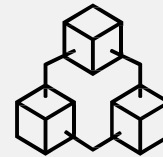# Modern Data Architectures

## *Vector Databases*

**HO GENT**

# Hype or the next big thing?

# Hype or the next big thing?

Hype?:

- "New kind of database for the AI era"
- For/by/focused on how AI models work
- Data → VectorDB → LLM

**swyx**
@swyx

$235m has been invested into Vector Databases in the past year:

- @qdrant_engine - $7.5m Seed

- @tryChroma - $18M Seed

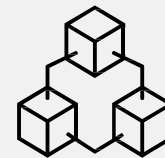- @weaviate_io - $50m Series A

- @milvusio - $60m Series B

- @Pinecone - $100m Series B

For reference, MongoDB raised $300m from start to $1.2b IPO.

Post vertalen

9:50 a.m. · 28 apr. 2023 · **81,3K** Weergaven

3

# The Theory

# The Theory

Put "simply": "store embedded data computed/retrieved using an *AI* model"

* Video:
  * https://www.youtube.com/watch?v=dN0lsF2cvm4
  * https://www.youtube.com/watch?v=ySus5ZS0b94

* Analogy in text: https://www.thdpth.com/p/the-vector-database-hype-explained
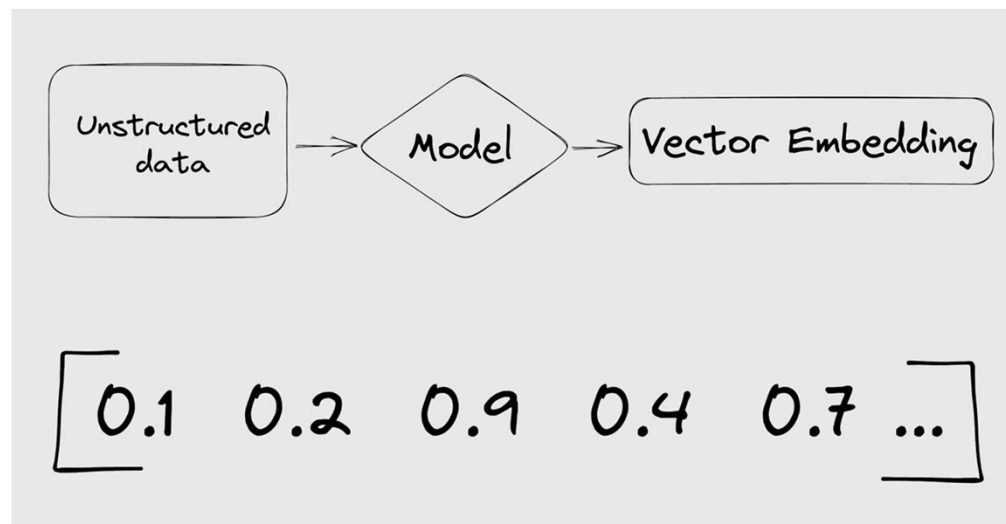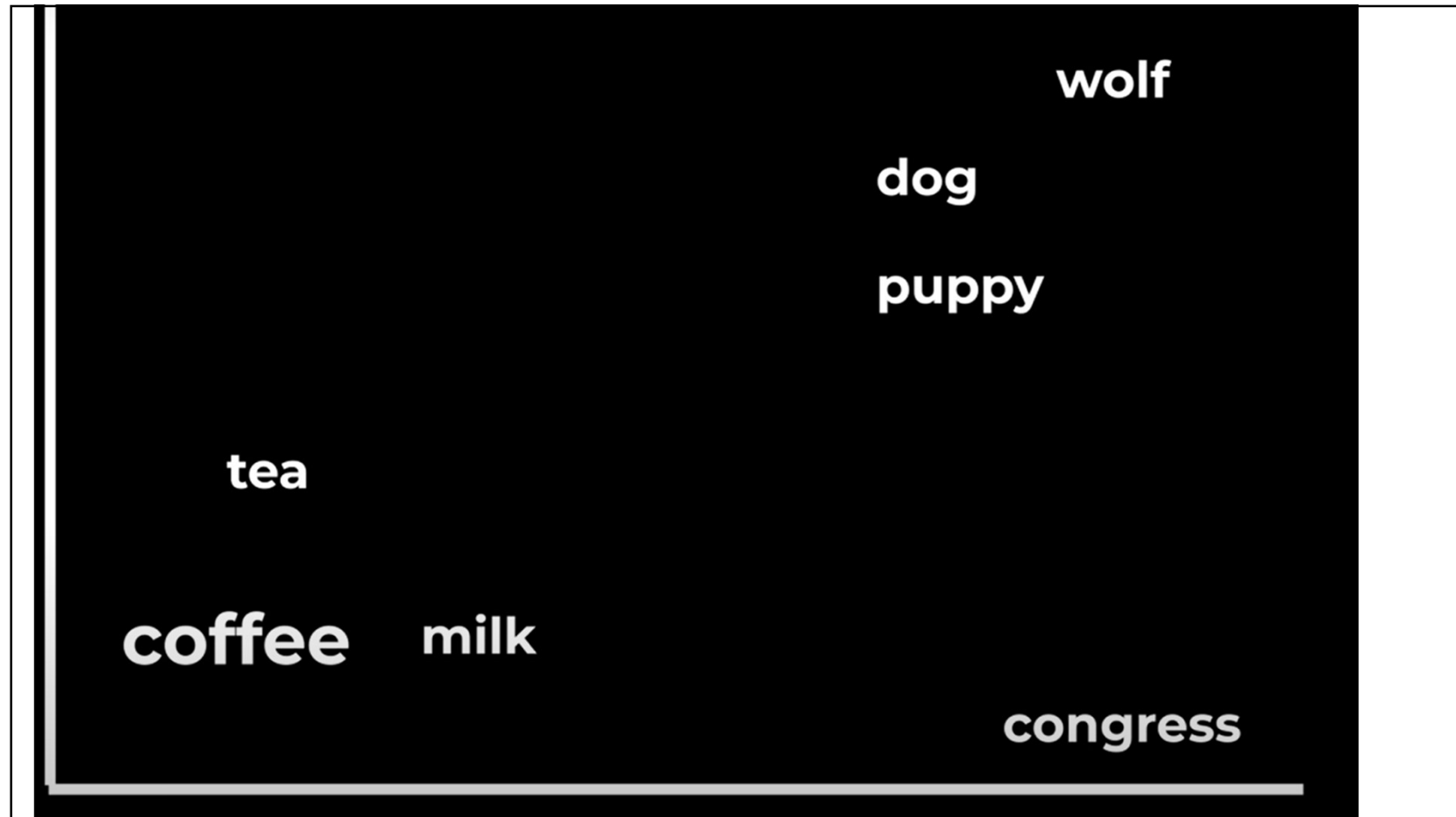
5

# The Theory

Why?
- Most data these days is unstructured and not easily "fittable"(= a good match) in a relational database:
  - Movies/video
  - Audio
  - Images
  - Social media data (= posts, tweets, …)

- For example:
  - The classic hello world example of AI: "Compare multiple pictures of a dog or a cat".
    - In a relational database → "color", "tags", "animal"
    - In non-relational → pixel values (?) or something else …

6

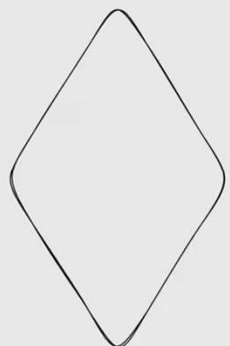# Vector embeddings ~ Vector databases

*Definition:*
"A vector database **indexes** and stores **vector embeddings** for **fast retrieval** and **similarity search**"
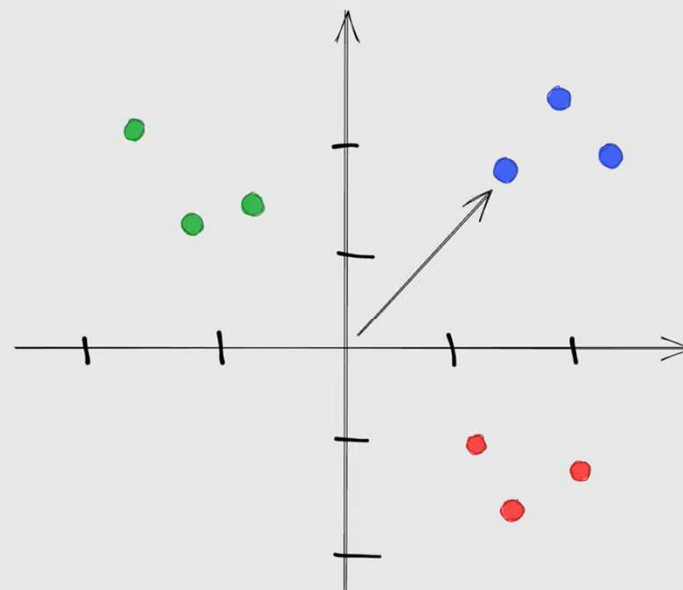


7

wolf

dog

puppy

tea

coffee  milk

congress

# vector embeddings (2D example)

king
man
woman

apple
banana
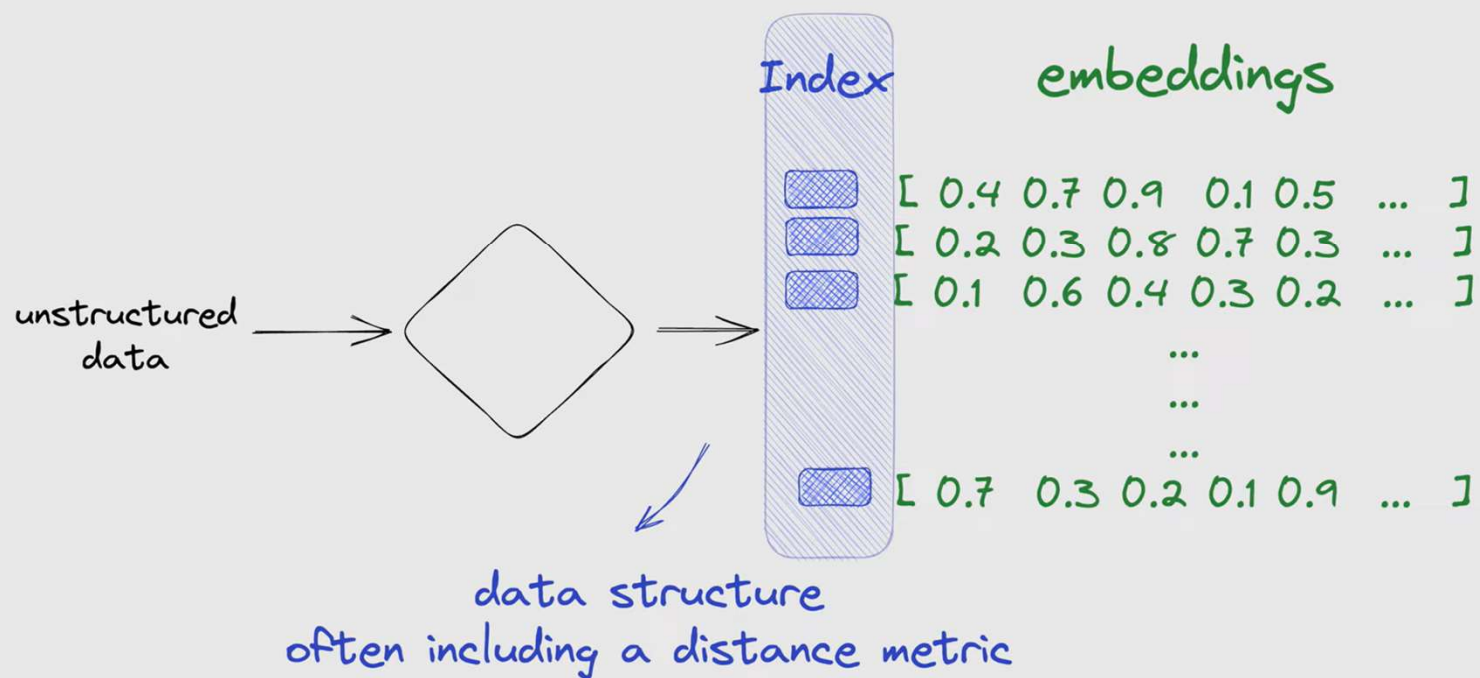orange

football
golf
tennis

[ 4   5   ]
[ 3   3.5 ]
[ 5   3   ]
[ 2.5  -2 ]
[ 2.5  -3 ]
[ 4   -2.5 ]
[ -3    4 ]
[ -1.5  3 ]
[ -2    2 ]
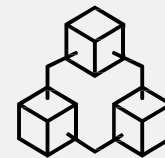
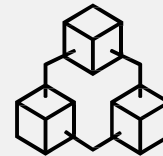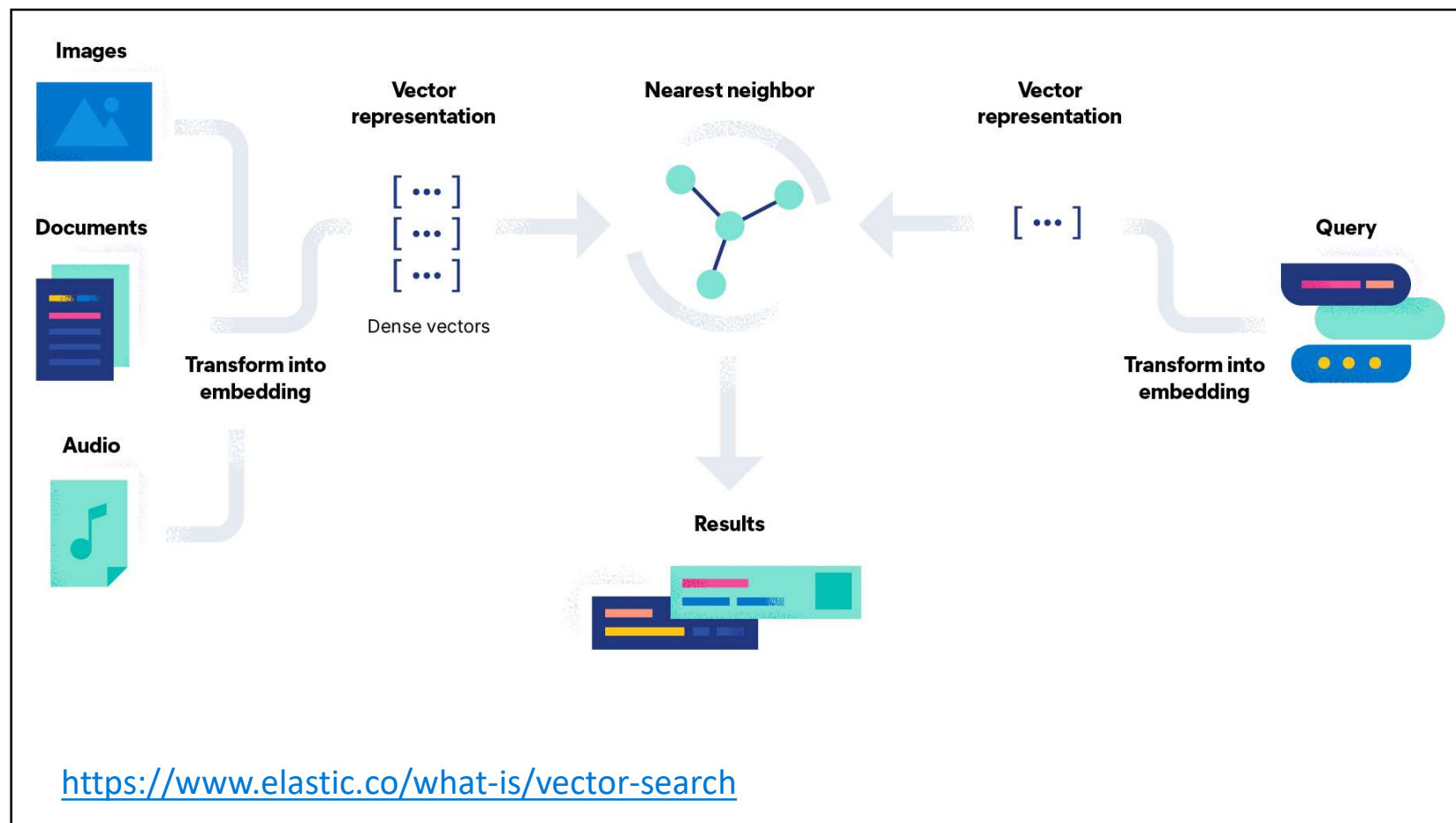$$d = \sqrt{(x2 - x1)^2 + (y2 - y1)^2}$$

9

# Use Cases

# Use Cases

- **Semantic search**: search based on the meaning or context of words and not the literal meaning or partial subset of words.

- Long term memory for LLMs.

- Similarity search for text, images, audio, video…

- Anomaly detection in datasets
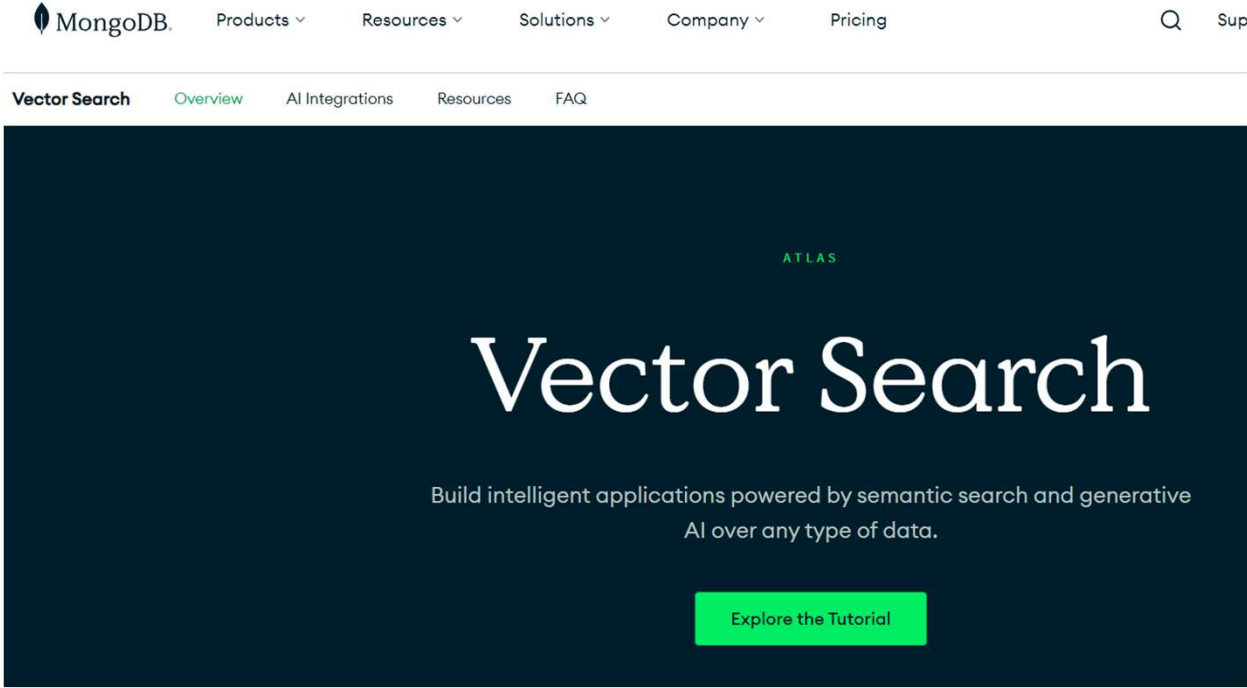
- Recommendation systems

# Different vector databases

# Vector Databases

- Pinecone
- Chroma
- Weaviate
- Qdrant
- Milvus
- Vespa
- SingleStore
- Redis
- Elastic Stack
- Mongo
- …
- (Every database that can store an n-array of numbers (?))

14

https://www.elastic.co/what-is/vector-search

# Vector Databases



https://www.mongodb.com/products/platform/atlas-vector-search

# Vector Databases

## Vector search

Query for data based on vector embeddings

This article gives you a good overview of how to perform vector search queries with Redis Stack. See the Redis as a vector database quick start guide for more information about Redis as a vector database. You can also find more detailed information about all the parameters in the vector reference documentation.

A vector search query on a vector field allows you to find all vectors in a vector space that are close to a given vector. You can query for the k-nearest neighbors or vectors within a given radius.
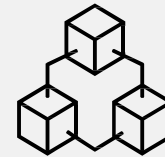
The examples in this article use a schema with the following fields:

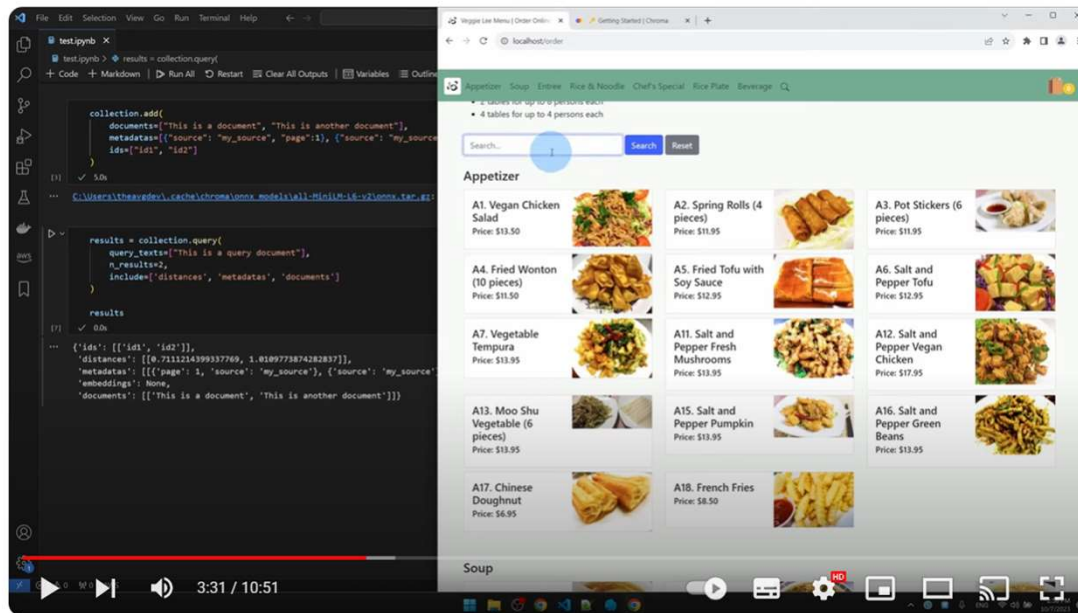| JSON field | Field alias | Field type | Description |
|---|---|---|---|
| $.description | description | TEXT | The description of a bicycle as unstructured text |
| $.description_embeddings | vector | VECTOR | The vector that a machine learning model derived from the description text |

## K-neareast neighbours (KNN)

The Redis command FT.SEARCH takes the index name, the query string, and additional query parameters as arguments. You need to pass the number of nearest neighbors, the vector field name

https://redis.io/docs/interact/search-and-query/query/vector-search/

17

# Demo use case

18

# Demo: Use case



Getting Started with ChromaDB - Lowest Learning Curve Vector Database & Semantic Search

https://www.youtube.com/watch?v=QSW2L8dkaZk

19