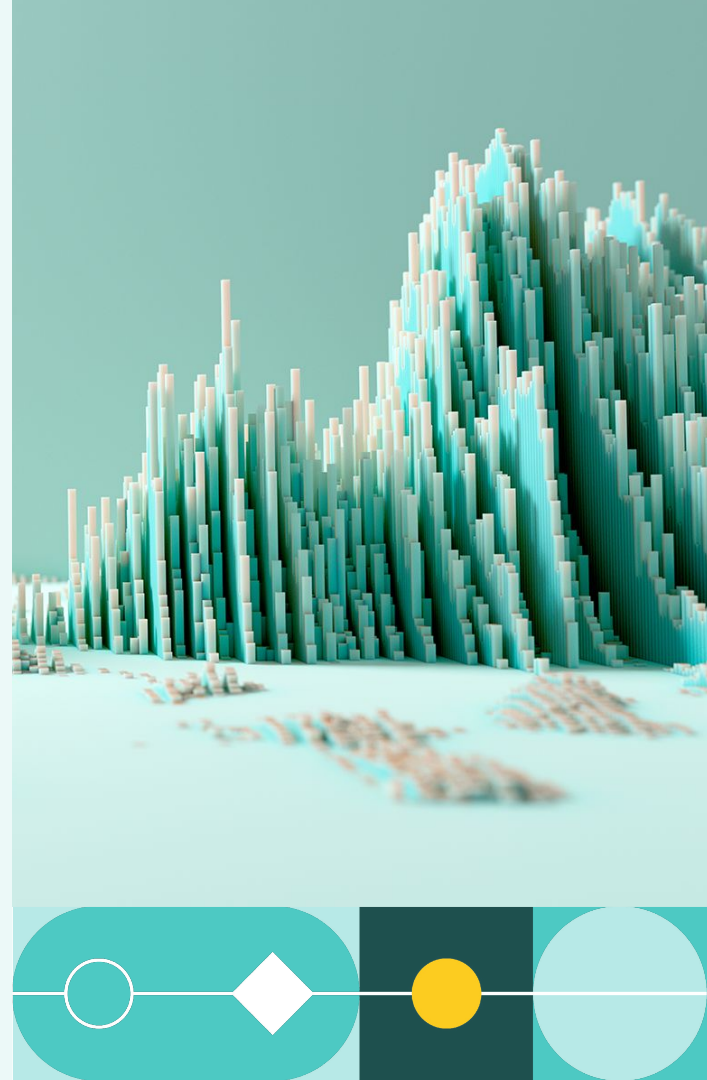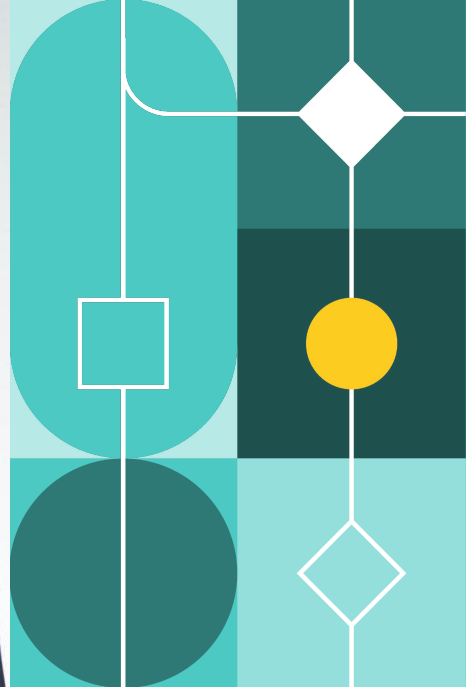# Storage in the AI era
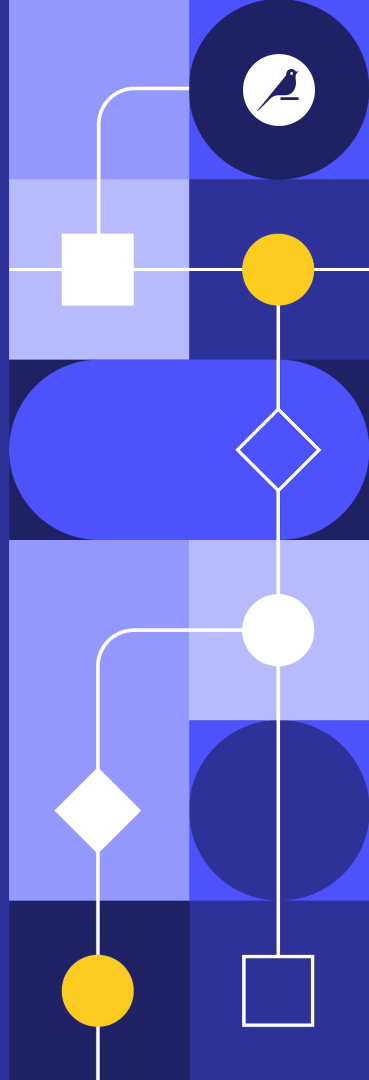
**Arnaud Pichery**

VP Engineering, Dataiku
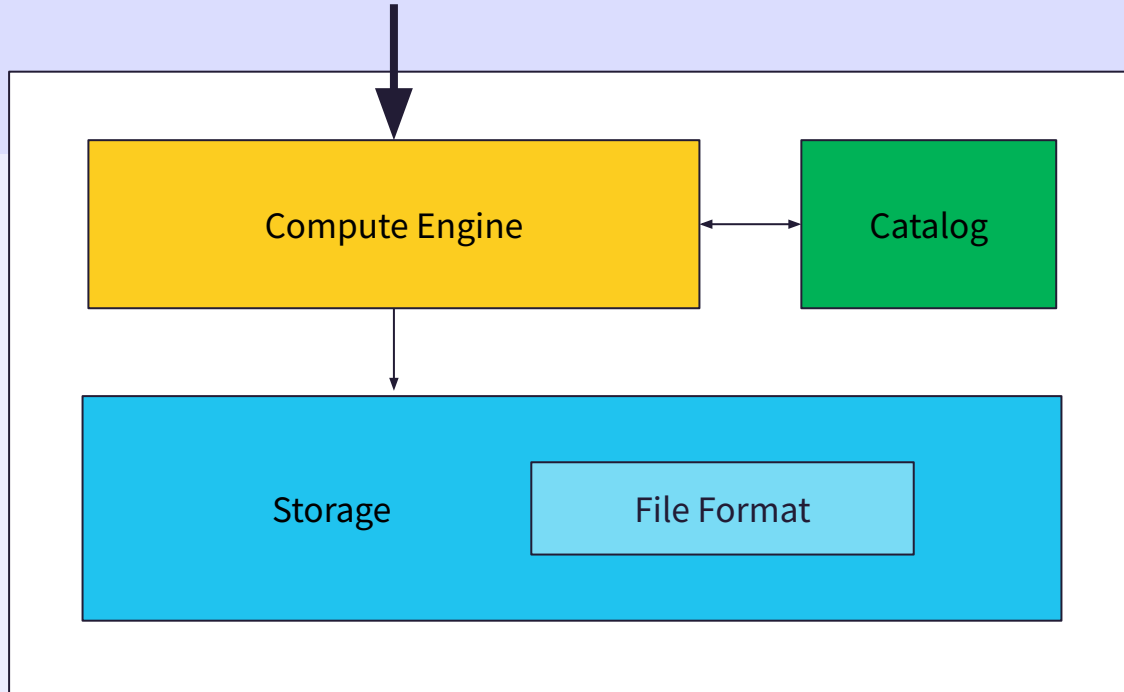
# Data Warehouses

# Data Warehouse

Analytics on Big Data



**(Semi-)structured data**

Separation of Storage and Compute

Massively Parallel Processing

# Cloud Data Warehouses

SQL champions

SQL support

Cloud native architecture

Separation of Storage and Compute

Massively Parallel Processing (MPP)

Support for Modern Data Formats

# Cloud Data Warehouses

SQL champions

**＋**

Simplicity of use and administration

Decouples storage and compute

Infinitely scalable

Contains a storage engine that optimizes data layout

**－**

Data stored in proprietary format
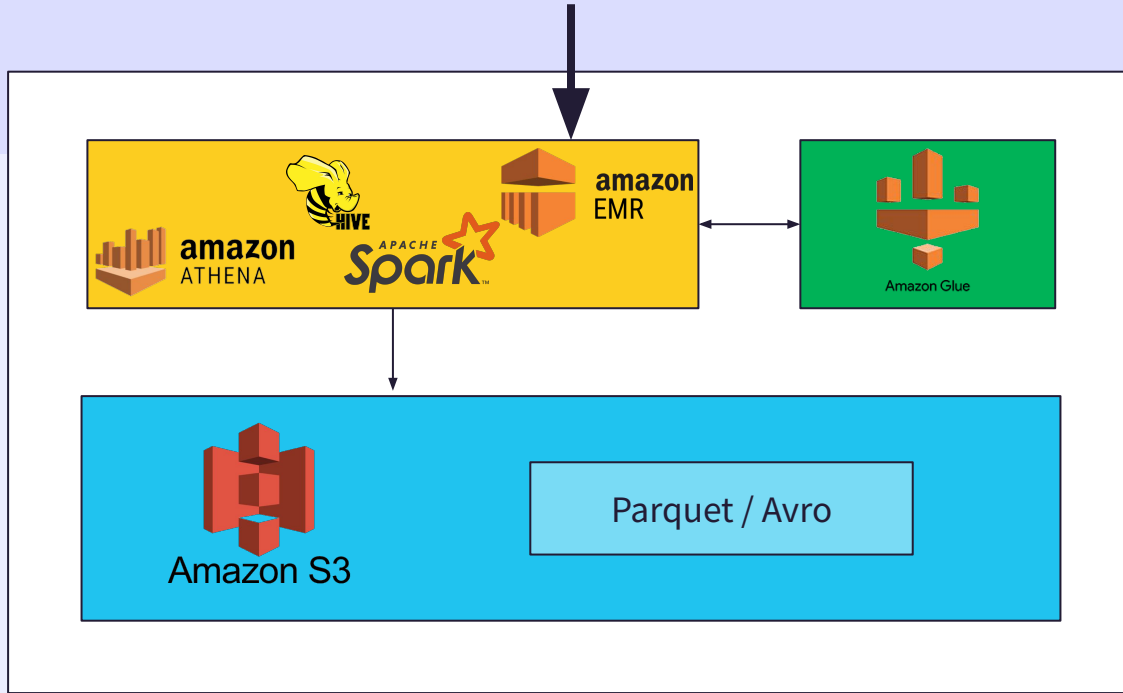
Need to land data in cloud storage to ingest

Expensive ($$$)

Cannot handle unstructured data nor non-SQL workloads

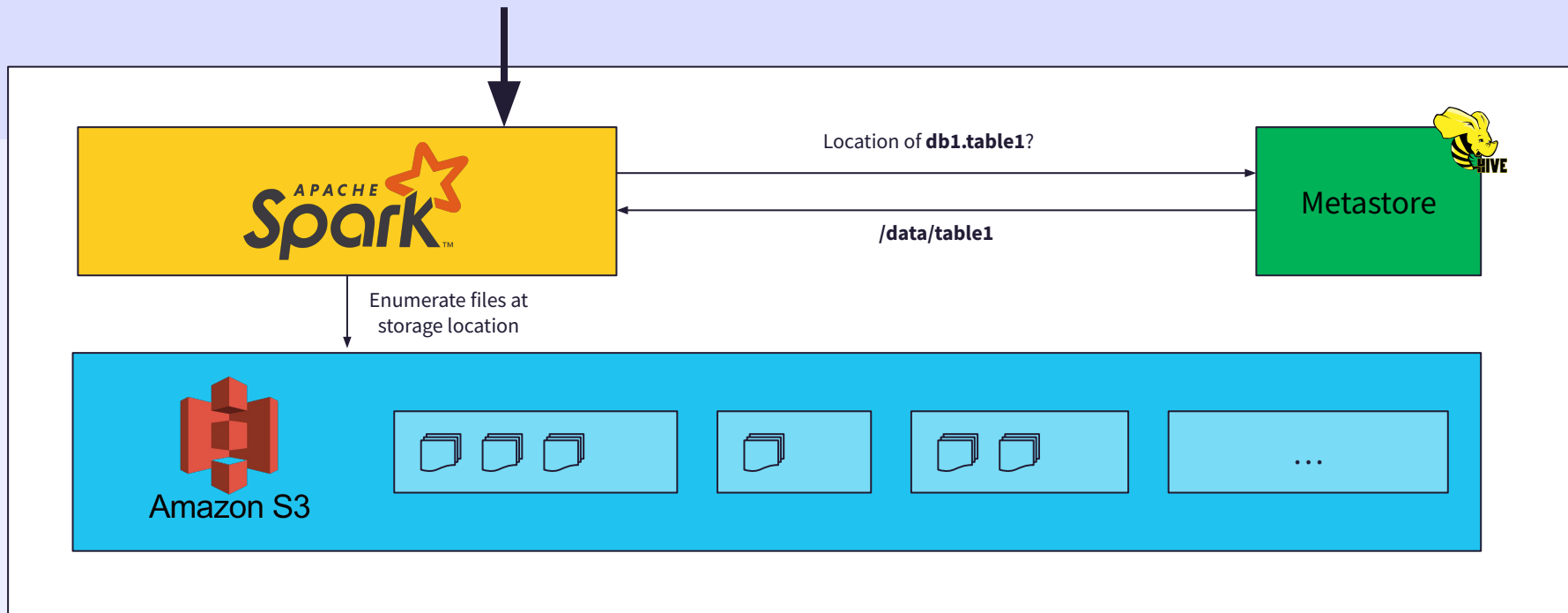# Data Lakes
Analytics on Big Data



**Hadoop-like architecture**

Separation of Storage and Compute

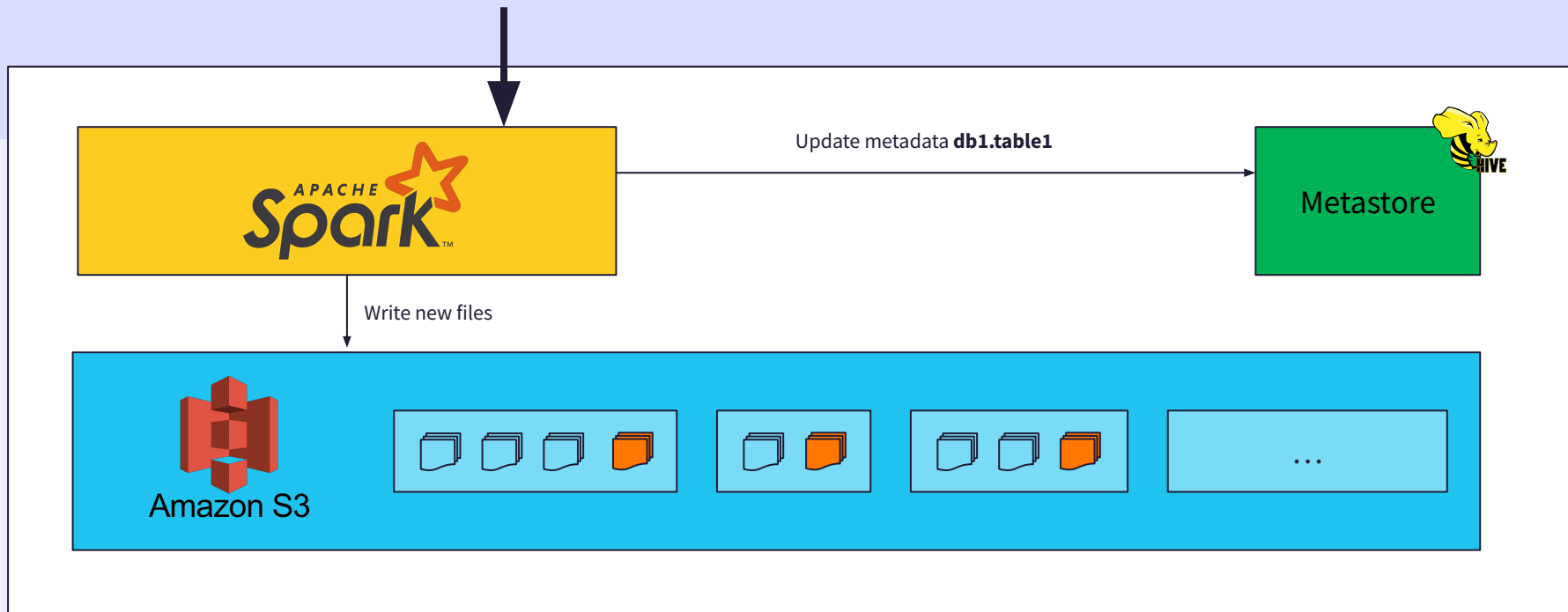Massively Parallel Processing

# Data Lakes

Analytics on Big Data

# Data Lakes

Analytics on Big Data

# Data Lakes

Analytics on Big Data

**+**

Fully open-source storage formats

Modular architecture

Lower cost ($$)

Can process unstructured data

Provides SQL and Python/R/Scala APIs for data processing

**−**

Complex tuning of computation engines

Generally less performant than equivalent data warehouses

Lack of ACID transaction guarantees

# The hot format

The new shiny thing
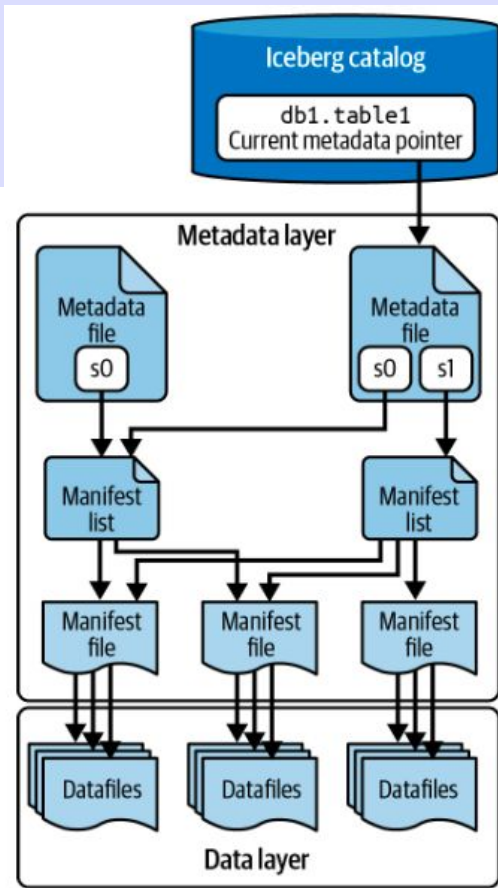
# Data Lakes on Iceberg

Open-source champion



Iceberg catalog points to <u>location of a metadata file</u>.

**Metadata file** contains references to one or more data snapshots, current snapshot, partition information and schema.
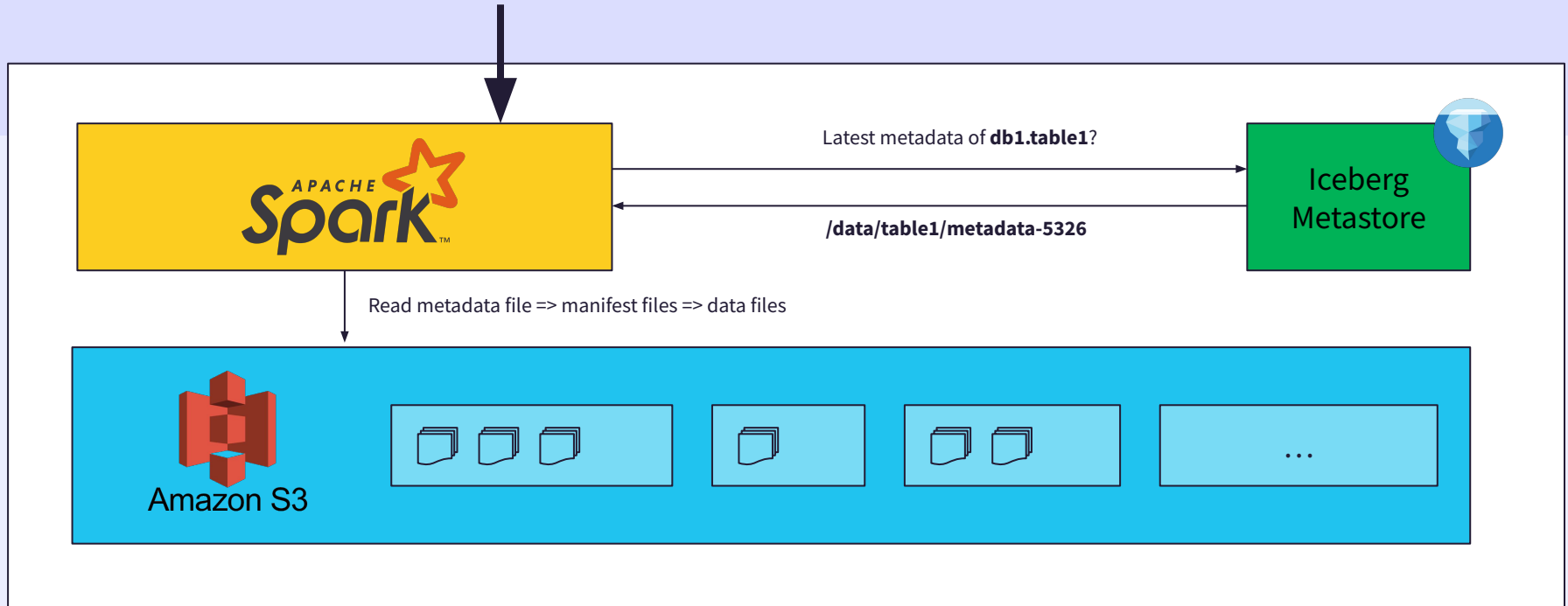
**Manifest list** is a snapshot of an Iceberg table at a point in time. Contains references to manifest files.

**Manifest files** keep track of datafiles and stats (e.g. partition membership, record counts, max/min, etc.).
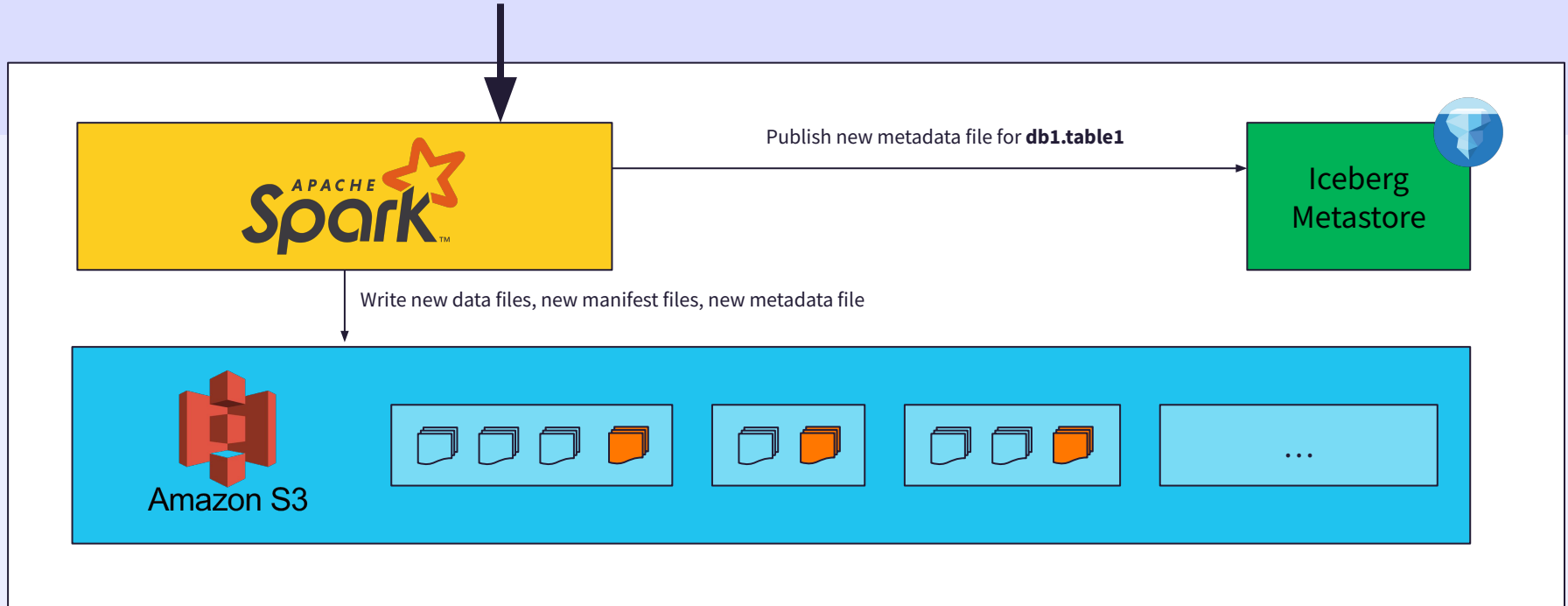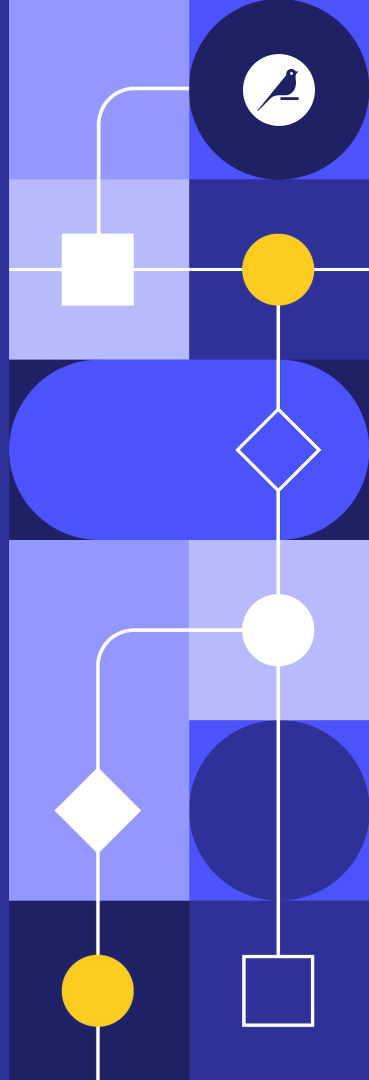
# Data Lakes on Ice(berg)

Upright spin

# Data Lakes on Ice(berg)

Triple axel jump

# Vector Stores

# Vector databases

New kids on the block

**RAG (Retrieval-Augmented Generation)**

➔ Store billions of vectors

➔ Hybrid storage

- ○ Unique ID
- ○ Vector (384/512 doubles)
- ○ Metadata (tags, timestamp, …)
- ○ Raw data (text, image, pdf, …)
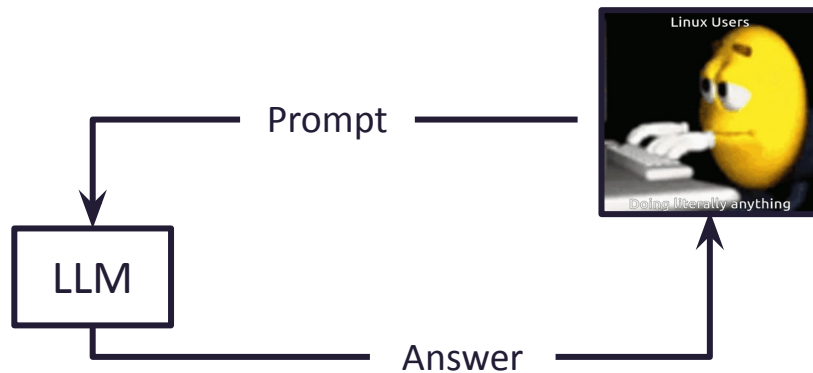
➔ Fast lookup of similar vectors

# Vector databases

New kids on the block

**RAG**



Prompt

LLM

Answer

# Vector databases

New kids on the block

## RAG



Search

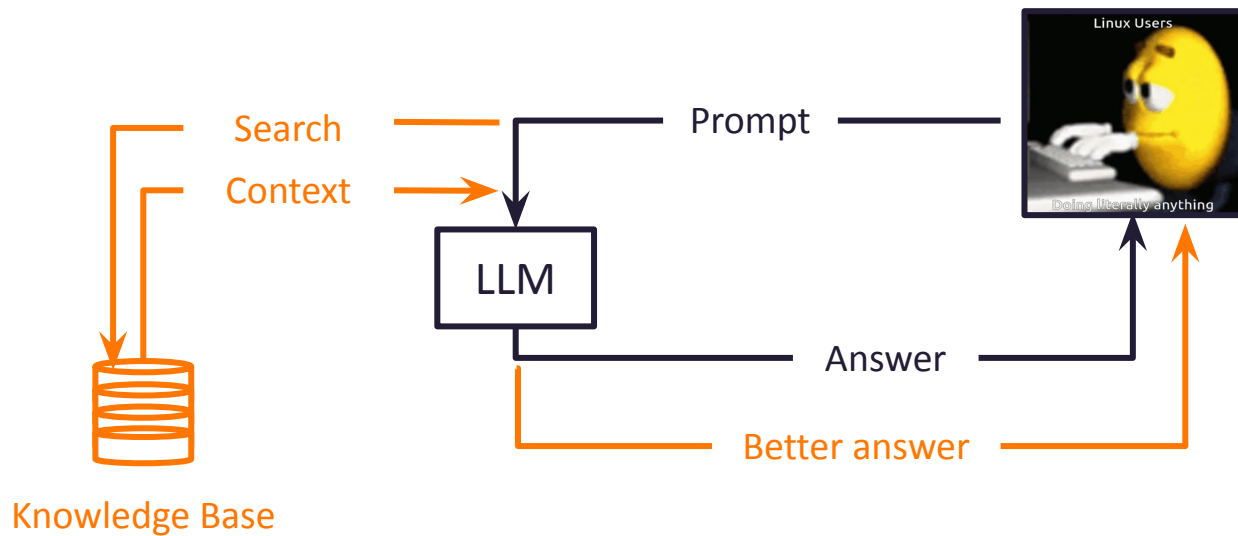Context

Prompt

LLM

Answer

Better answer

Knowledge Base
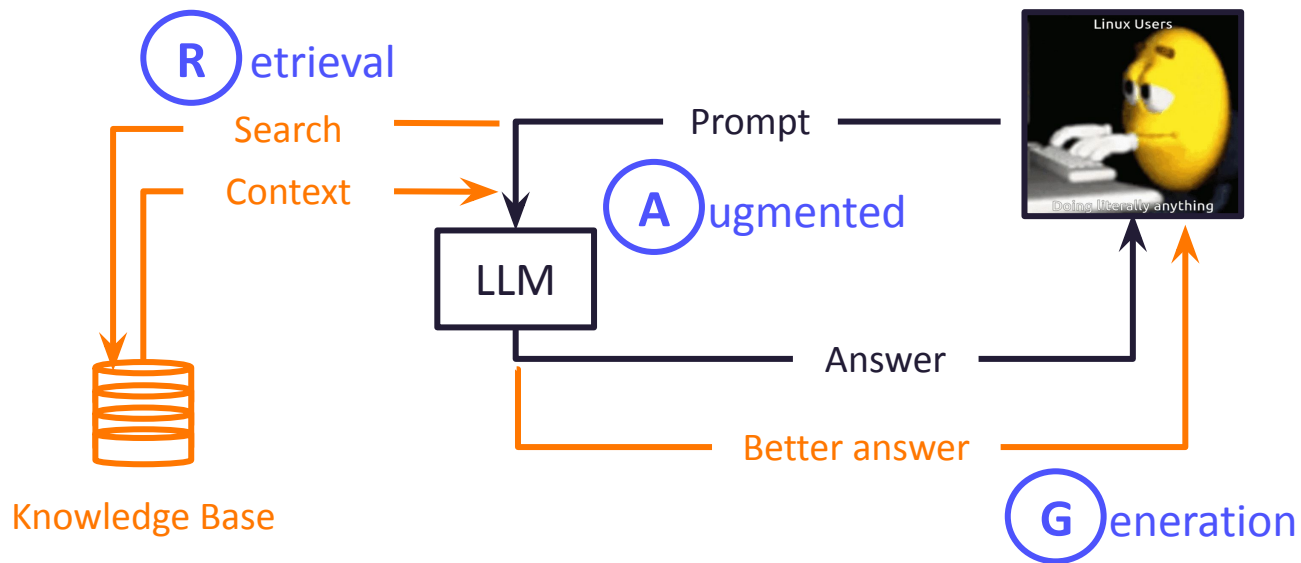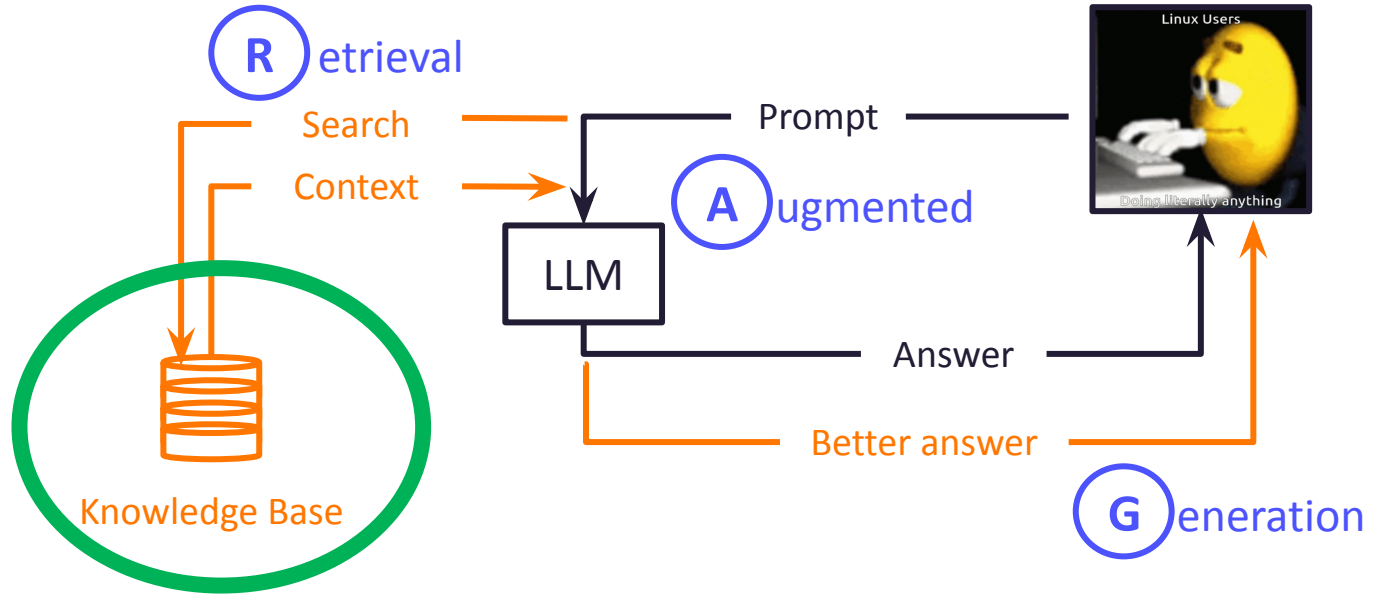
# Vector databases

New kids on the block

**RAG**
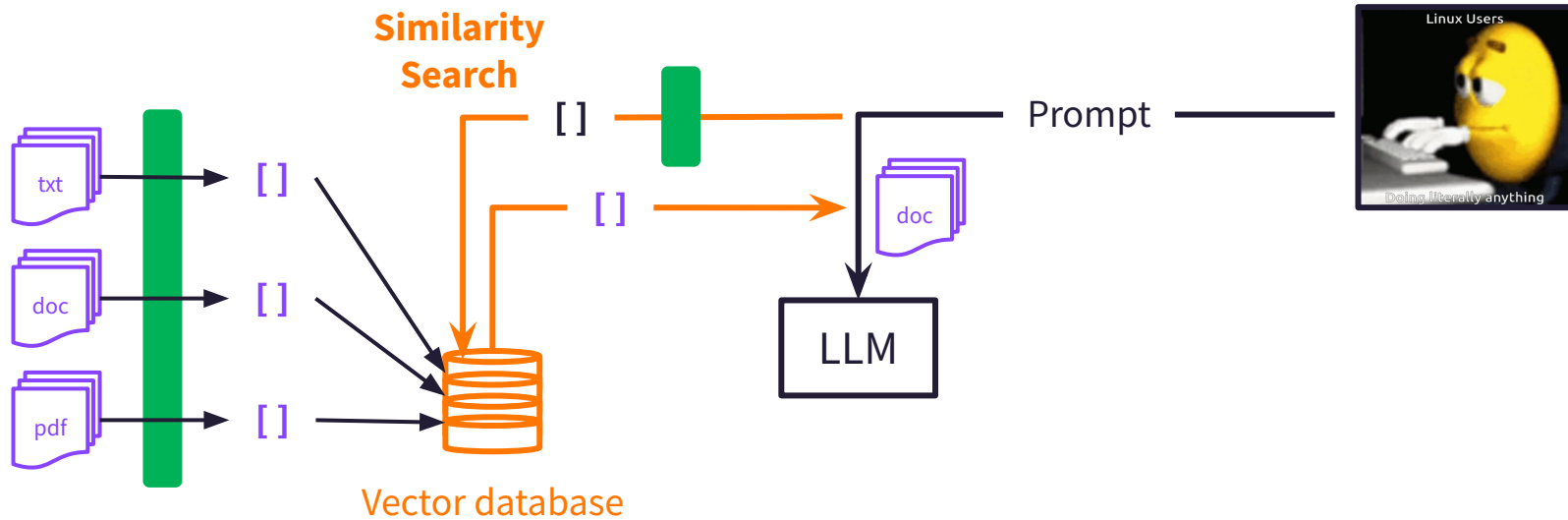
# Vector databases
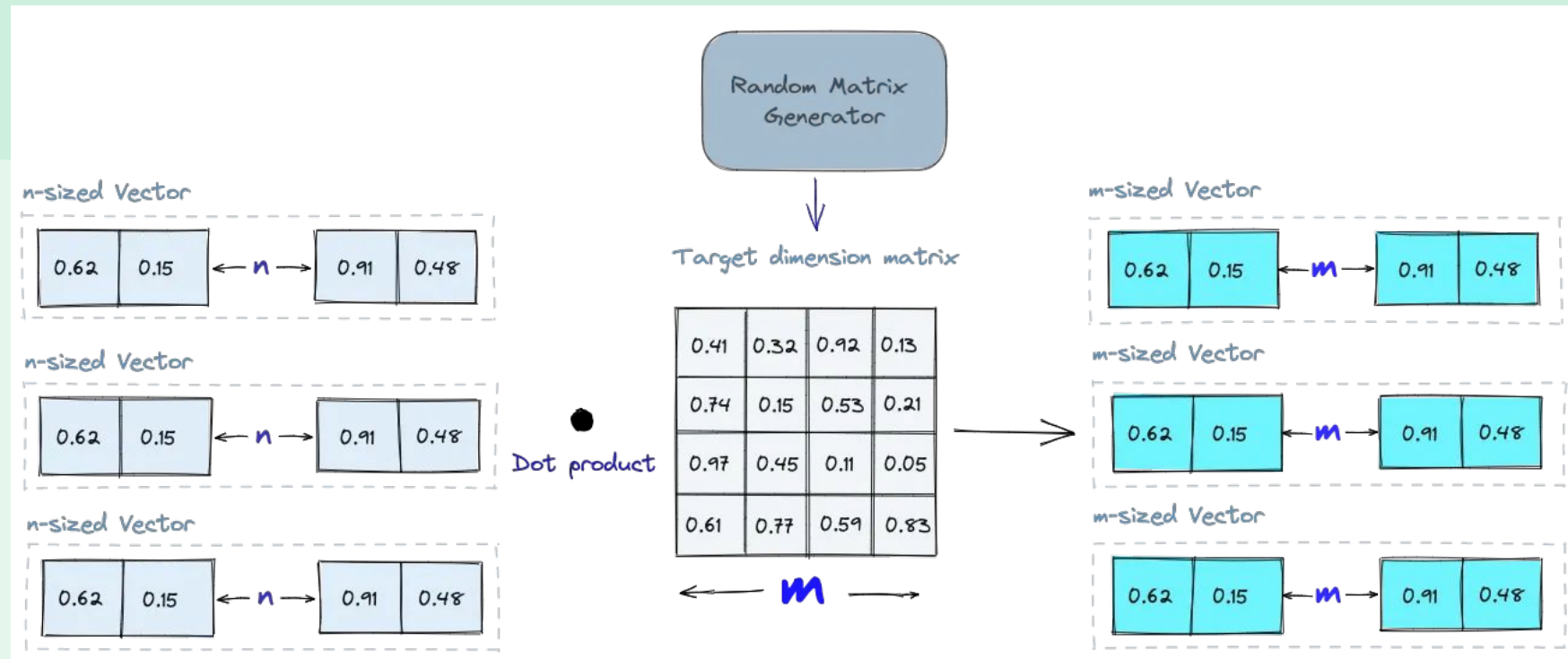
New kids on the block

## RAG

# Vector databases

New kids on the block

**RAG**
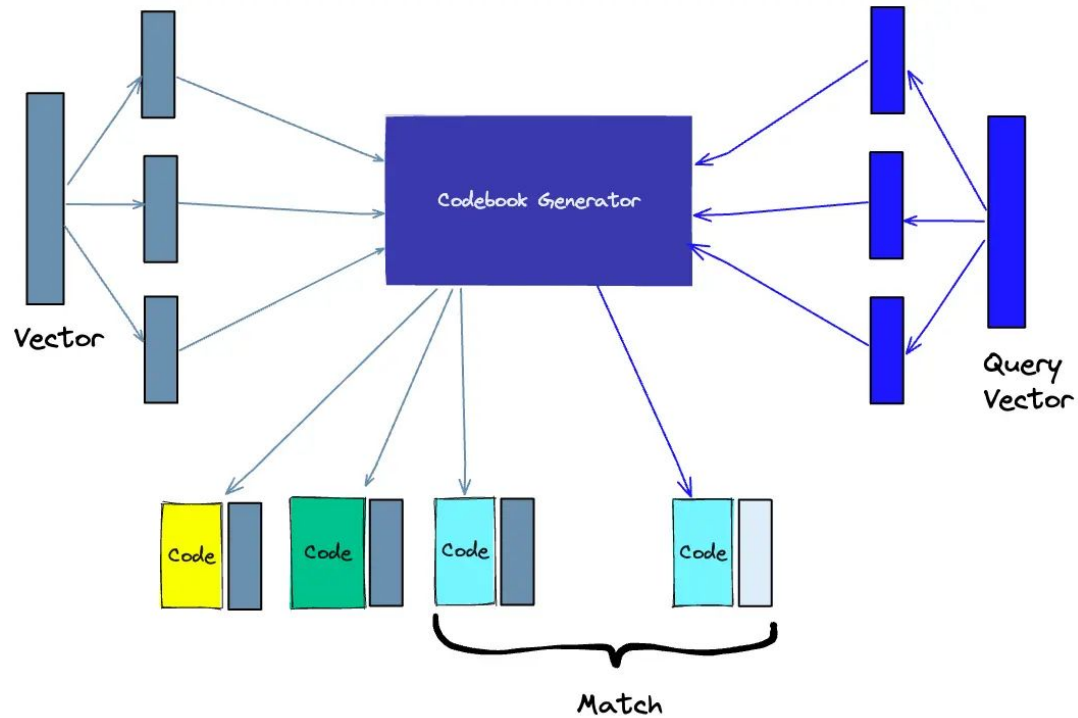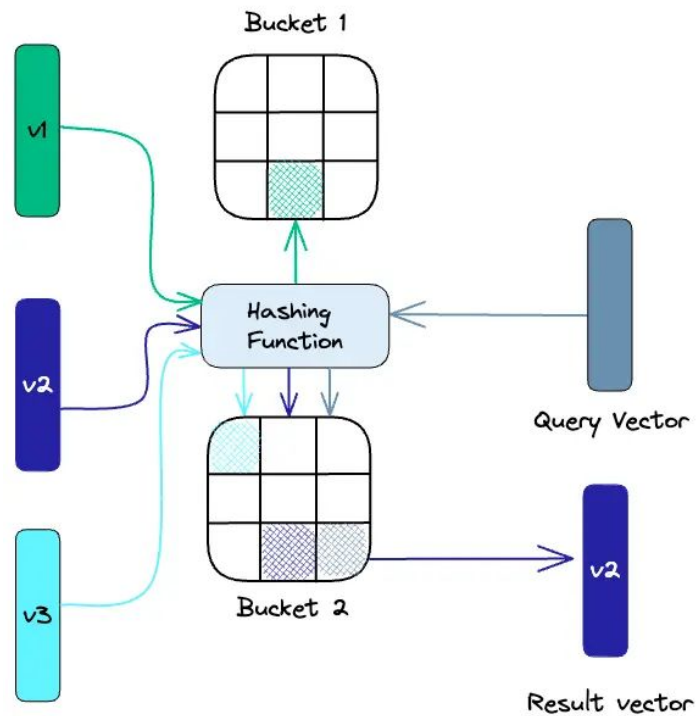
# Vector databases

Random Projection

# Vector databases

Product Quantization

# Vector databases

Locality-sensitive hashing

# Thank You

Do you have any questions?