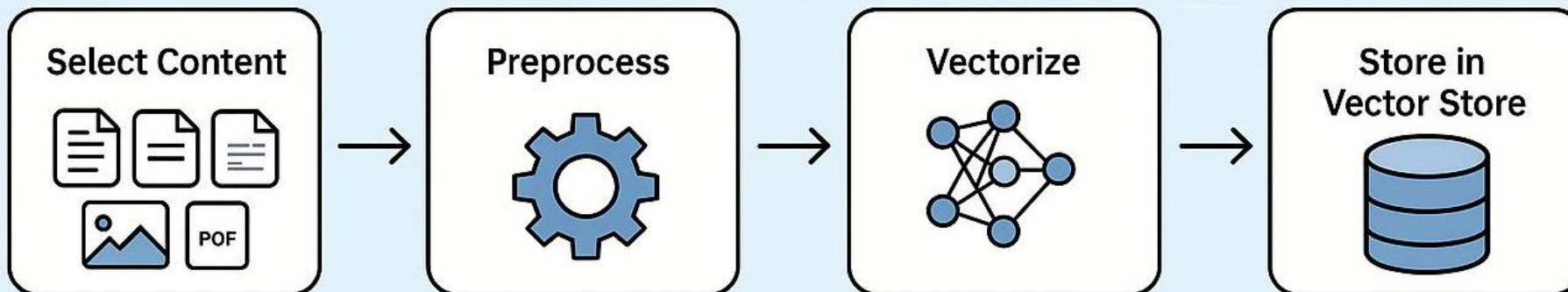
Combine query + retrieved context



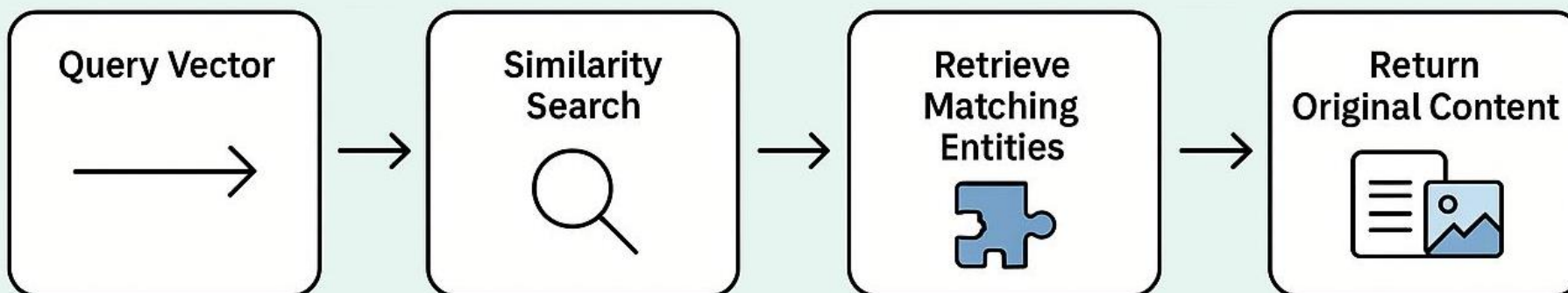
### Generation

LLM generates grounded response

Dynamic Knowledge, Domain-Knowledge, Source citation



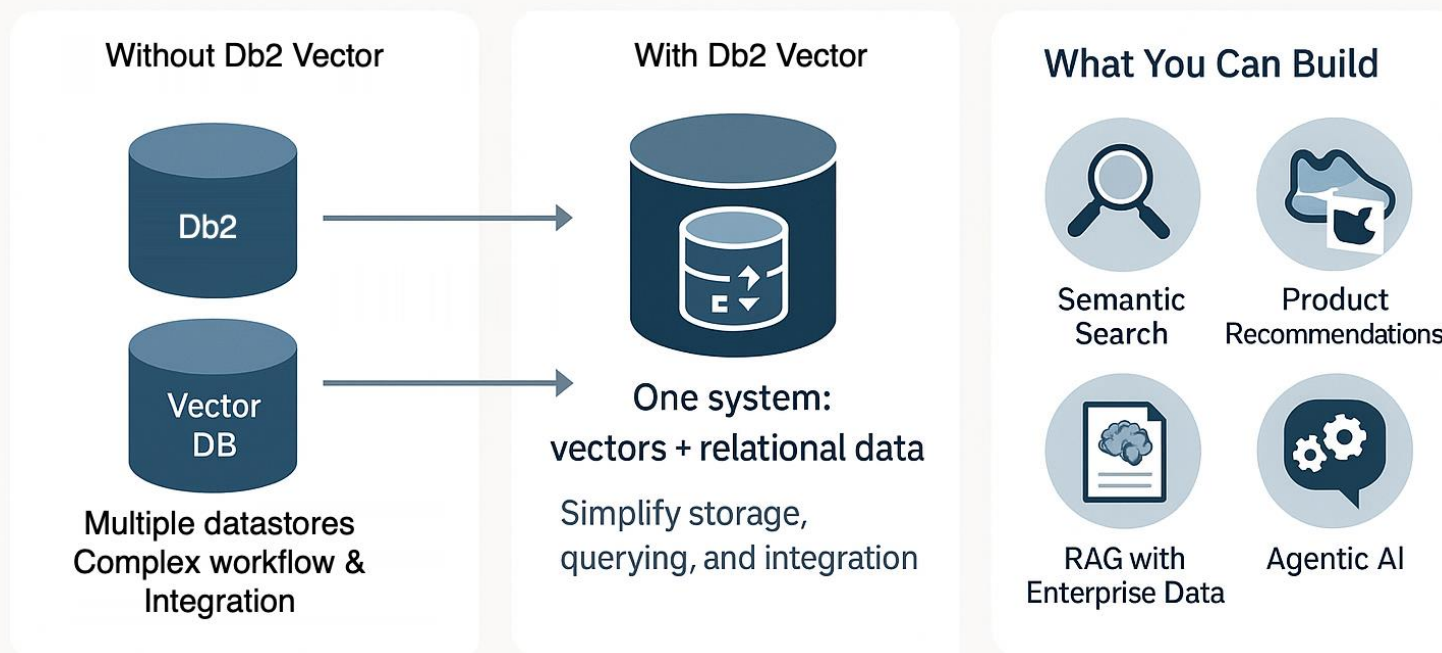
## PHASE 1: DATA INGESTION & VECTORIZATION



## PHASE 2: QUERY & RETRIEVAL

# Db2 Unlocks Semantic Search and GenAI Workflows

Built-in Vector Support for Modern AI Use Cases



Your data. Your vectors. One Db2 platform.

## **VECTOR type:**

- FLOAT32 vectors
- INT8 vectors
- Max dimension:
  - 8168 (FLOAT 32)
  - 32672 (INT8)

## **VECTOR Functions:**

- VECTOR\_DISTANCE – between vectors
- VECTOR constructor: string -> vector
- VECTOR\_SERIALIZE – vector -> string
- *VECTOR\_NORM* – vector magnitude
- VECTOR\_DIMENSION\_COUNT – number of dimensions

**LLM Apps Dev Support:  
LangChain Connector**



 NEW RELEASE

# Db2 12.1.3

## Production-Ready Vector Support

Data movement utilities • SQL routines • LLM framework integration



### LangChain & LlamaIndex

Native Python packages for RAG apps



### SQL Routines

Custom AI logic at database layer



### Data Movement

LOAD, backup, schema migration

Columnar Engine

VECTOR Data Type

Enterprise Scale

# Build LLM Apps with Db2



LangChain



LlamaIndex



Native Python



Quick Start



RAG Apps



AI Agents

## Db2 VECTOR Data Type

Cosine • Dot Product • Euclidean

✓ Enterprise Reliability   ✓ Python-Native APIs   ✓ No Low-Level SQL

# Custom AI Logic in SQL Routines

Encapsulate AI application logic at the database layer where your data lives



**Domain Logic**  
Custom metrics



**Reusable Operations**  
Normalization, validation



**Hybrid Queries**  
Vector + structured filters



**Batch Processing**  
Compute statistics



**Consistency**  
Same logic across apps



**Columnar Engine**  
Native optimization

✓ Write logic once in SQL • Execute where data resides • No network overhead

UDFs & Stored Procedures with VECTOR data types

# Production-Grade Data Movement

Vector support across Db2's data movement utilities

**Db2 12.1.2**

IMPORT • EXPORT



**Db2 12.1.3**

+ 6 More Utilities



**LOAD**

High-volume insertion



**ADMIN\_COPY\_SCHEMA**

Clone schemas



**ADMIN\_MOVE\_TABLE**

Relocate tables



**db2move**

Migrate databases



**External Tables**

Query in data lakes



**Logical Backup**

Schema-level backups

## Key Capabilities



Full & Incremental



Table-Level



External File

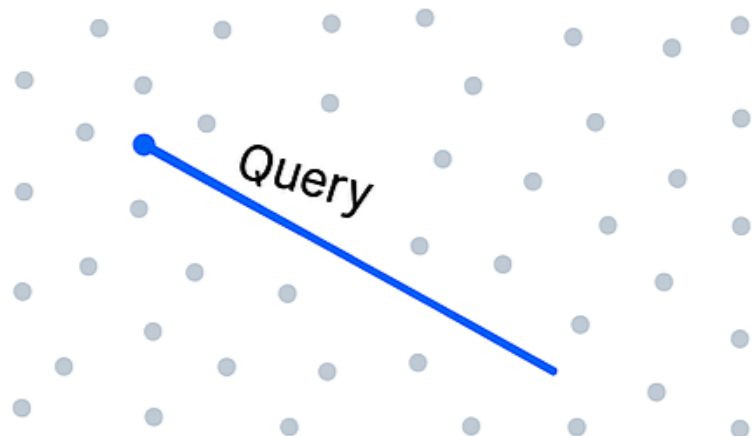


Production-Ready

✓ Move vector data using the same tools as traditional data types

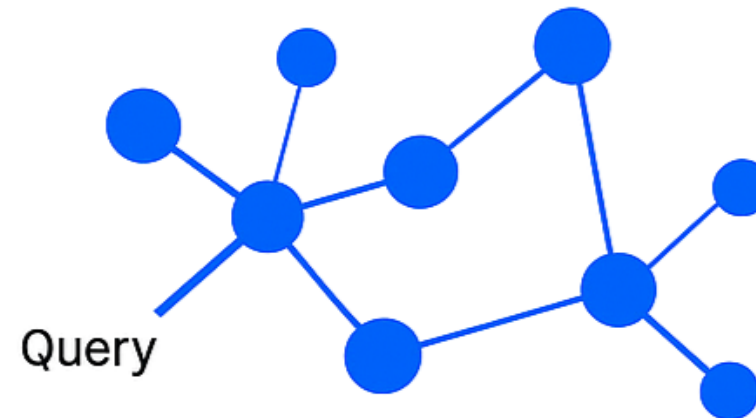


## Exact Search



Exact Search  
– slower for large collections

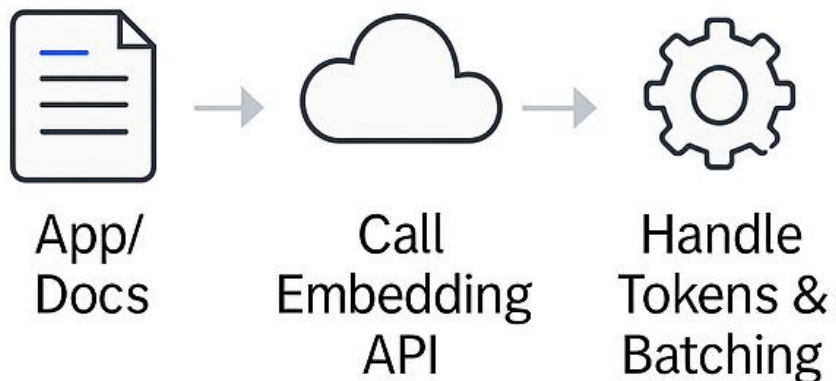
## Vector Index



Indexed search  
– faster at million / billion scale

**Faster vector search with Db2 Vector Index**

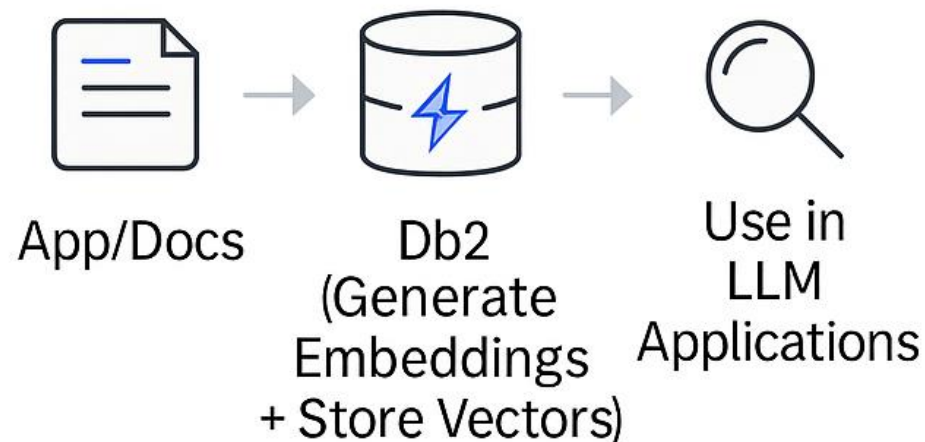
## The Problem: Too Much Glue Code



## Challenges:

**Time • Fragility • Cost**

## The Solution SQL-Native Embeddings in Db2

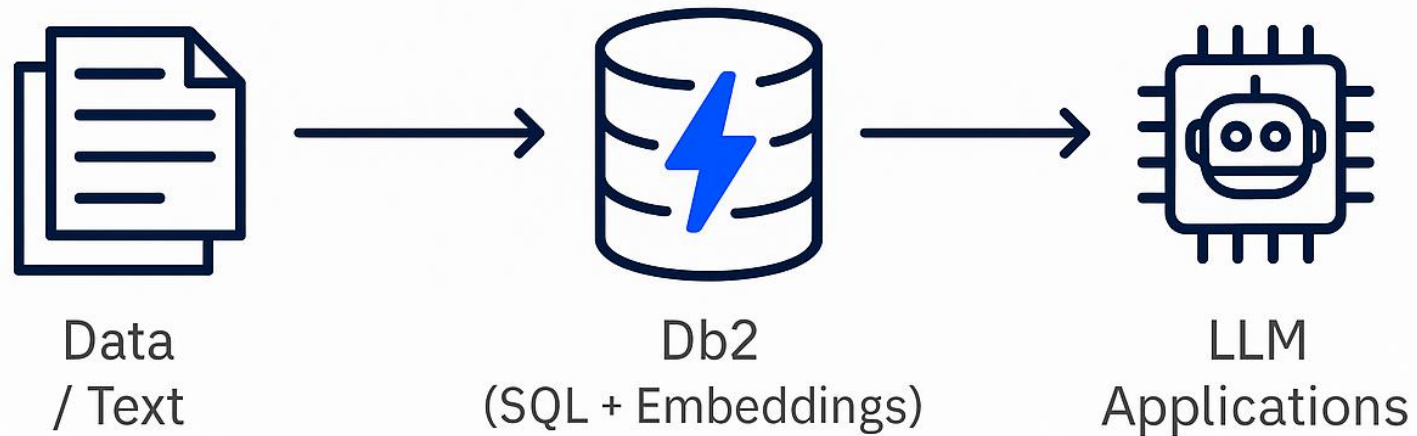


With Db2, register your embedding endpoint once. Then generate, store, and search embeddings SQL.

## Demo

# SQL-Native Embeddings in Db2

Simplifying LLM application development with fewer moving parts



**[ibm.biz/learndb2ai](https://ibm.biz/learndb2ai)**





**Shaikh Quader**

AI Architect @ IBM Db2 | Ph.D. Candidate in AI | IBM  
Master Inventor | Sharing lessons from building & d...

