# Vector Databases

# Data Ingestion Pipeline:
# Introduction
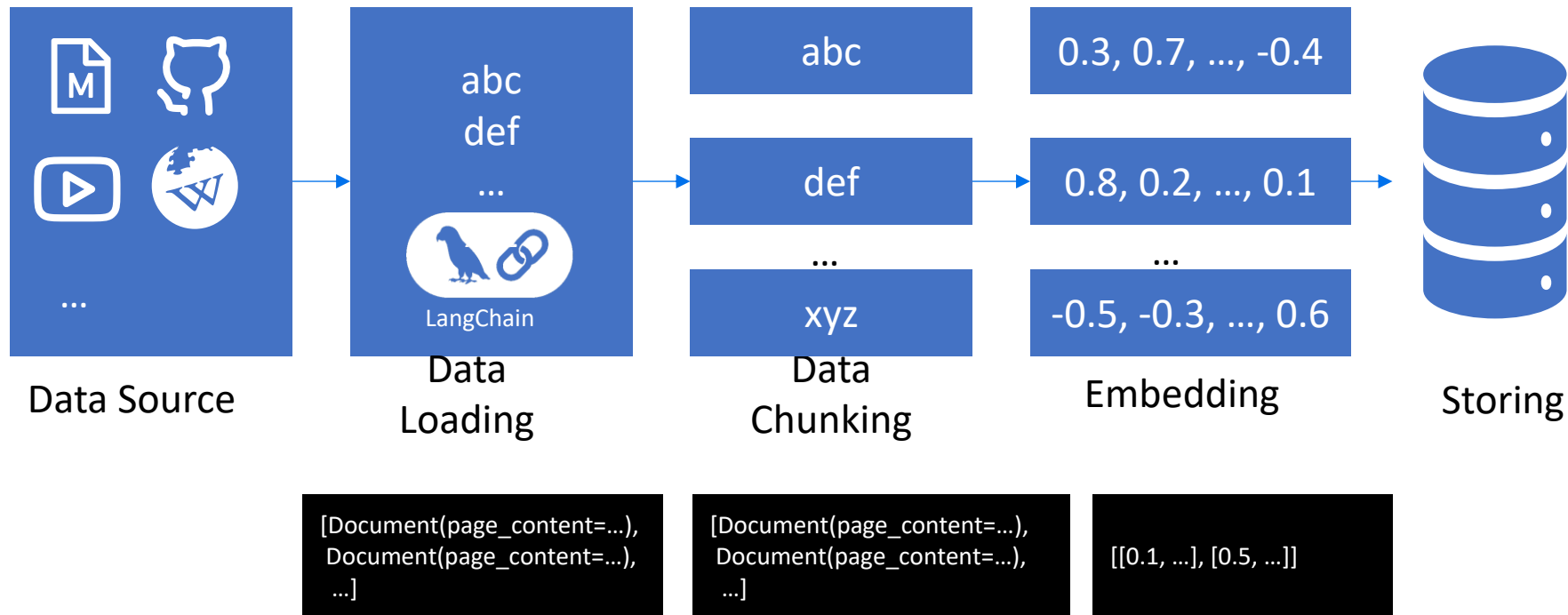
# Vector Database

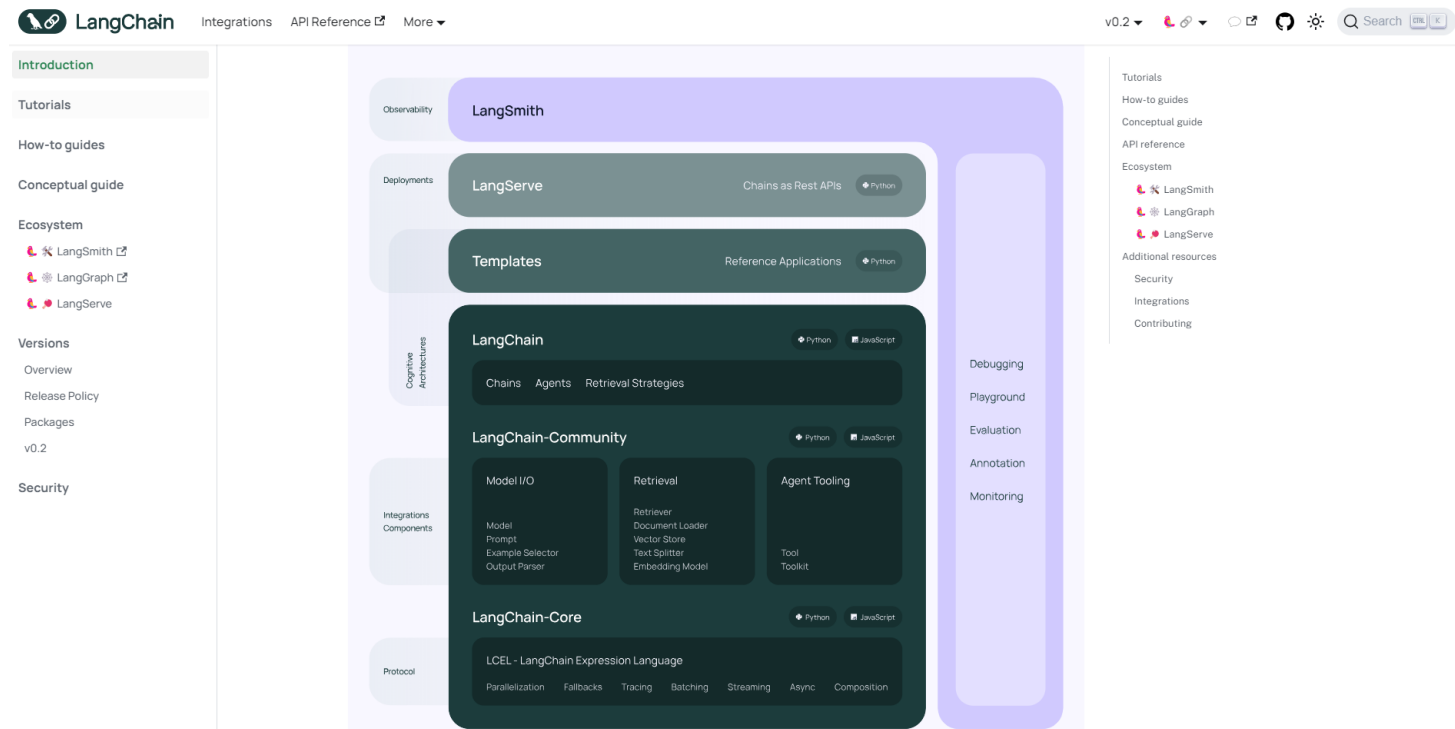Introduction



| Data Source | Data Loading | Data Chunking | Embedding | Storing |

# Vector Database

# Vector Database

## Additional Resources



Source: https://python.langchain.com/
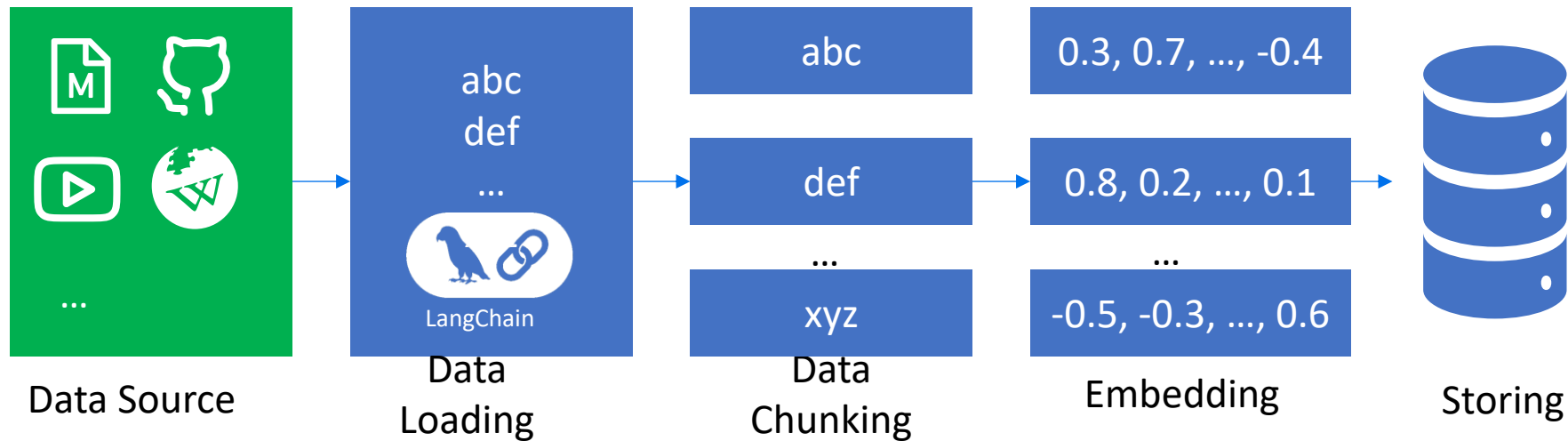
# Data Ingestion Pipeline:

# Data Source and -Loading

# Vector Database
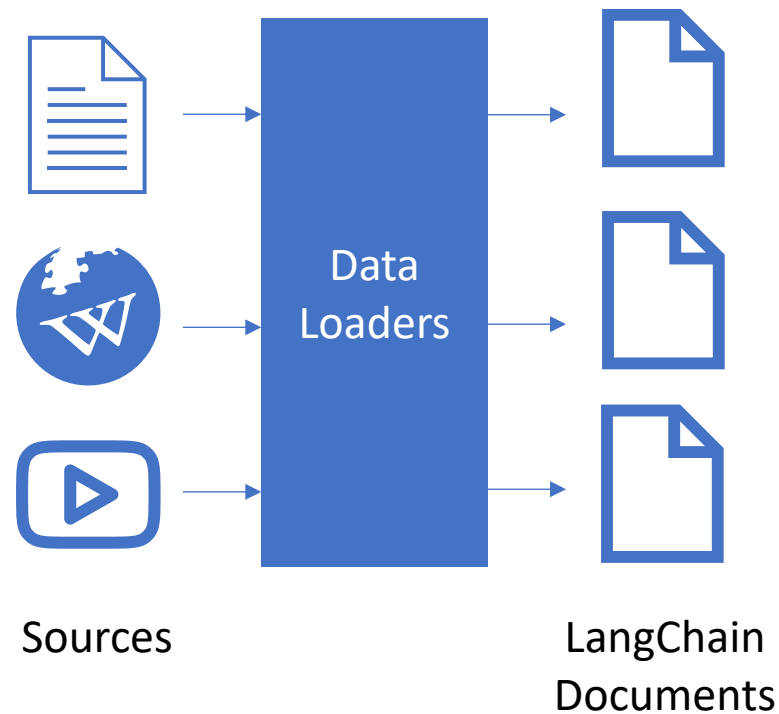
# Vector Database
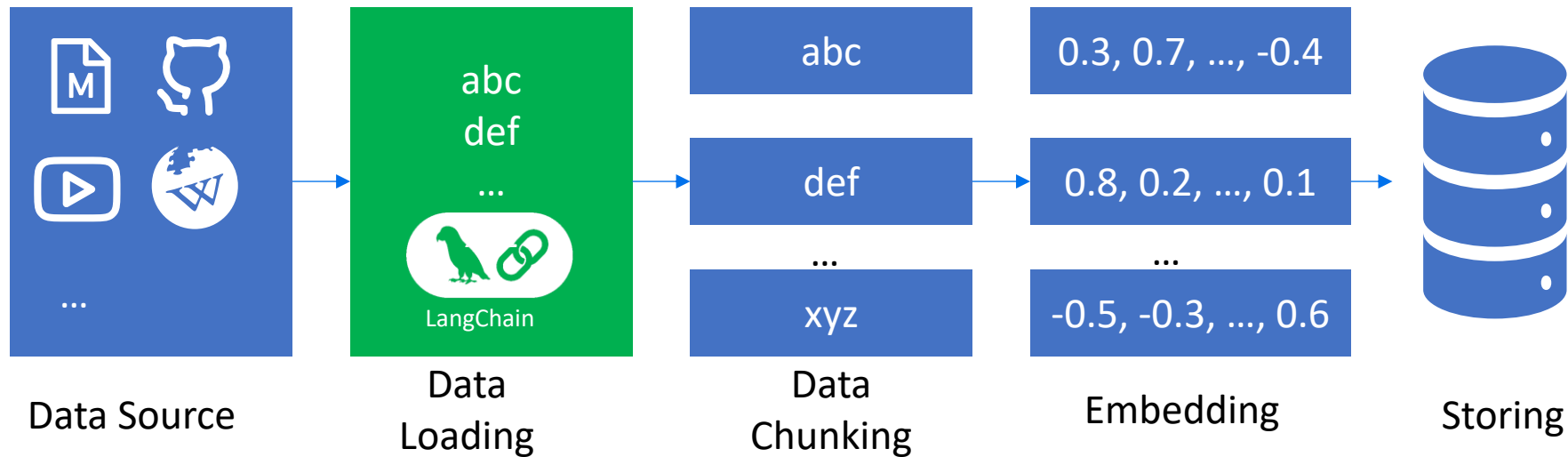
- Hundreds of different data sources are supported by LangChain

- DataLoader returns list of LangChain documents

- Documents have two attributes

  - Metadata

  - page_content



Sources                Data Loaders          LangChain Documents

# Data Ingestion Pipeline:

# Data Chunking
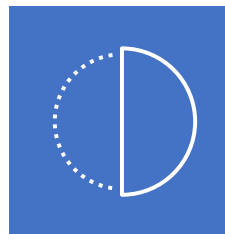
# Vector Database
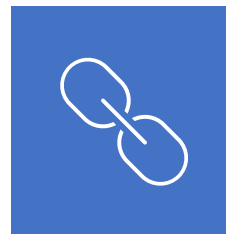
Data Source → Data Loading → Data Chunking → Embedding → Storing
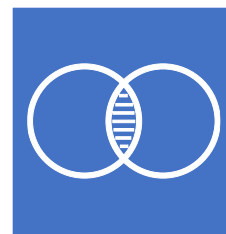
What is Data Chunking?

- Dividing larger pieces of information into smaller, manageable units

- These units called „chunks"

- Required to fit model context window

- Chunks should be:
  - Small
  - Semantically meaningful

Split    Combine    Overlap

# Vector Database

Overlap

This is my example

Chunk 1   Chunk 2

This is my

my example

Sentence 1                          Sentence 2

The cat chased the mouse. The mouse

Sentence 3

was afraid.   Trees provide oxigen for

us to breath.

- Sentence 1 and 2 are very similar
  → same chunk
- Sentence 3 different → new chunk
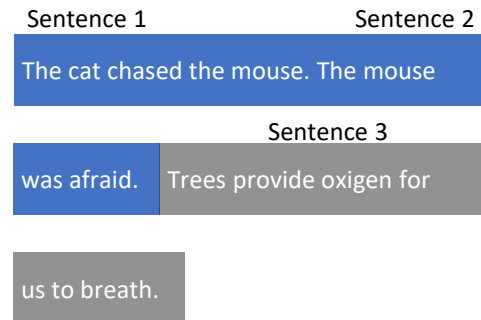
**Fixed Chunk-Sizes**
**Identical**
**pre-defined**

**Structure-Based Chunk-Sizes**
- e.g. chat messages should be consistent, no mix of users and chunks

**Semantic Chunking**
- based on semantic similarity
- e.g. when semantic break is observed

# Vector Database

Text

JSON

HTML

Code

# Vector Database

Data Chunking: Splitter Types

- chunk_size…defines maximum size of chunks [characters]
- chunk_overlap…possible overlap of max 5 characters

The quick brown fox jumps over the lazy dog.\n This is a simple example to show text splitting.\n.

RecursiveCharacterTextSplitter(
  chunk_size=20,
  chunk_overlap=5
  separators=["\n", " ", ""]
)

The quick brown fox

brown fox jumps

jumps over the lazy

the lazy dog.

This is a simple

simple example to

to show text

text splitting.

# Vector Database

Data Chunking: Best Practices for Chunk Sizes

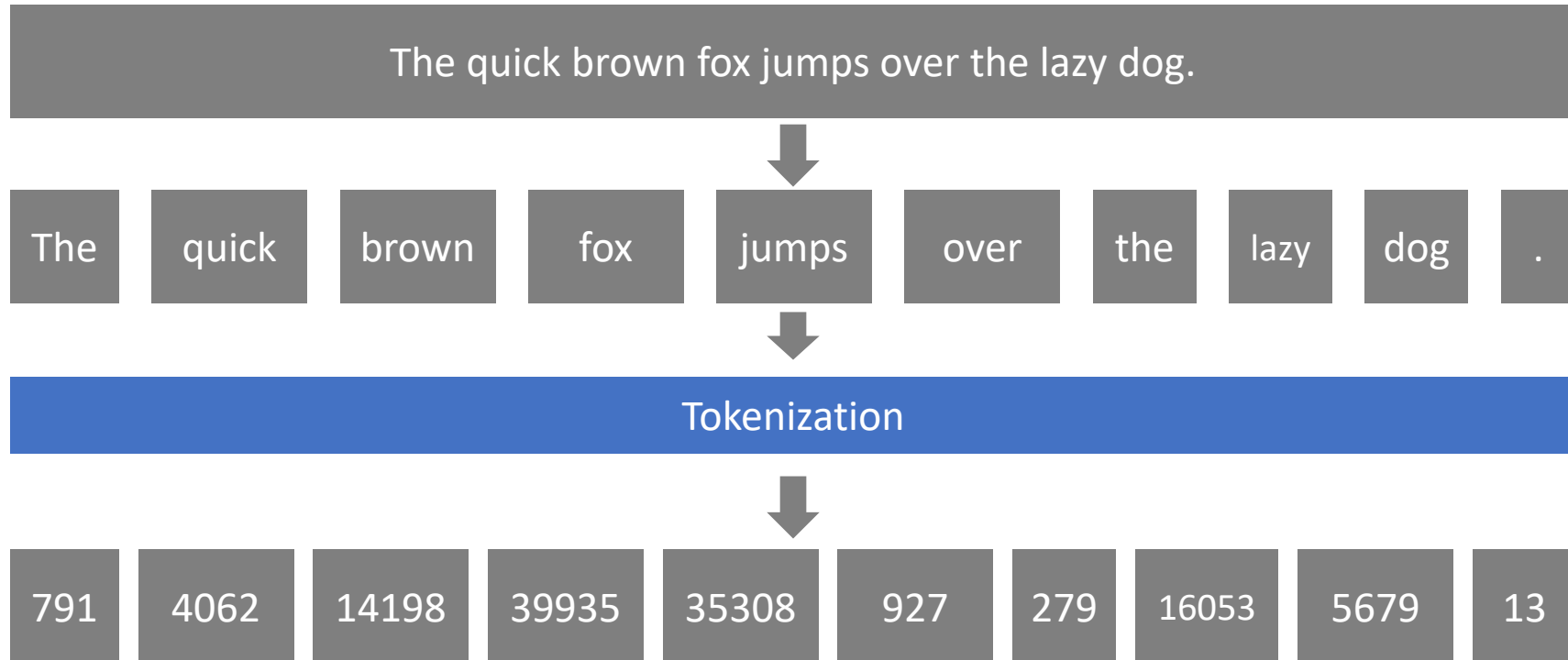| Texttyp | Granularität | Empfohlene Chunk-Größe (Tokens) | Empfohlener Overlap (Tokens) | Begründung |
|---|---|---|---|---|
| **FAQs / Kurze Q&A** | Hoch (Fein) | **128 - 256** Tokens | **10 - 30** Tokens | Ermöglicht hochpräzise Antworten; minimiert "Rauschen" (irrelevante Infos). |
| **Code-Snippets / Logs** | Hoch (Fein) | **200 - 400** Tokens | **50 - 80** Tokens | Kurz genug, um eine einzelne Funktion/Logik-Einheit abzubilden; Overlap sichert Kontext (z.B. Funktionssignatur). |

# Vector Database

| | | | | |
|---|---|---|---|---|
| **Technische Dokumente** (z.B. Specs, Ingenieursberechnungen) | Mittel | **500 - 800** Tokens | **100 - 150** Tokens | Genug Platz für Rechenschritte, Formeln oder einen kompletten technischen Absatz; Overlap bewahrt den Fluss. |
| **Allgemeine Artikel / Webseiten** (Standard) | Mittel | **512 - 1024** Tokens | **100 - 200** Tokens | Dient als **ausgewogener Startpunkt** für die meisten Anwendungsfälle; sichert einen umfassenden Kontext. |
| **Forschungspapiere / Verträge** (Langform, hohe Kohärenz) | Niedrig (Grob) | **1024 - 2000** Tokens | **200 - 400** Tokens | Erlaubt die Aufnahme längerer Argumentketten, gesamter Abschnitte oder juristischer Klauseln, um den *globalen* Kontext zu erhalten. |

# Vector Database

The quick brown fox jumps over the lazy dog.

| The | quick | brown | fox | jumps | over | the | lazy | dog | . |
|-----|-------|-------|-----|-------|------|-----|------|-----|---|

Tokenization

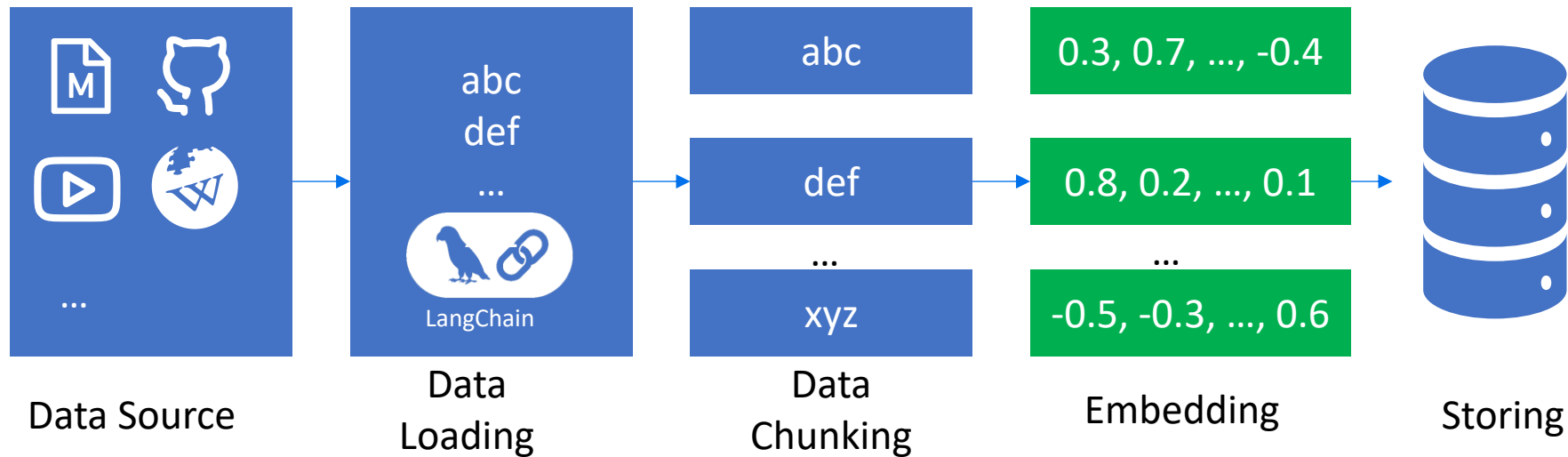| 791 | 4062 | 14198 | 39935 | 35308 | 927 | 279 | 16053 | 5679 | 13 |
|-----|------|-------|-------|-------|-----|-----|-------|------|----|

# Vector Database

- Embedding model works with tokens, NOT words

- Model can cover only specific sequence lengths

- Too long text (longer than context window) will be truncated

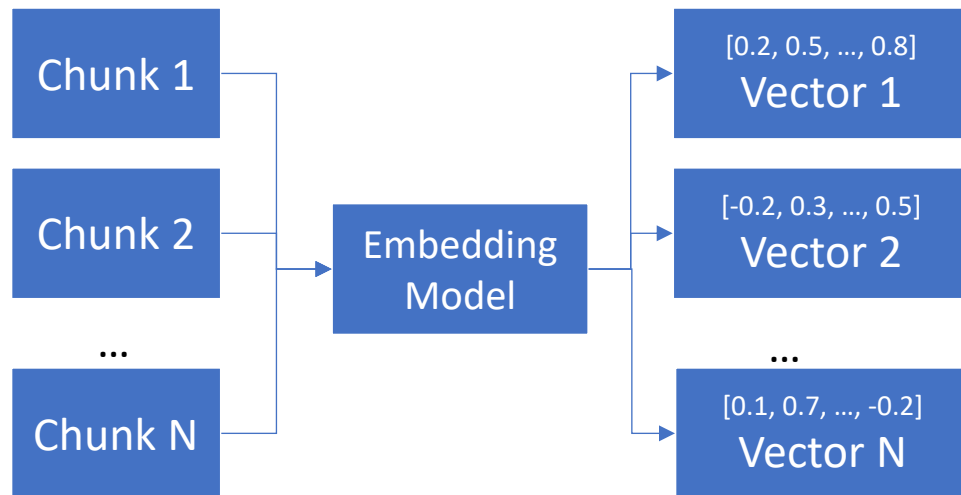# Data Ingestion Pipeline:

# Embeddings

# Vector Database

- Conversion of text data into numeric vectors

- Each word / sentence is represented as vectors

- Vector has „low" number of dimensions

| Chunk 1 |  | Embedding Model |  | [0.2, 0.5, …, 0.8] Vector 1 |
|---------|--|-----------------|--|-----------------------------|
| Chunk 2 |  |  |  | [-0.2, 0.3, …, 0.5] Vector 2 |
| … |  |  |  | … |
| Chunk N |  |  |  | [0.1, 0.7, …, -0.2] Vector N |

# Vector Database

Embeddings: What?



Source: https://www.youtube.com/watch?v=wjZofJX0v4M&t=814s

# Vector Database

- Convert words to numbers

- Representation of words as unique tensors in high-dimensional space

- Relationships to other words are captured

- Ideally similar words are close

- Usually Deep Learning applied to get embeddings
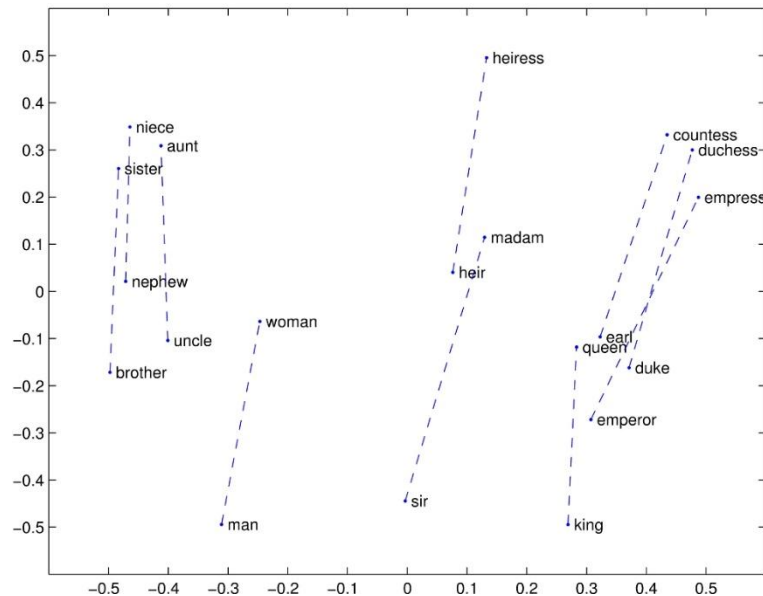
- Embeddings represent meaning



GloVe Word Embeddings and Categories

Word Embeddings represent words as **low-dimensional vectors** in mathematical space and **capture** their semantic and syntactic **meaning**.

# Vector Database

- Semantic representation
  - capture meaning of data
  - enable comparison and analysis

- Lower dimensionality
  - computational complexity is reduced
  - high-dimensional data can be represented in lower dimensions

- Reusability
  - usable across different applications



Source: https://nlp.stanford.edu/projects/glove/

# Vector Database

From Words to Tensors

Input sentence

Tokenization

Tensor

"It was a bright cold day"

$word_0$ — | It | 28 | → $[x_0$

$word_1$ — | was | ... | → $x_1$

$word_2$ — | a | | → $x_2$

$word_3$ — | bright | | → $x_3$

$word_4$ — | cold | | → $x_4$

$word_5$ — | day | | → $x_5$ ]

How?

# Vector Database

Word Embedding Approaches

One-Hot Encoding

Frequency-Based

Neural Network

# Vector Database

One-Hot Encoding

| Index: | 0 | 1 | 2 | 3 | 4 | 5 |
|--------|---|---|---|---|---|---|
| Word:  | It | was | a | bright | cold | day |

|        | 0 | 1 | 2 | 3 | 4 | 5 |
|--------|---|---|---|---|---|---|
| It     | 1 | 0 | 0 | 0 | 0 | 0 |
| was    | 0 | 1 | 0 | 0 | 0 | 0 |
| a      | 0 | 0 | 1 | 0 | 0 | 0 |
| bright | 0 | 0 | 0 | 1 | 0 | 0 |
| cold   | 0 | 0 | 0 | 0 | 1 | 0 |
| day    | 0 | 0 | 0 | 0 | 0 | 1 |

# Vector Database
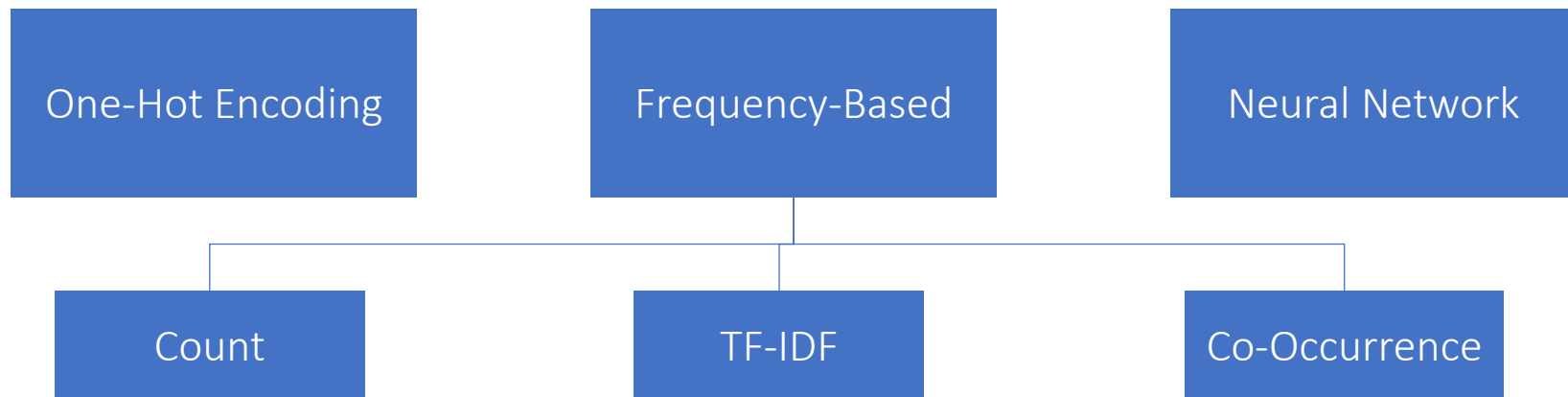
Problems

- Curse of dimensionality → memory issues
- Matrix very sparse

- Words are isolated from each other

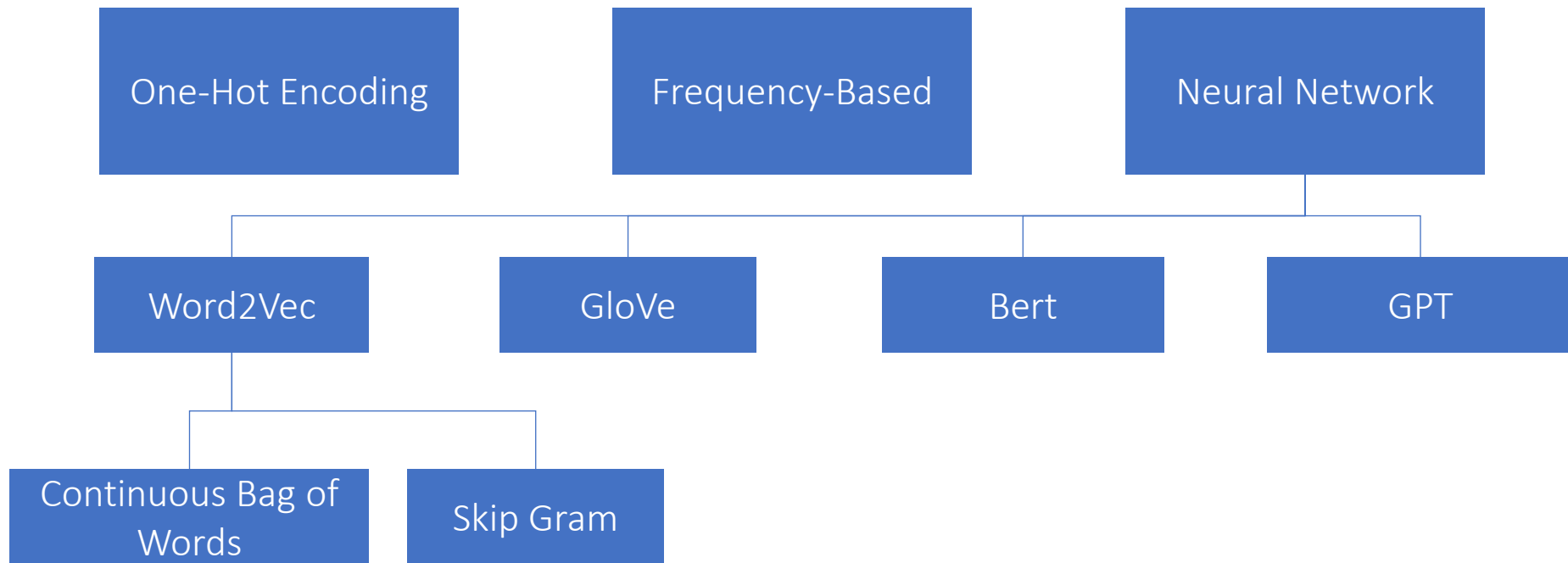- All words have the same distance to each other

# Vector Database

| One-Hot Encoding | Frequency-Based | Neural Network |
|---|---|---|

| Count | TF-IDF | Co-Occurrence |
|---|---|---|

- Very similar to OHE
- Gets count of words in document

- Term-Frequency/Inverse Term Freq.
- Gets count of words in document AND corpus
- Words frequent in a doc → important
- Words frequent in corpus → not important
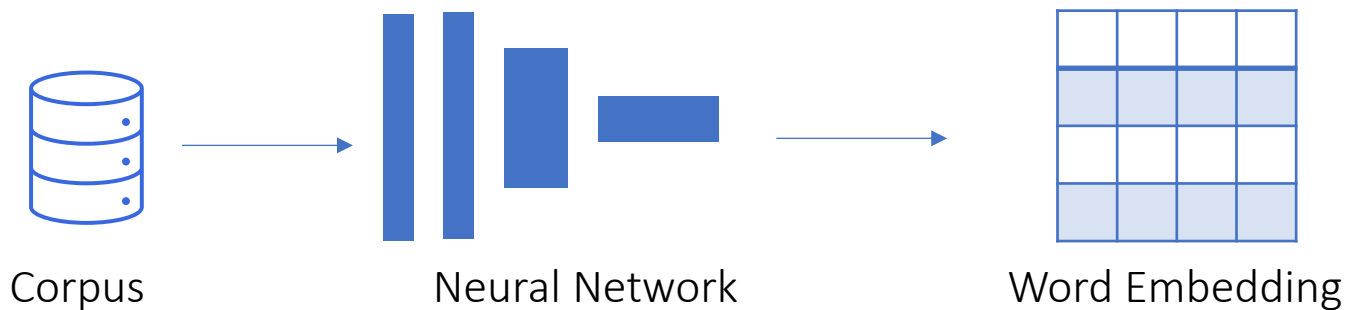
- Gets similarity of words

# Vector Database
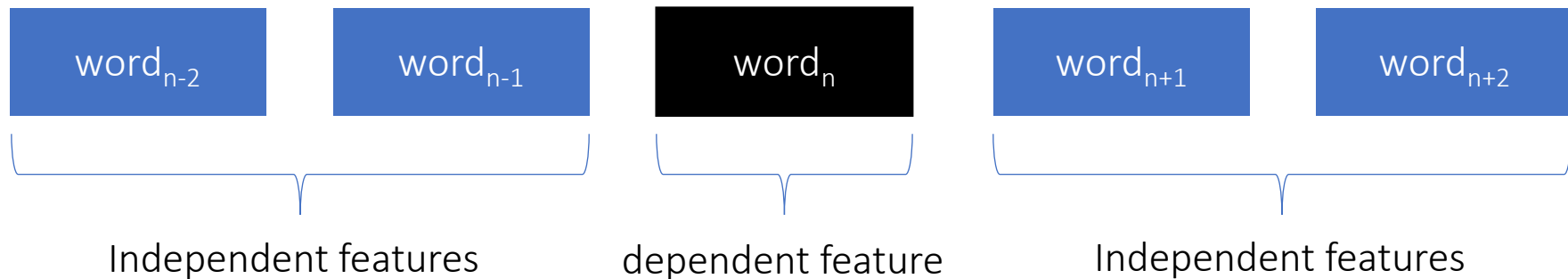
- Aim to
    - Capture context / meaning
    - Capture similarity to other words
    - Reduce dimension
    - Avoid memory issues
- Developed based on Neural Networks

Corpus       Neural Network       Word Embedding

# Vector Database

Word2Vec: Continuous Bag of Words

| word$_{n-2}$ | word$_{n-1}$ | word$_n$ | word$_{n+1}$ | word$_{n+2}$ |

Independent features      dependent feature      Independent features

| Independent Features | Dependent Feature" |
|---|---|
| ["It", "was", "bright", "cold"] | "a" |
| ["was", "a", "cold", "day"] | „bright" |
| ... | |

# Vector Database

Word2Vec: Skip Gram

Inputs

Model

Outputs

a → Embedding → Averaging → Linear → Softmax → It, was, bright, cold

# Vector Database

Word Embedding Approaches

One-Hot Encoding

Frequency-Based

Neural Network

Word2Vec

GloVe

Bert

GPT

Continuous Bag of Words
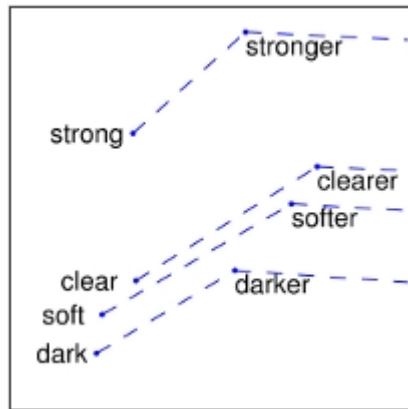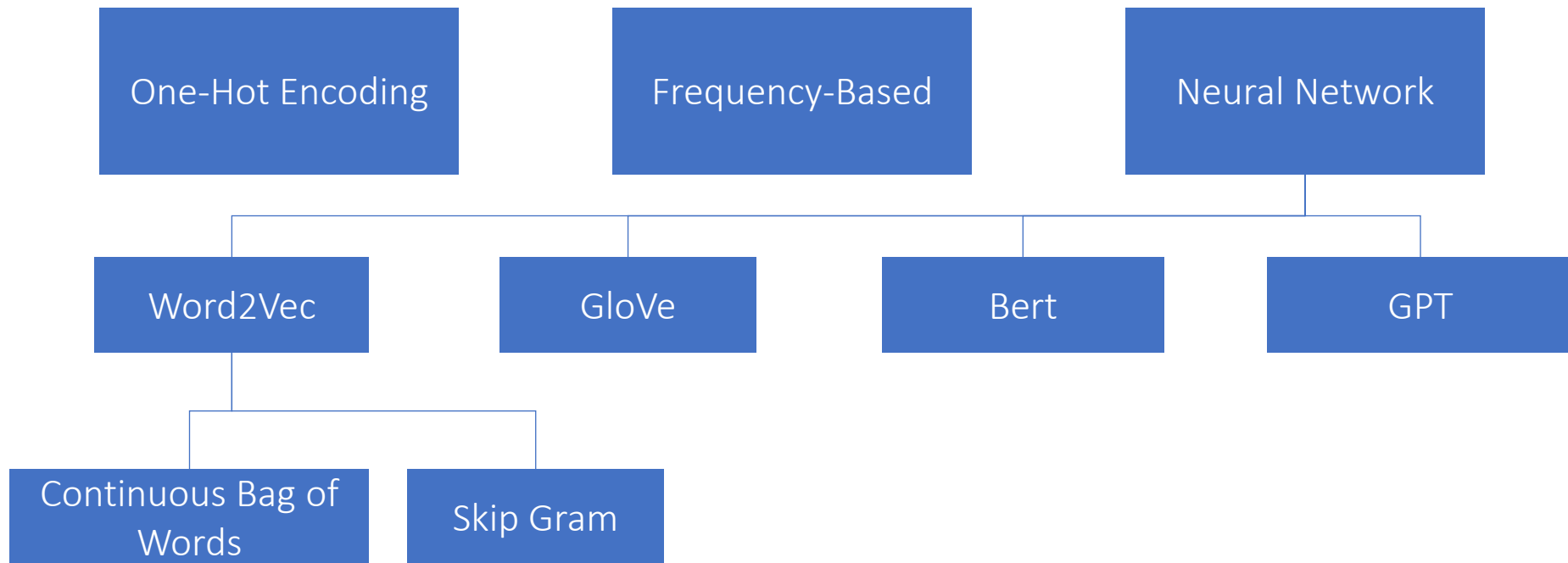
Skip Gram

# Vector Database

- Global Vectors for Word Representations
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation
- based on co-occurrence matrix of words in a corpus, which counts how often words appear together in the same context.
- constructs a matrix of word co-occurrence counts and then factorizes this matrix to obtain word embeddings
- factorization based on singular value decomposition (SVD)
- resulting embeddings are dense, low-dimensional vectors
- Encode words as vector of other words



Source: https://nlp.stanford.edu/projects/glove/
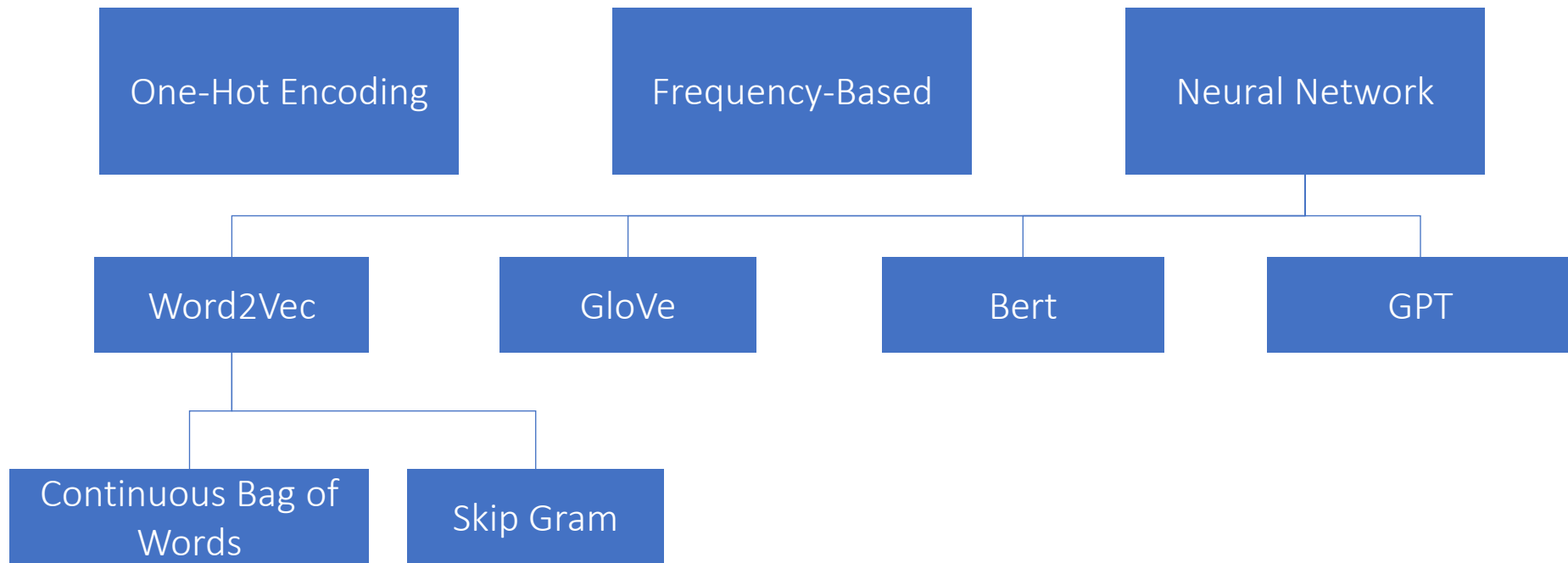
# Vector Database

- **B**idirectional **E**ncoder **R**epresentations from **T**ransformers
- Developed by Google in 2018
- Pre-trained word embedding
- Based on Transformers
- Applies „masked language modeling" – masking some words in sentence and learn to predict them
- Applies „next sentence prediction" – model predicts whether two sentences are similar in a text

- Original variants: BERT-base (110m parameters, 440MB) and BERT-large (340m parameters, 1.3GB)
- Other variants: RoBERTa, ALBERT, ELECTRA, …

# Vector Database

Word Embedding Approaches

# Vector Database

- **G**enerative **P**re-trained **T**ransformers
- Developed by OpenAI
- Not strictly a word embedding, but contextualized word embedding
- Unique embedding for each occurrence of a word based on surrounding words in text
- Applies Transformer architecture
- GPT-3 has 175 billion parameters

# Vector Database

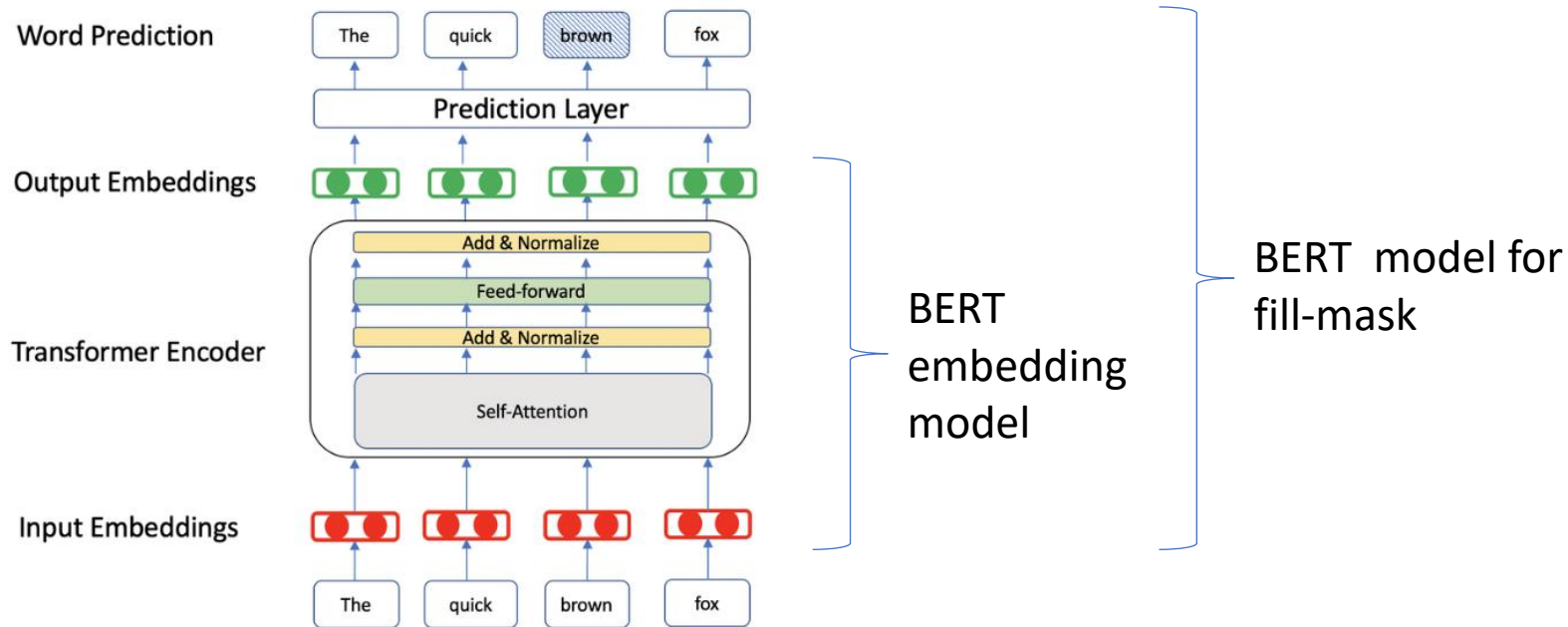| Parameter | Embedding Model | LLM |
|---|---|---|
| Base architecture | transformers | transformers |
| Process | texts, words, … $\rightarrow$ numerical vectors | predict next words |
| Target | find semantic similarities of texts | generate outputs depending on context, e.g. next words |
| Applications | semantic search, clustering, representations for ML | text generation, QA systems, chatbots, translations, code generation |

# Vector Database

Difference Embedding Model vs. Large-Language Model

| Parameter | Embedding Model | LLM, LMM |
|---|---|---|
| Inputs | Text, words, sentences, images, … | Text, words, sentences, images, … |
| Output | vector | human-readable text/code |
| Focus | representation of data | processing and generation of data |
| Model Size | smaller, more specific (narrow AI, e.g. sentence transformers) | larger, e.g. GTP, Llama, … |
| based on | pre-trained language models, uses architecture only for vector creation | uses transformers |

# Vector Database
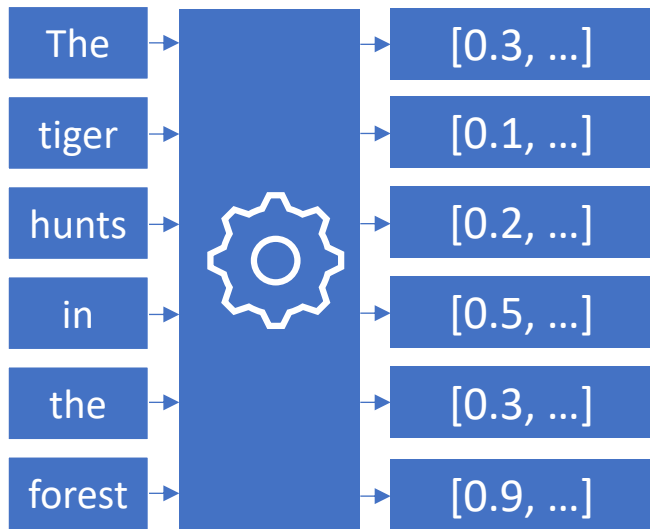
Difference Embedding Model vs. Large-Language Model

# Vector Database

Embeddings: How?

## Word Embeddings

| The | | [0.3, …] |
|-----|-----|-----|
| tiger | | [0.1, …] |
| hunts | ⚙ | [0.2, …] |
| in | | [0.5, …] |
| the | | [0.3, …] |
| forest | | [0.9, …] |

## Sentence Embeddings

The tiger hunts in the forest → ⚙ → [0.8, …]

# Vector Database

Embeddings: Which types are available?

| Type | Model | Provider | Price | Vector Size |
|------|-------|----------|-------|-------------|
| Online | text-embedding-3-small | OpenAI | 0.02$ / 1M tokens | 1536 |
| Online | text-embedding-3-large | OpenAI | 0.13$ / 1M tokens | 3072 |
| Online | mistral-embed | MistralAI | 0.10$ /1M tokens | 1024 |
| Offline | all-MiniLM-L6-v2 | Open Source | --- | 384 |
| Benchmark: https://huggingface.co/spaces/mteb/leaderboard … | | | | |

# Vector Database

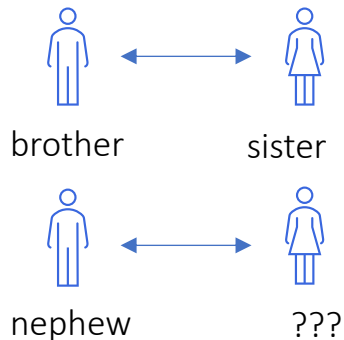Embeddings: Factors to consider

Price

Speed

Off-/Online

Benchmark Performance

# Vector Database

- Find closest words

```
get_closest_words_from_word('chess')
✓ 7.5s

[('chess', 0.0),
 ('backgammon', 4.379469394683838),
 ('grandmasters', 4.56368350982666),
 ('grandmaster', 4.613785743713379),
 ('scrabble', 4.677640438079834)]
```

- Find word analogies



brother    sister

nephew    ???

```
analogy = get_word_analogy(word1='sister',
                           word2='brother',
                           word3='nephew')
analogy
✓ 7.3s

'niece'
```
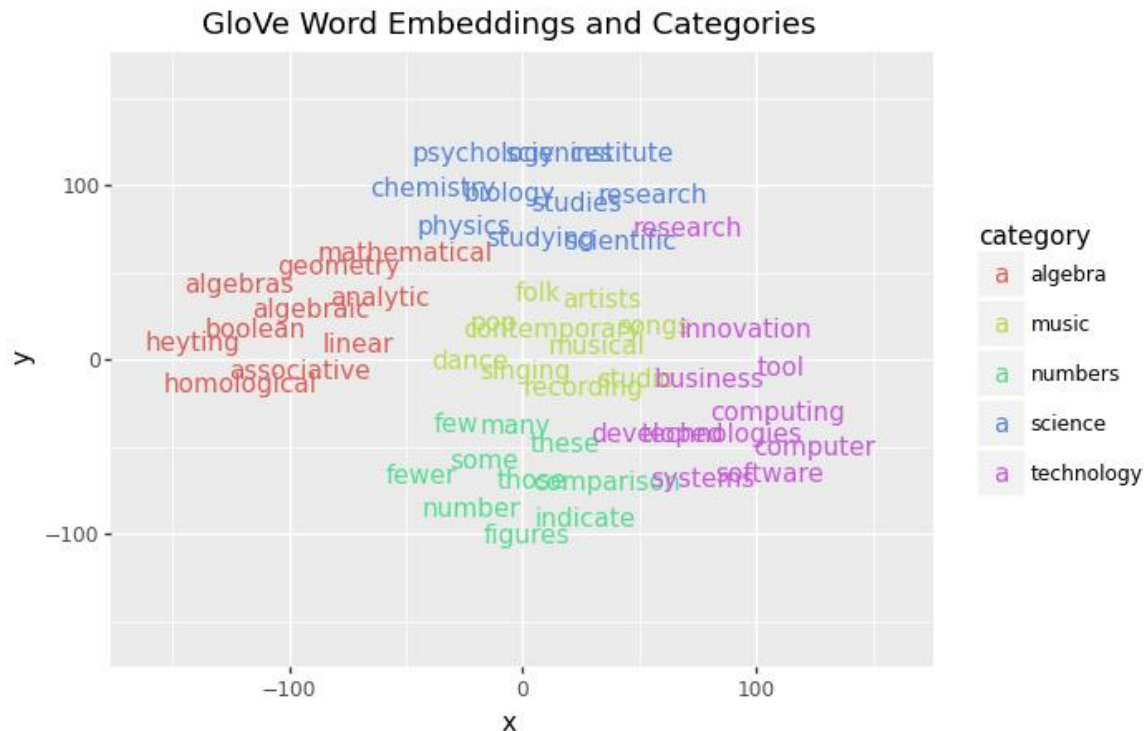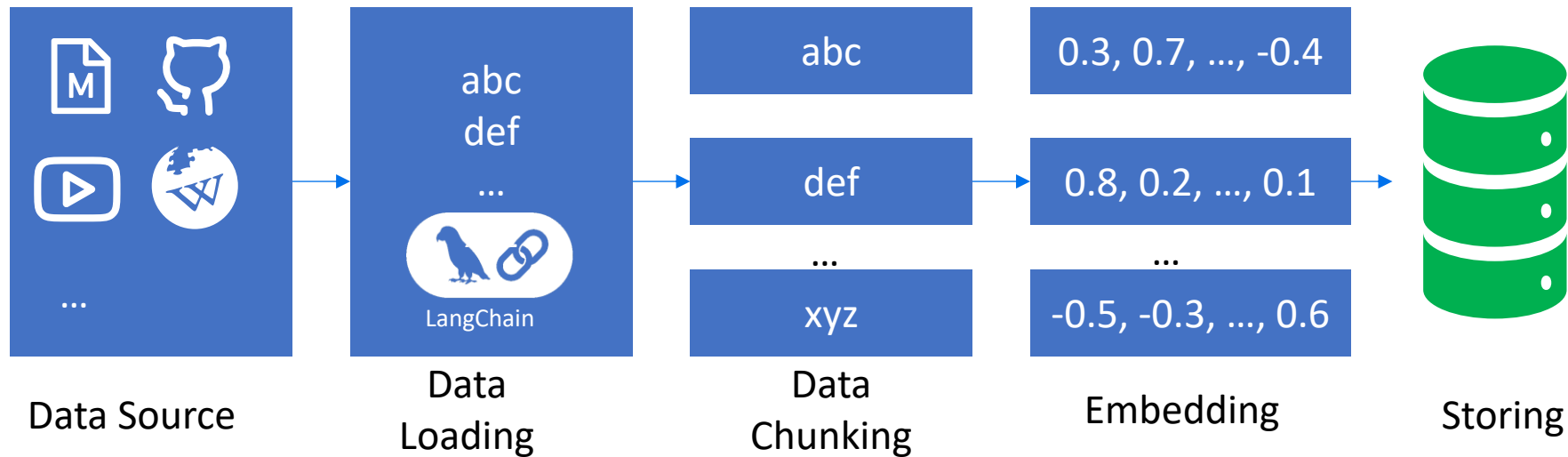
# Vector Database

- Given some categories
- Find words for the categories
- Check if they are „close" (similar)

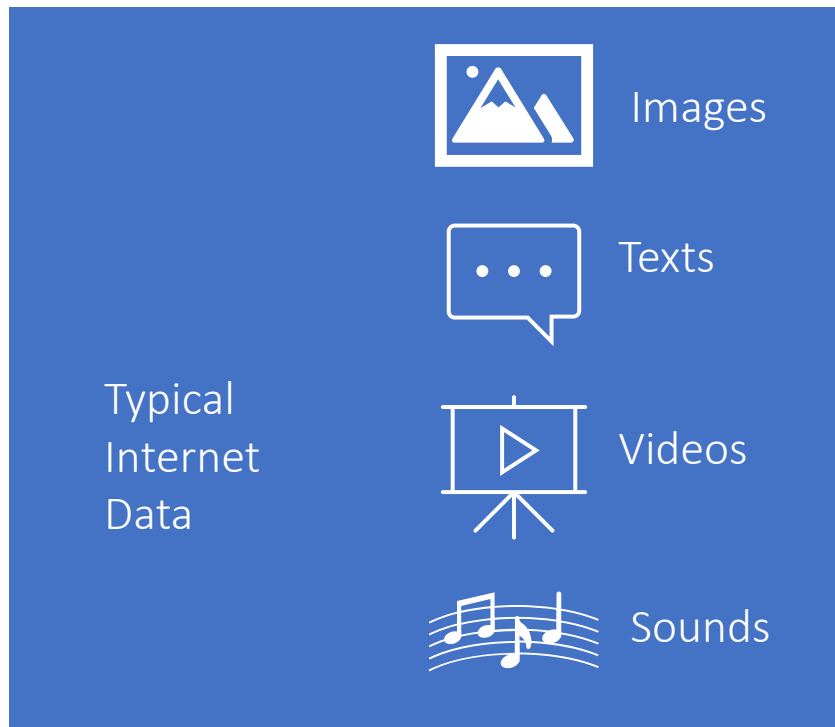# Data Ingestion Pipeline:

# Data Storing

# Vector Database

A vector database stores low-dimensional data (embeddings) for fast querying and similarity analysis.
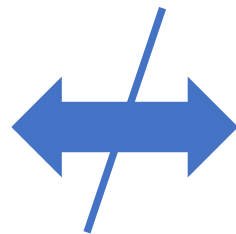
## Features

- Special type of database

- Allows to store, manage, and query data which is represented in geometric formats

- Enables similarity search, clustering, real-time analytics

# Vector Database

Data Storing: Why is a vector database needed?



Images

Texts

Videos

Sounds

Typical Internet Data

Unstructured Data

Typical Databases

SQL

Structured Data

# Vector Database

Source: https://www.datacamp.com/blog/the-top-5-vector-databases

# Data Ingestion Pipeline:

# Data Querying

# Vector Database

Input Prompt → Embedding Model → Vector DB → Output Docs

## Practical Implementation

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

Query  CLIP Embeddings  Vector DB  Result

```
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
```

# Vector Database

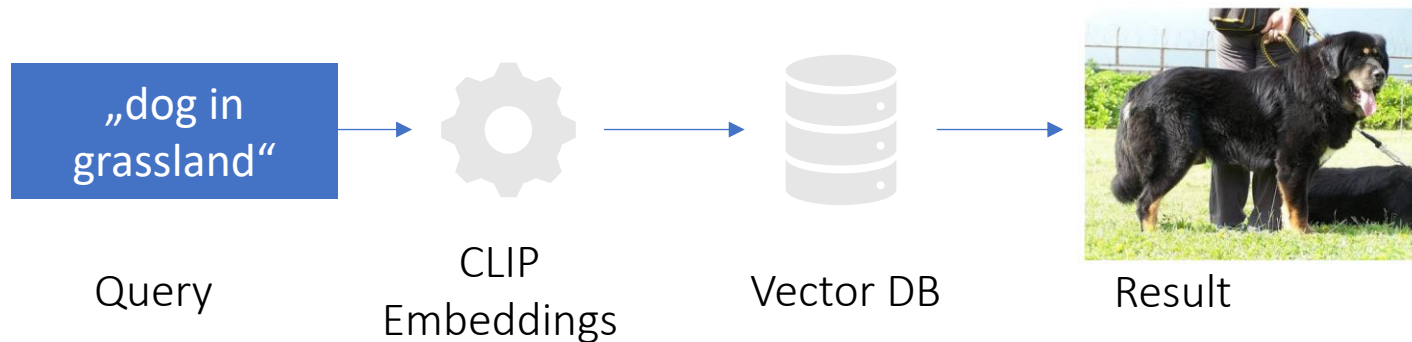| „dog in grassland" | ⚙ | 🗄 | 🖼 |
|---|---|---|---|
| Query | CLIP Embeddings | Vector DB | Result |

```
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
```
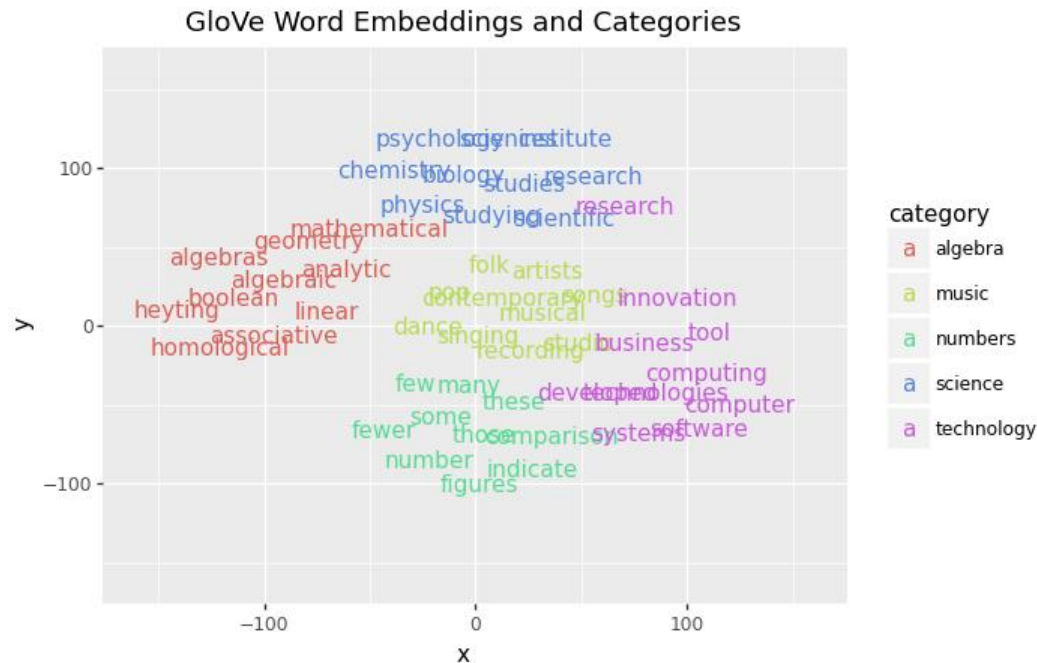
# Data Ingestion Pipeline:

# Similarity Search

# Vector Database

- Vector DB needs to analyze similarity of query-embedding compared to document embeddings.

- Approaches:
    - Cosine Similarity
    - Maximum Margin Relevance

# Vector Database

GloVe Word Embeddings and Categories

Example: word embeddings reduced to 2 dimensions

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

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

# Vector Database

| 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 |
|-----|------|-----|-----|-----|-----|

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

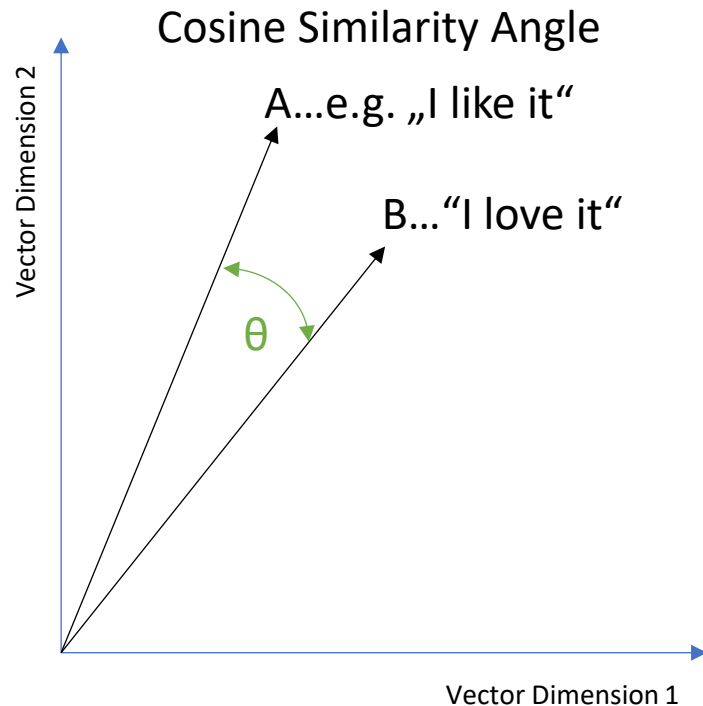For small data, go with np.array().

For large dataset not feasible!

# Vector Database
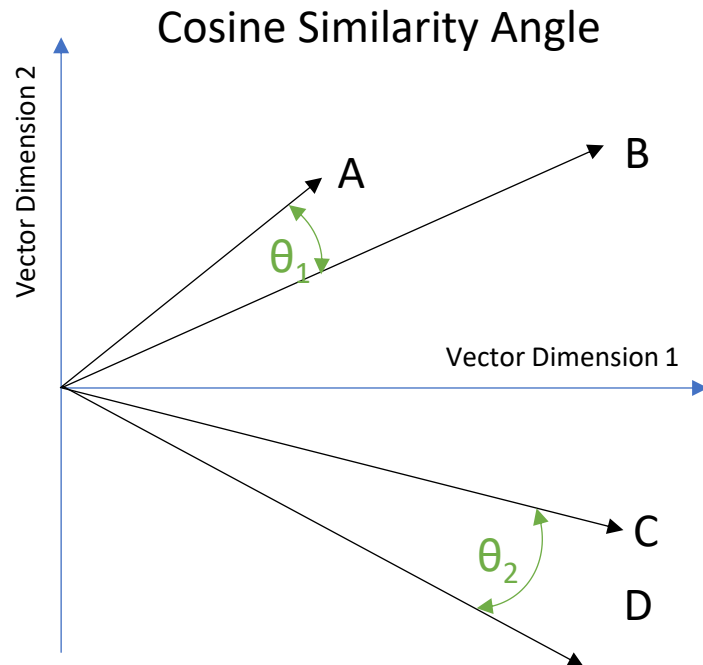
- Measures similarity between Embedding-Vectors based on angle θ.
  - Vectors maximally dissimilar
    - → vectors perpendicular (θ = 90°)
  - Vectors completely similar
    - → vectors parallel (θ = 0°)

Cosine Similarity Angle

A…e.g. „I like it"

B…"I love it"

θ

Vector Dimension 2

Vector Dimension 1

66

# Vector Database

- 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."
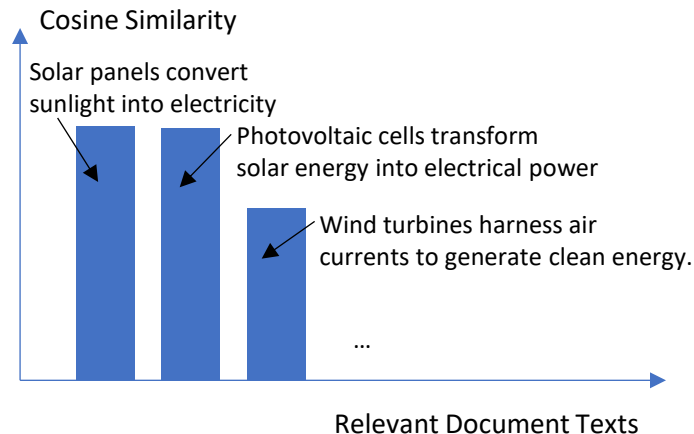  - C: "Trees grow leaves in spring."
  - D: "Fish swim in the ocean."



Cosine Similarity Angle

67

# Vector Database

Topic: Renewable Energies

- Aproach: reduce redundance 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?

Cosine Similarity

Solar panels convert sunlight into electricity

Photovoltaic cells transform solar energy into electrical power

Wind turbines harness air currents to generate clean energy.

…

Relevant Document Texts

68