

O'REILLY®

Data Querying

Bert Gollnick





Learning Objectives

By the end of this module, you will:

- Know how to pull information from the database



Querying Introduction

Input Prompt

What did the fox do
with the dog?

Output Prompt

The fox jumps over
the lazy dog.

Embedding 1010
Model 1010



The fox jumps over
the lazy dog.

1010
1010

[0.5, 0.3, ..., 0.2]



{„author“: „...“, ...}

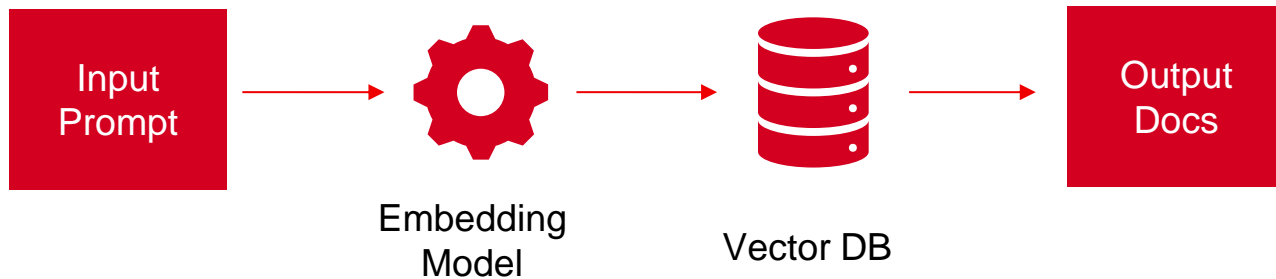


unique id1

Vector Database



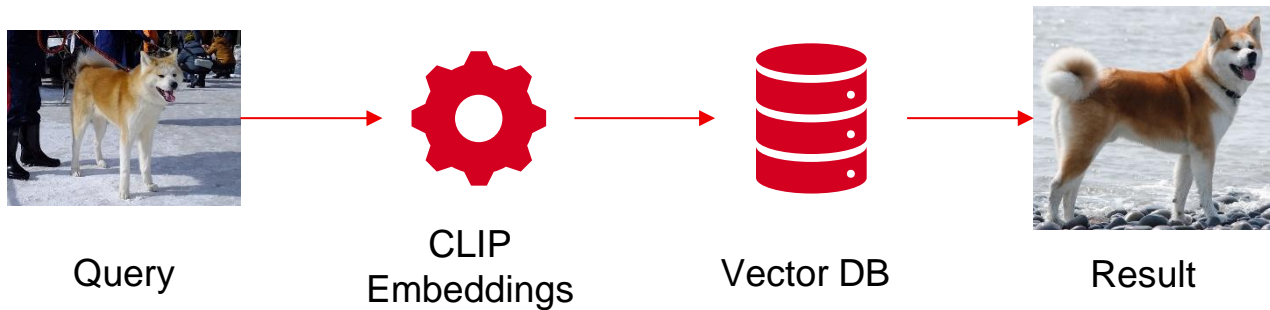
Text Querying



Practical Implementation

```
collection.query(query_texts=["This is my input text"])
```

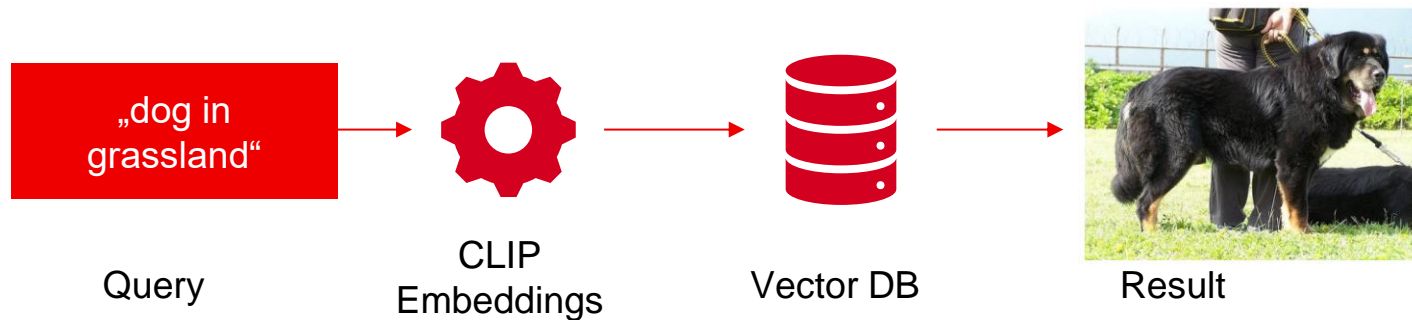
Image Querying



```
query_list = ["../data/dogs/akita_1.jpg"]
query_result = chroma_collection.query(
    query_images = query_list,
    n_results=3,
    include=['documents',
             'distances',
             'metadatas', 'data',
             'uris'],)
```

Result 1: ../data/dogs/akita_3.jpg
with distance: 0.17

Image Querying



```
query_list = ["dog in grassland"]
query_result = chroma_collection.query(
    query_texts = query_list,
    n_results=3,
    include=['documents',
             'distances',
             'metadatas', 'data',
             'uris'],)
```

Query: dog in grassland Result 0:
../data/dogs/mastiff_1.jpg
with distance: 0.85

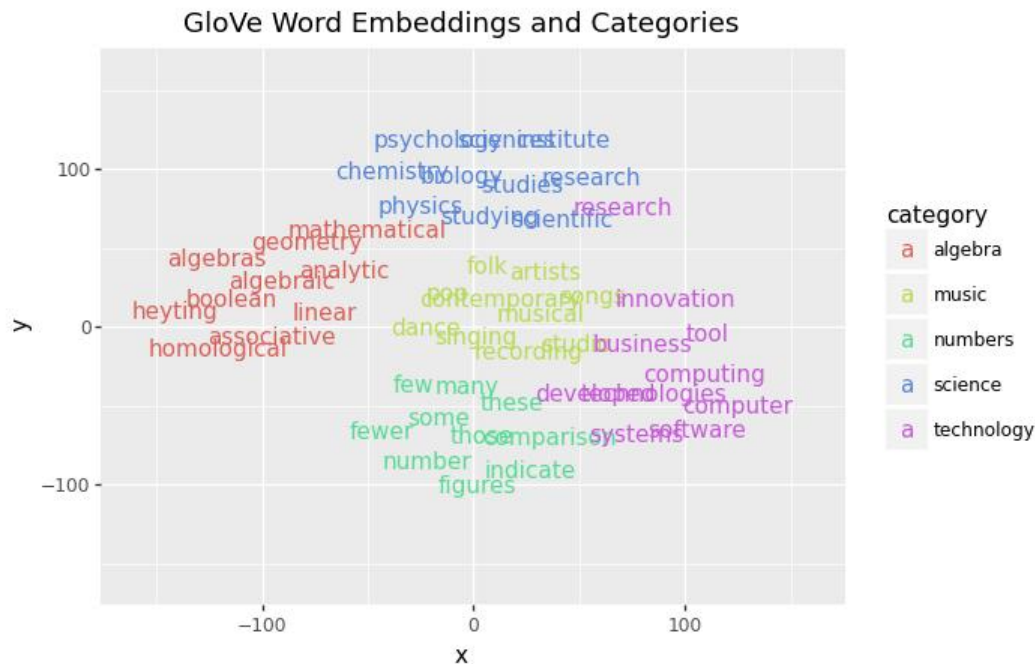


Similarity

- Vector DB needs to analyze similarity of query-embedding compared to document embeddings.
- Approaches:
 - Cosine Similarity
 - Maximum Margin Relevance



Similarity Search



$$dist = \sqrt{(x_1 - y_1)^2 + (x_n - y_n)^2}$$

For an embedding vector of 768 embeddings, there are 768 distance terms

Example: word embeddings reduced to 2 dimensions



Similarity Search

Imagename	Embedding					
dog1	0.3	0.02	0.8	0.6	...	0.4
dog2	0.1	0.52	0.7	0.6	...	0.4
...						
dogN	0.3	0.62	0.9	0.2	...	0.3

Vector Database



dogTest

0.3	0.02	0.8	0.6	...	0.4
-----	------	-----	-----	-----	-----

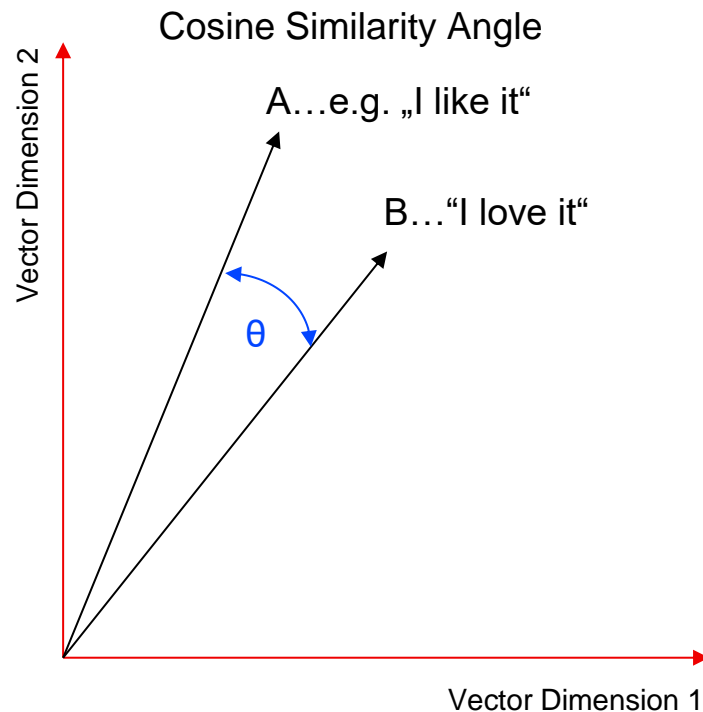
$$dist = \sqrt{\sum (x_i - y_i)^2}$$

For small data, go with np.array().

For large dataset not feasible!

Cosine Similarity

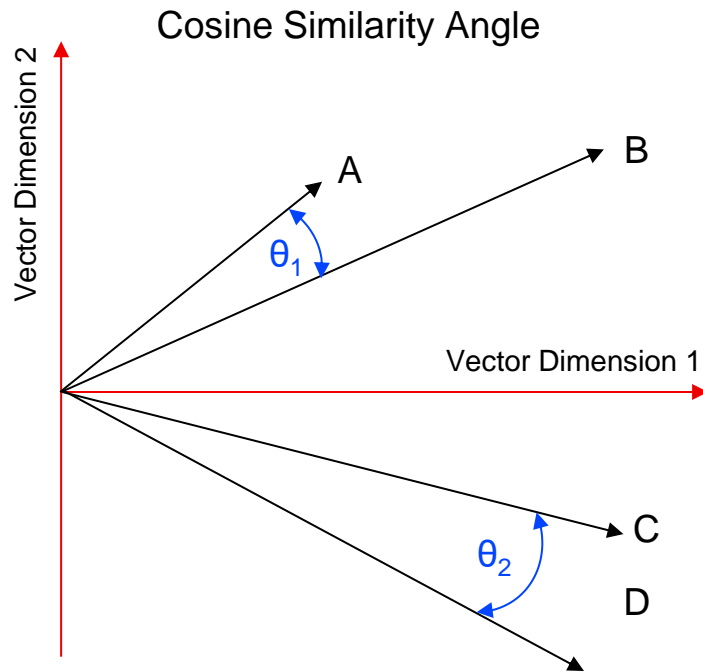
- Measures similarity between Embedding-Vectors based on angle θ .
 - Vectors maximally dissimilar
→ vectors perpendicular ($\theta = 90^\circ$)
 - Vectors completely similar
→ vectors parallel ($\theta = 0^\circ$)





Cosine Similarity

- Only the angle defines the similarity
- NOT the euclidean distance or magnitude of a vector
- Example
 - A: "The cat sleeps."
 - B: "The feline slumbers peacefully on the soft cushion."
 - C: "Trees grow leaves in spring."
 - D: "Fish swim in the ocean."



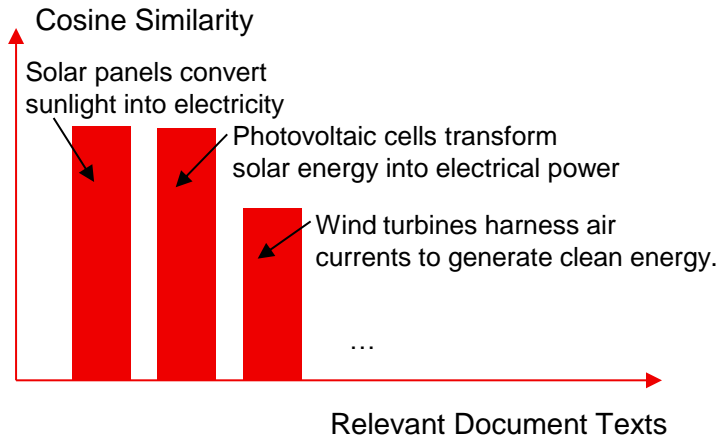
Maximum Margin Relevance

Topic: Renewable Energies

- Approach: reduce redundancy while maintaining relevance and diversity
- Redundancy...similar vectors
- Relevance...how closely do query and documents match
- Avoid clustering effect



What are the main types of renewable energy sources and how do they work?



The background is a gradient from red-orange on the left to yellow on the right. There are three large, semi-transparent circles of varying shades of orange and red. The text "O'REILLY" is centered in white, with a registered trademark symbol (®) at the end.

O'REILLY®