VALENTINA ALTO

# Getting Started with MM RAG

**An implementation with GPT-4v and Llama-Index**

# Agenda

**Tip:** Download the workshop material to start familiarizing with code and slides.

**How:** Scan the QR code below!

| | |
|---|---|
| **1** | Intro to RAG |
| **2** | Embeddings and VectorDB |
| **3** | Multimodal RAG |
| **4** | Demo Time |

**LLMs as "brains" for our applications**

Conversational Front-end

How can I assist you today?

Your Data

Web Search

Math

Code executor

DB connector

Tools

Reasoning engine

# Anatomy of an LLM-powered application

LLM-powered applications open the way to a new landscape of components

Frontend

↑ ↓

Orchestration

Prompt & response filtering

Metaprompt

↑ ↓

Data grounding

Tools

↑ ↓

Large Language Models

↑

AI infrastructure
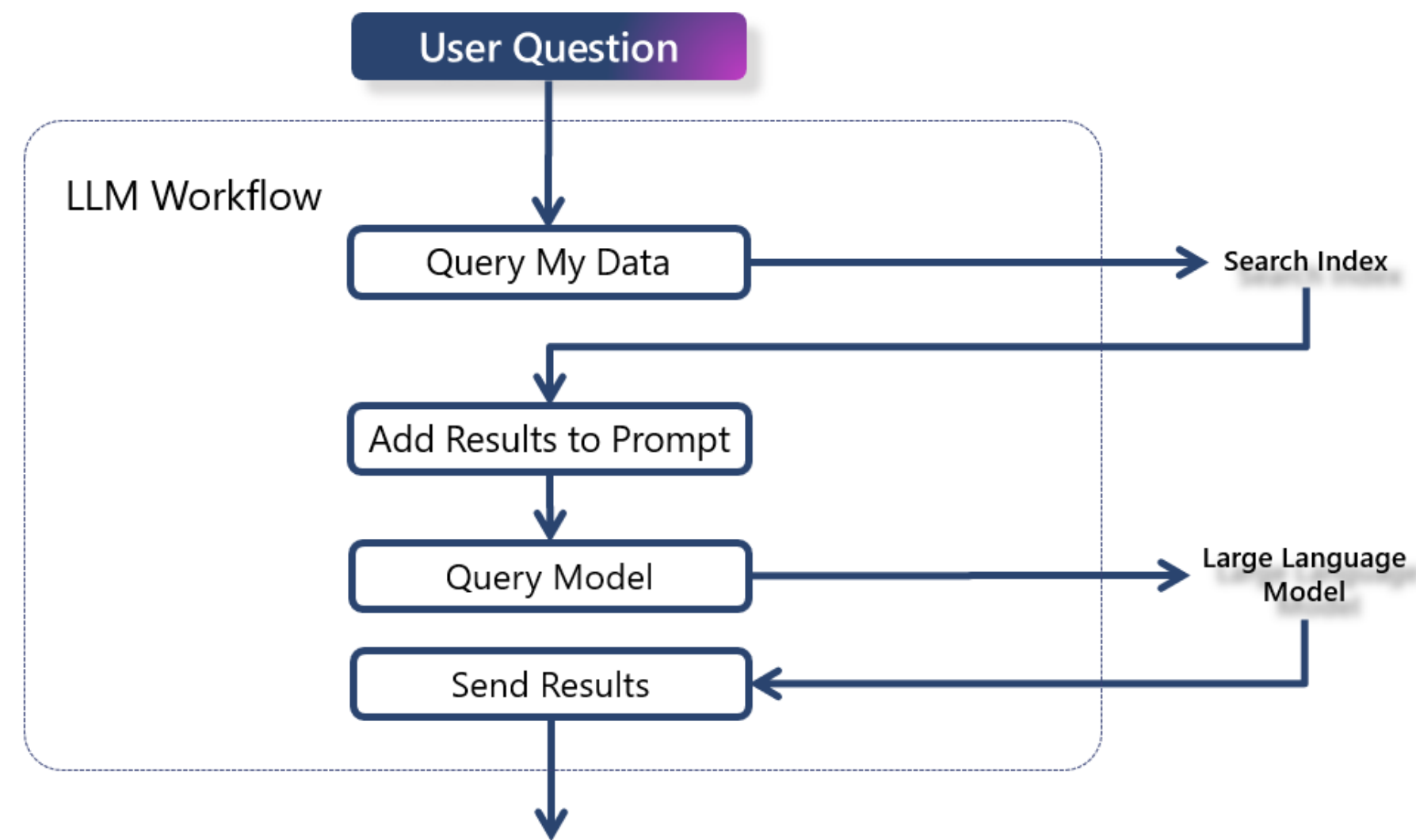
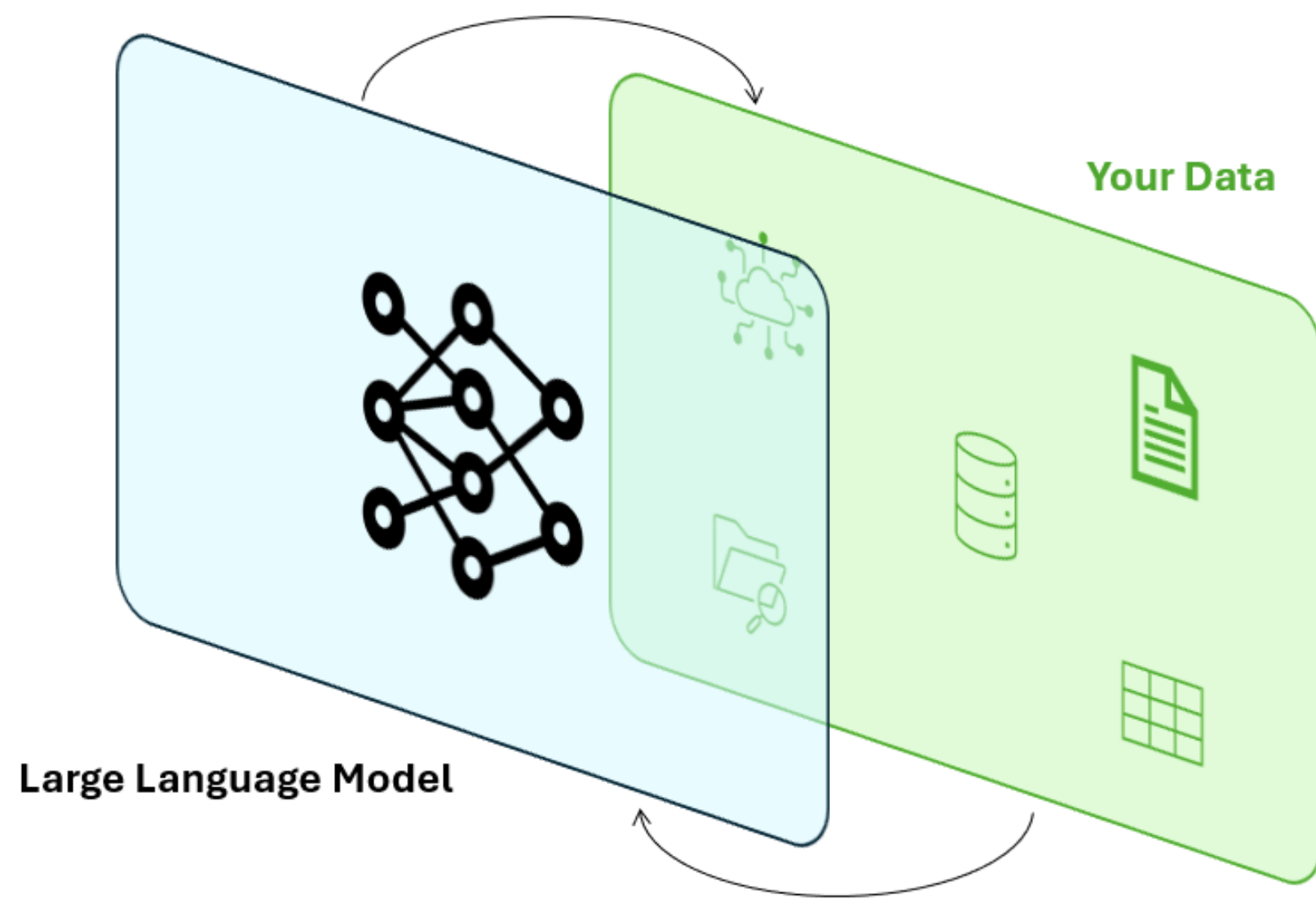# Problem: Generative AI doesn't know about your data

What if data you are interested in are not part of the training dataset?
- Personal Data (confidential, not public...)
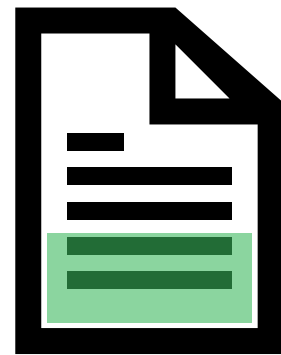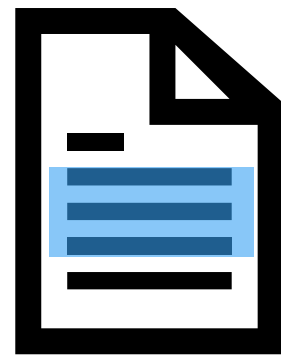- Up to date data
- Application data

Does my insurance plan cover eye exams?

I'm sorry, but as an AI assistant I don't have access to personal information.

Training data

Parametric Knowledge

# Introducing Retrieval Augmented Generation



Your Data

Large Language Model

**LLM Workflow**

User Question

Query My Data → Search Index

Add Results to Prompt

Query Model → Large Language Model

Send Results

# Augmentation

## Context

Tips for beginners: start…

Pegasus Shoe 123: best experience…

Milano Shoe 456: solid yet…

## Prompt Engineering

Instructions    Grounding

Safety    Personality    …

## System Message

You are an AI assistant that helps users answering their query.

# Documentation
The following documentation should be used in the response.
{retrieved_docs}

Tips for beginners: start…

Pegasus Shoe 123: best experience…

Milano Shoe 456: solid yet…

#Safety
You **should always** reference factual statements to search results based on retrieved docs.

# Generation

**System message + retrieved documents**

You are an AI assistant that helps users answering their query.

# Documentation
The following documentation should be used in the response.
{retrieved_docs}

Tips for beginners: start…

Pegasus Shoe 123: best experience…

Milano Shoe 456: solid yet…

#Safety
You **should always** reference factual statements to search results based on retrieved docs.

**User's query**

What is the best equipment for beginner climbers?

**+**

Generative Model (e.g. GPT-4

"According to the catalogue, if you are about to start climbing…"

# How do we retrieve relevant documents?

<3d style illustration of two cats playing a card game, one cat holding the deck of cards, the other cat thinking about which card to draw from the deck>
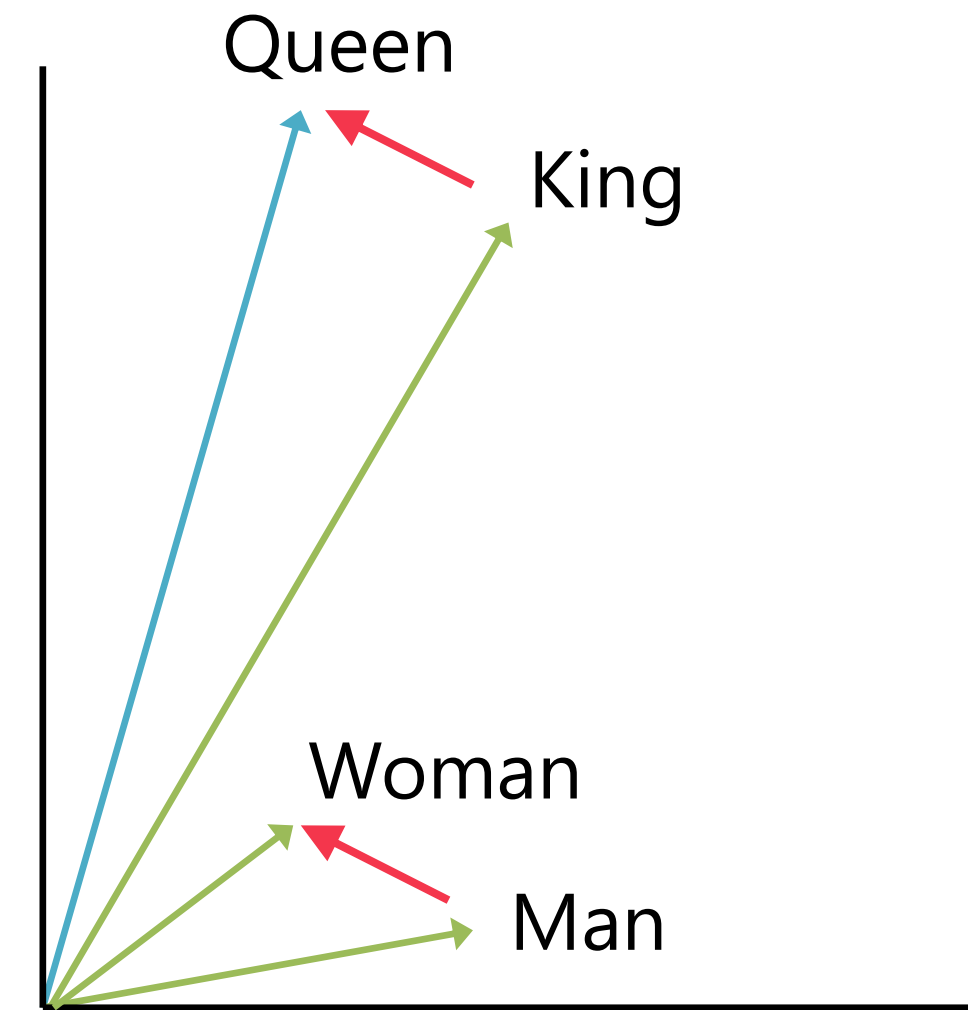
An embedding is a way of representing high-dimensional, non-numeric data, such as words or sentences, in a lower-dimensional space, such as vector.

A text embedding can capture the semantic and syntactic features of the text, such as meaning, context, and similarity.

Each embedding is a vector of floating-point numbers, such that the distance between two embeddings in the vector space is correlated with semantic similarity between two inputs in the original format.

For example, if two concepts are similar, then their vector representations should also be similar.

King-Man+Woman ≈ Queen

Embeddings represent your data, and each dimension represents a feature of that data.
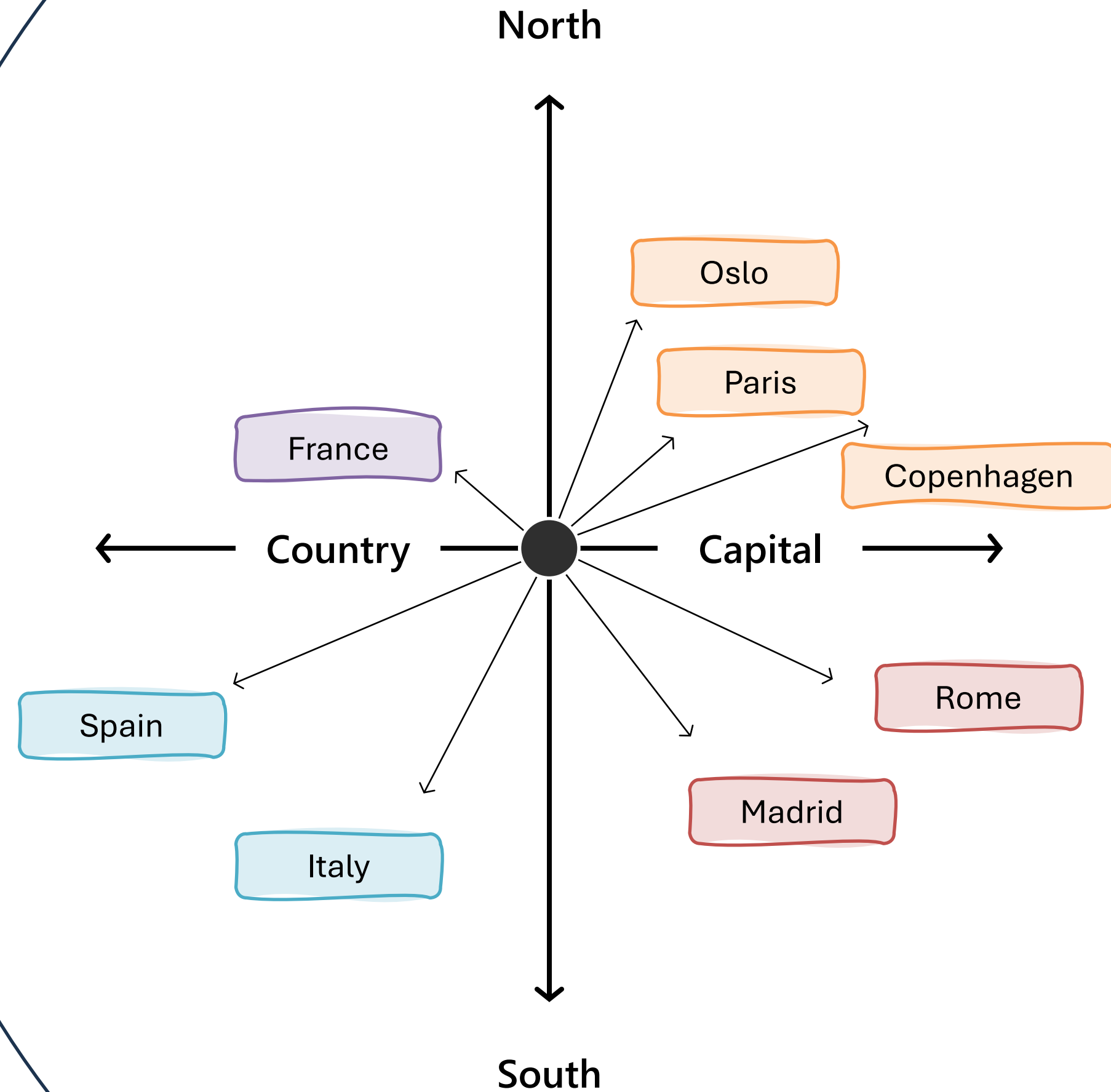
France

Paris

← Country ●── Capital →

For example, one dimension could be the geographic connotation (country vs capital), another one the geographic position (north vs south).
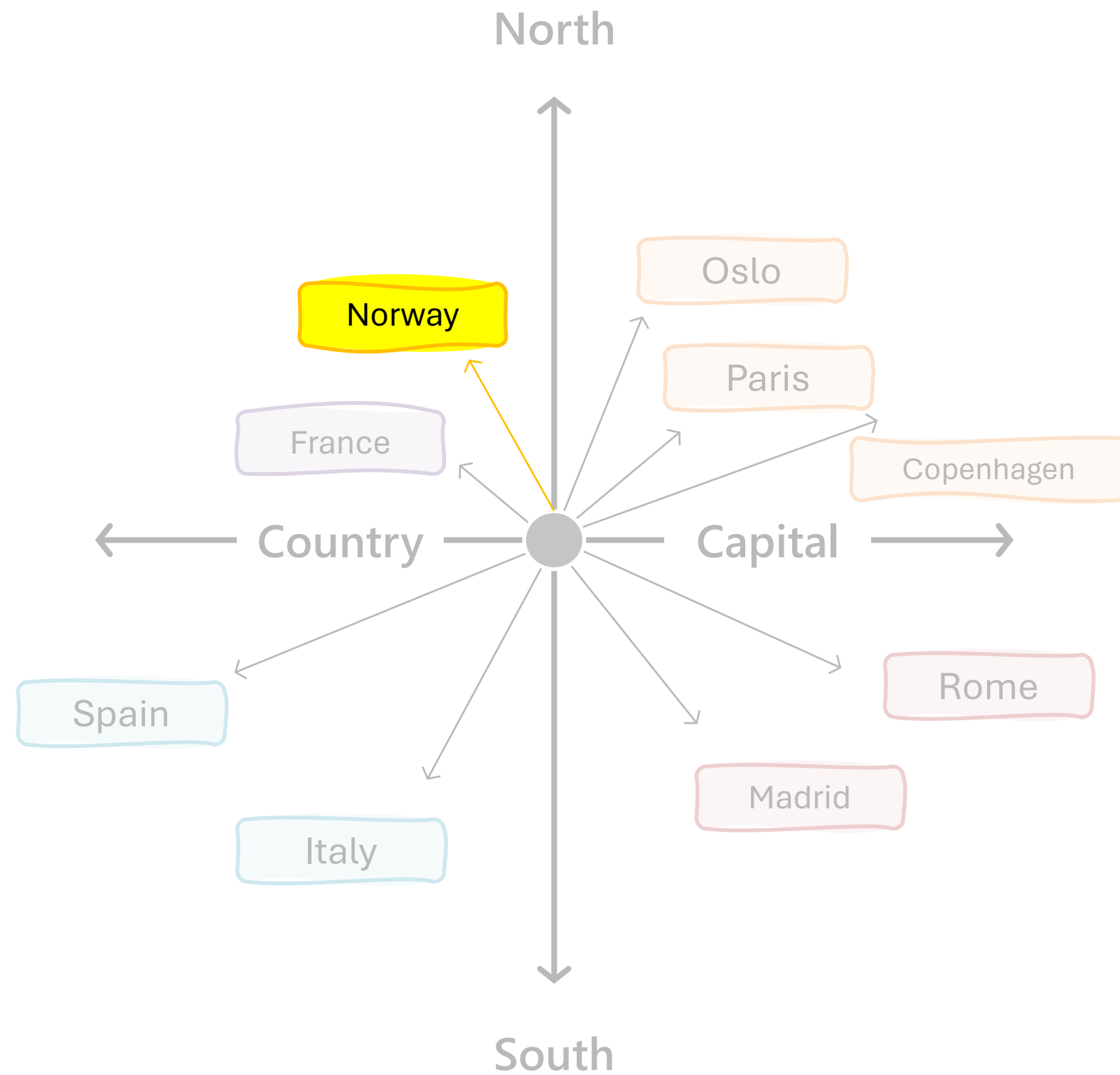
North

France

Paris

Country

Capital

South

In the embedding space, similar concepts (words, sentences, documents) should be close in mathematical distance.
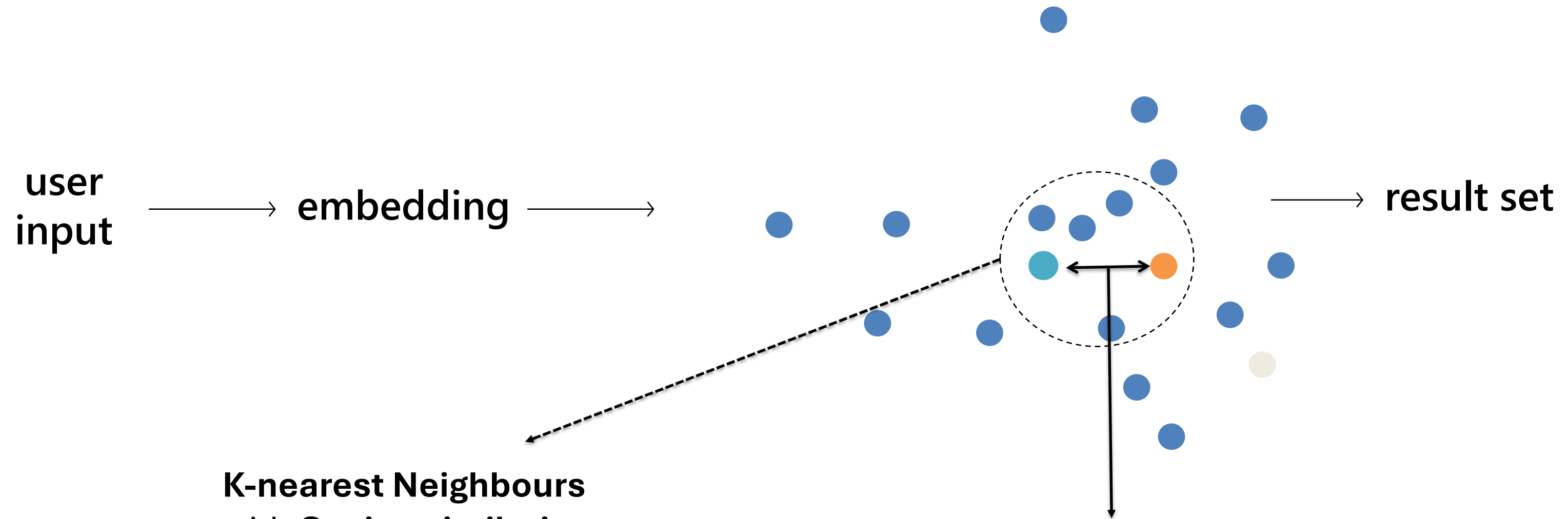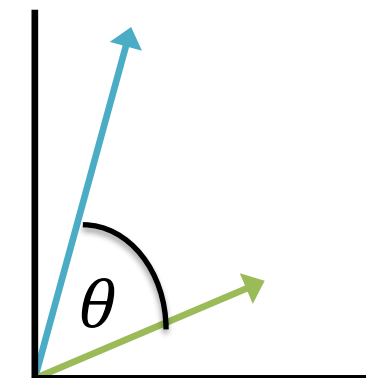
**Vector search ranks objects by similarity (relevance) to the query**

Norway

| Relevance | Result |
|-----------|--------|
| Query | Norway |
| 1 | France |
| 2 | Oslo |
| 3 | Copenhagen |
| 4 | ... |

# Similarity Search with Embeddings



user
input → embedding →                                     → result set

**K-nearest Neighbours**
with **Cosine similarity**
distance

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2}\sqrt{\sum\limits_{i=1}^{n} B_i^2}}$$

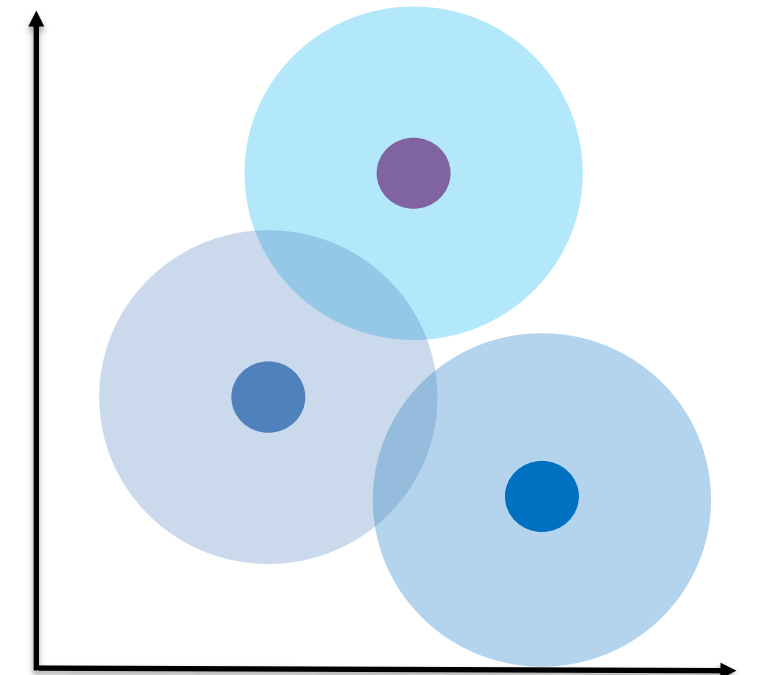A new entry in the landscape of DB

Documents DB

Key-value DB

Wide-columns DB

Graph DB
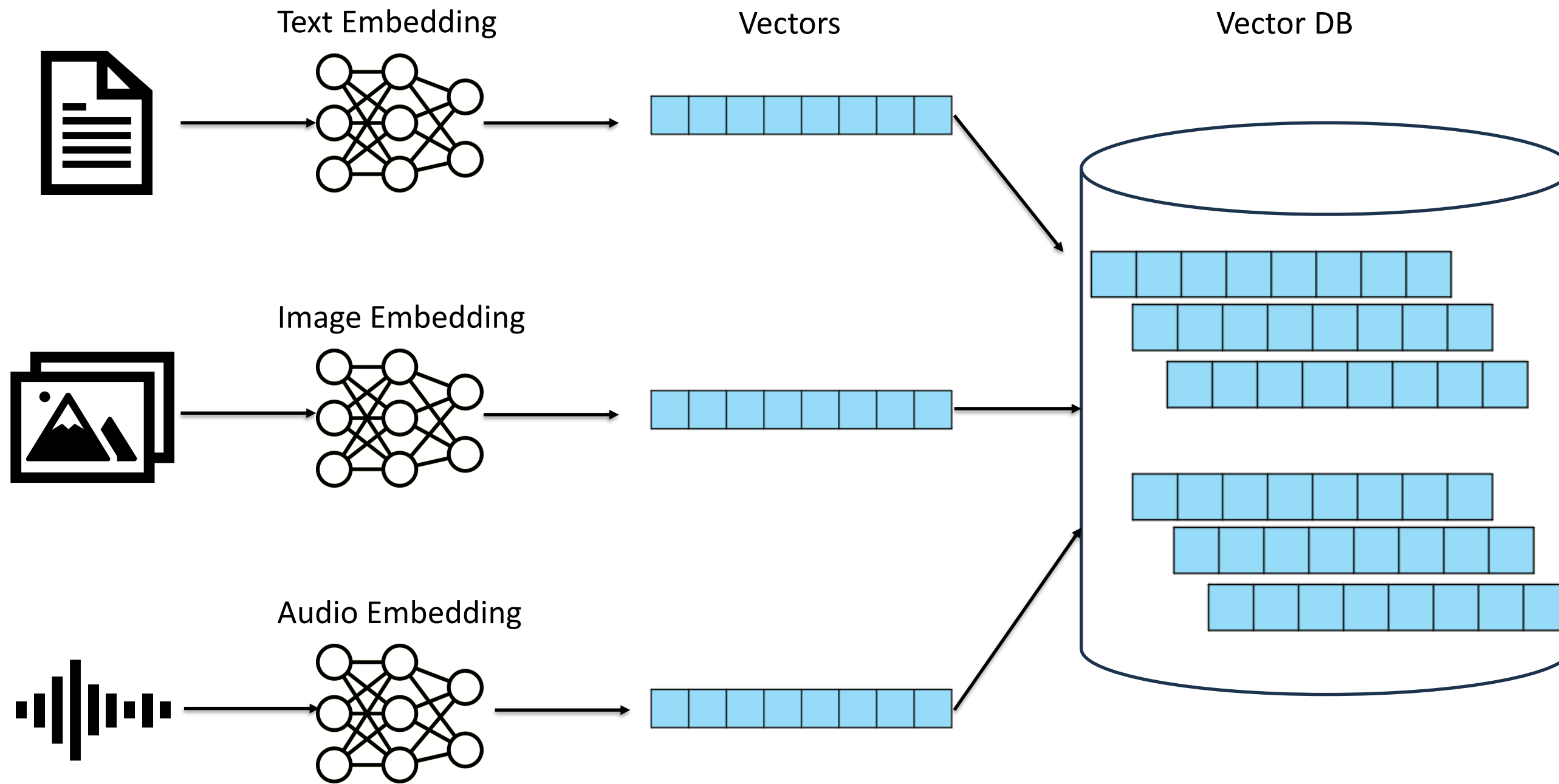
Vector DB

Key → Value

Key → Value

Key → Value

# Monomodal RAG (text only)

What action movie do you recommend?

Embedding

[[1 2 … 12 9],
[6 2 … 1 9],
[7 5 … 10 9],
[1 0 … 0 7],
[2 9 … 7 8]]

Cosine Similarity

| Embedding | Text |
|-----------|------|
| [1 4 … 4], [...] | The movie starts... |
| [2 3 … 9], [...] | Tom Cruise is... |

Embedding

Text

Vector DB

Retriever

Mission Impossible!

Large Language Model

{context}

Top K results

**Humans communicate in multiple ways**

# Introducing Large Multimodal Models

# Example of Multimodality: VATT

# Multimodal RAG (text + images)

What action movie do you recommend?

Embedding

$$[[1\ 2\ \dots\ 12\ 9],$$
$$[6\ 2\ \dots\ 1\ 9],$$
$$[7\ 5\ \dots\ 10\ 9],$$
$$[1\ 0\ \dots\ 0\ 7],$$
$$[2\ 9\ \dots\ 7\ 8]]$$

Cosine Similarity

Same embedding space!!

| Embedding | Item |
|---|---|
| [1 4 … 4], […] | The movie starts… |
| [2 3 … 9], […] |  |

Text Embedding

Image Embedding

Text + Images

Vector DB

Retriever

Top K results

{context}

Large Language Model

Mission Impossible!

**Option 1: leverage an image embedding model**

Vision Transformer! E.g. OpenAI's CLIP

Multimodal embedding

[0 34 21 4 0 13 … 3]

# Attention is all you need – once more

# Breaking down images into patches



1D flattened patch

16x16

16x16 patch = 256 tokens

Transformer encoder

Positional embedding

Linear projection of 1D Flattened patches

# Option 2: Using an LMM to generate images embeddings

Large Multimodal Model
(E.g. OpenAI's GPT-4-
vision)

Text embedding (E.g.
OpenAI's *text-ada-002)*

Image
captioning

*An orange cat taking
a selfie with an
astonished
expression […]*

Text
embedding

| | "An orange cat […]" | [0 34 21 4 0 13 … 3] |

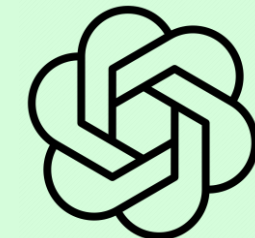**Demo Time!**

**Our Ingredients**

**MODELS**
GPT-4
GPT4-TURBO-VISION

**EMBEDDING**
CLIP
TEXT-EMBEDDING-ADA-002

**VECTORDB**
QDRANT

**ORCHESTRATOR**
LLAMA-INDEX
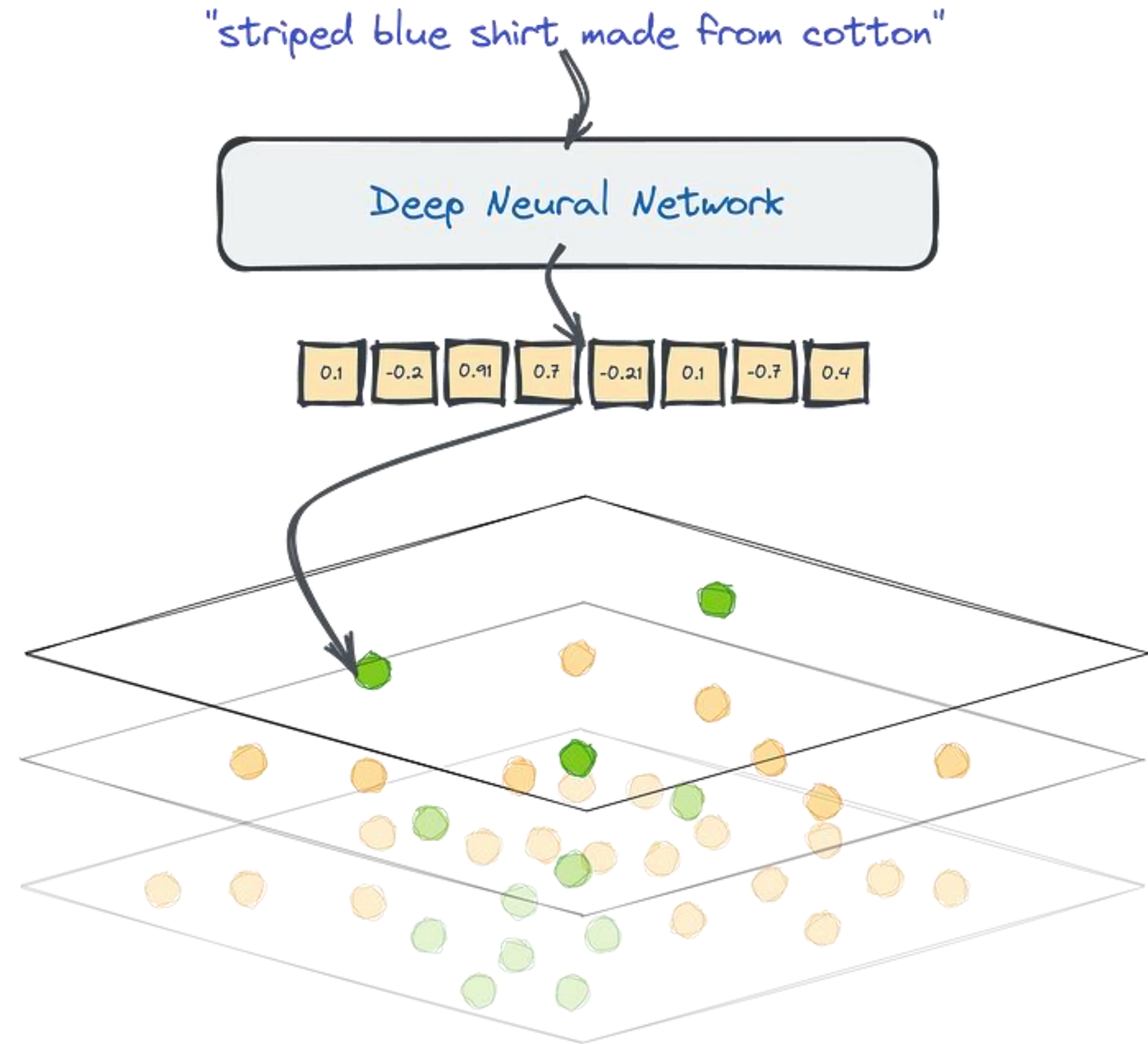
# Qdrant

🚨 **Challenge**→find similar documents in a big set of objects

💡 **Solution**→using a graph-like structure to find the closest object, so that we compute the distance for some candidates only.
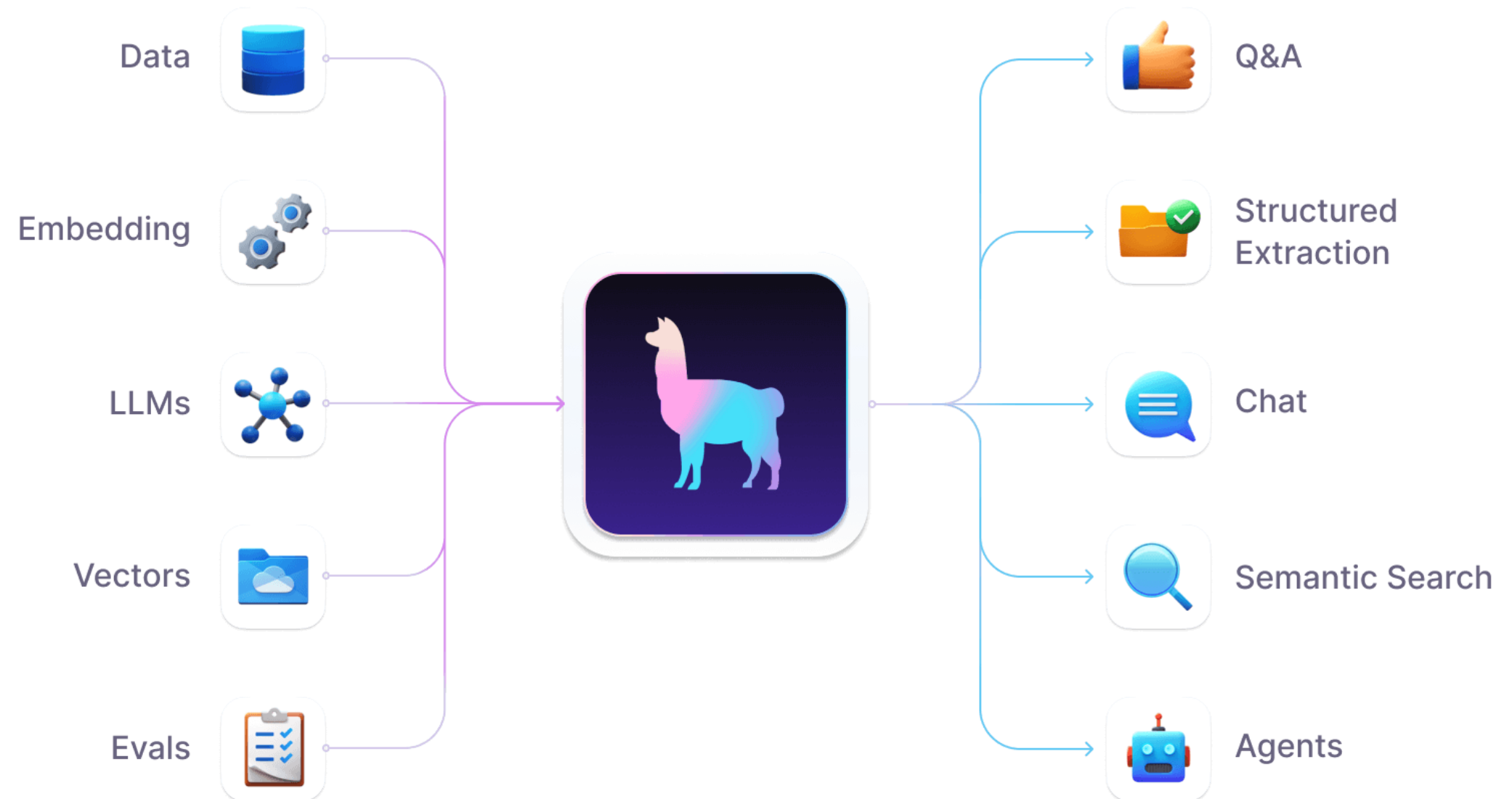


"striped blue shirt made from cotton"

Deep Neural Network

| 0.1 | -0.2 | 0.91 | 0.7 | -0.21 | 0.1 | -0.7 | 0.4 |

*Source: https://qdrant.tech/documentation/overview/vector-search/*
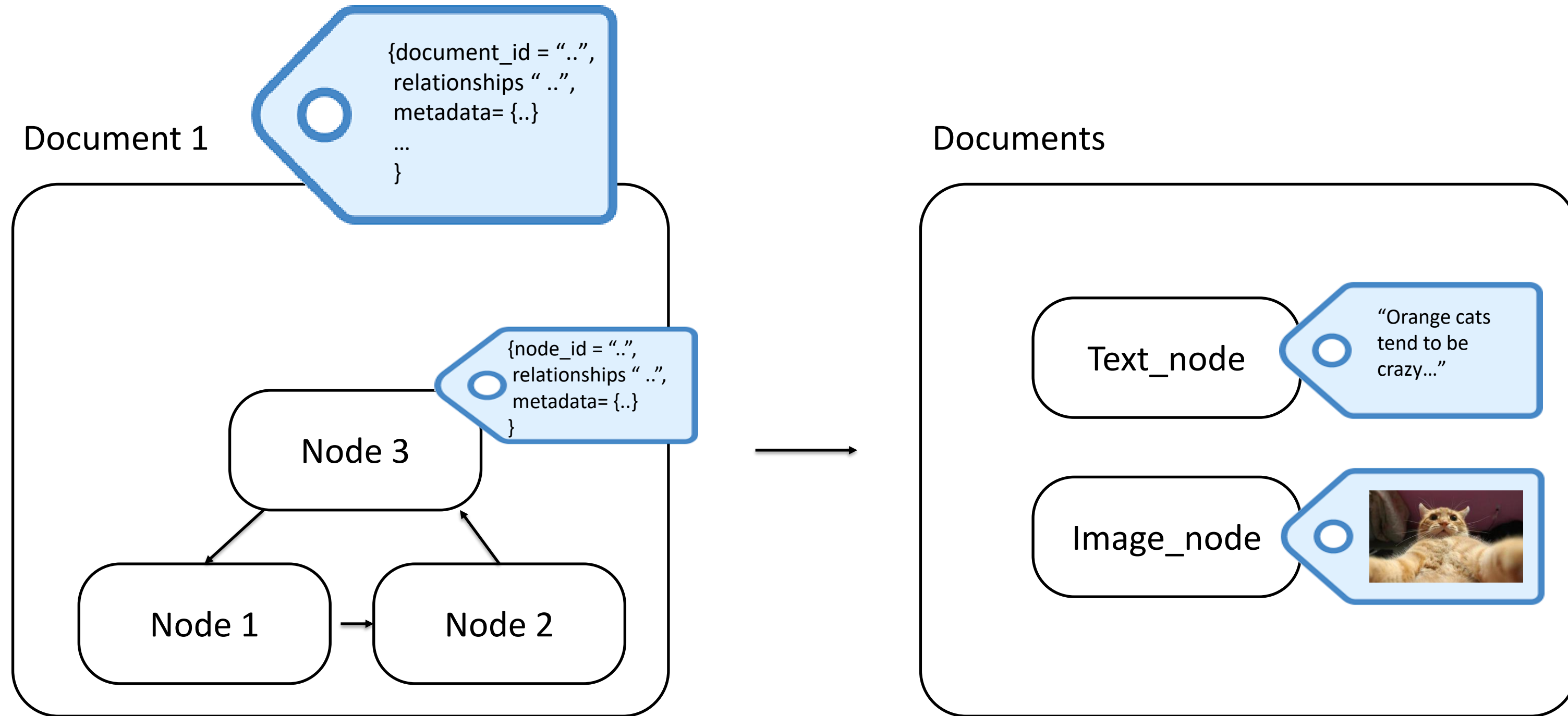
# Llama-Index

🔌 **Data Connectors** to ingest data from various sources

📄 **Data Indexes** to structure your data in intermediate representations

🤖 **Engines to interact** with your data in natural language

Data

Embedding

LLMs

Vectors

Evals

Q&A

Structured Extraction

Chat

Semantic Search

Agents

*Source: https://www.llamaindex.ai/*

# Nodes and Documents are first citizens in Llama-index

Document 1

{document_id = "..",
 relationships " ..",
 metadata= {..}
 …
 }

{node_id = "..",
 relationships " ..",
 metadata= {..}
 }

Node 3

Node 1

Node 2

Documents

Text_node

"Orange cats tend to be crazy…"

Image_node

# THANK YOU!

Let's Keep in touch!