# Real-World RAG:
# Eigene Daten & Dokumente mit semantischer Suche & LLMs erschließen

Sebastian Gingter

sebastian.gingter@thinktecture.com

Developer Consultant

Überschneidung
mit dem Montag-Workshop

think
tecture

## Real-World RAG:
# Eigene Daten & Dokumente mit semantischer Suche & LLMs erschließen

- Was Sie ERWARTET
    - Hintergrundwissen und Theorie zu RAG
    - Überblick über Semantische Suche
    - Probleme die auftreten können
    - Pragmatische Methoden für die Verwendung eigener Daten im RAG
    - Demos (Python)

- Was Sie NICHT erwartet
    - ChatGPT, CoPilot(s)
    - Grundlagen von ML
    - Deep Dives in LLMs, Vektor-Datenbanken, LangChain

# Sebastian Gingter

Developer Consultant @ Thinktecture AG

- Generative AI in business settings
- Flexible and scalable backends
- All things .NET

- Pragmatic end-to-end architectures
- Developer productivity
- Software quality

Special Day
# Generative AI für Business-Anwendungen

| Thema | Sprecher | Datum, Uhrzeit |
|---|---|---|
| **Large Language Models: Typische Use Cases & Patterns für Business-Anwendungen - in Action** | Christian Weyer | DI, 17. September 2024, 10.45 bis 11.45 |
| **Real-World RAG: Eigene Daten & Dokumente mit semantischer Suche & LLMs erschließen** | Sebastian Gingter | DI, 17. September 2024, 12.15 bis 13.15 |
| **Von 0 zu Smart: SPAs mit Generative AI aufwerten** | Max Marschall | DI, 17. September 2024, 15.30 bis 16.30 |
| **Deep Dive in OpenAI Hosted Tools** | Rainer Stropek | DI, 17. September 2024, 17.00 bis 18.00 |

BASTA!

# Was Euch erwartet (und was nicht):

- Ein bisschen Hintergrund-Info & Theorie

- Überblick über das Themengebiet Semantische Suche

- Probleme und mögliche Strategien

- Pragmatische Ansätze für die eigenen Daten

- Kein C#, sondern Python 😱

- Kein Deep-Dive in
    - LLMs
    - Vektor-Datenbanken
    - LangChain

# Agenda

- Short Introduction to RAG

- Embeddings (and a bit of theory 😱 )

- Vector-Databases

- Indexing

- Retrieval

- Not good enough? – Indexing II
  - HyDE & alternative indexing methods

- Conclusion

# Introduction

Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG

# Use-case: Talk to my internal data

# Retrieval-augmented generation (RAG)
## Indexing & (Semantic) search



*Indexing / Embedding*

QA

# Semantic Search

- Classic search: lexical
    - Compares words, parts of words and variants
    - Classic SQL: WHERE 'content' LIKE '%searchterm%'
    - We can search only for things where we know
      that its somewhere in the text

- New: Semantic search
    - Compares for the same contextual meaning
        - "Das Rudel rollt das runde Gerät auf dem Rasen herum"
        - "The pack enjoys rolling a round thing on the green grass"
        - "Die Hunde spielen auf der Wiese mit dem Ball"
        - "The dogs play with the ball on the meadow"

# Semantic Search

- How to grasp "semantics"?

- Computers only calculate on numbers
  - Computing is "applied mathematics"

- AI also only calculates on numbers

- We need a numeric representation of meaning
  - ➔ "Embeddings"

# Embeddings

| Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG |

# Embedding (math.)

- Topologic: Value of a high dimensional space is "embedded" into a lower dimensional space


- Natural / human language is very complex (high dimensional)
  - Task: Map high complexity to lower complexity / dimensions


- Injective function


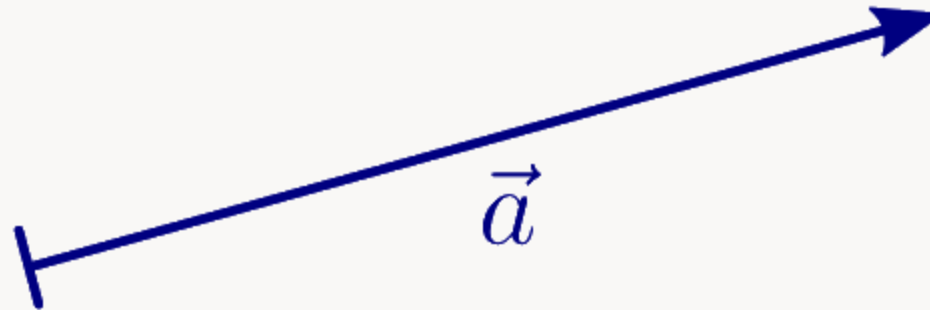- Similar to hash, or a lossy compression

# Embeddings

- Embedding model (specialized ML model) converting text into a numeric representation of its meaning

- Representation is a vector in an n-dimensional space
    - n floating point values
    - OpenAI
        - "text-embedding-ada-002" uses 1536 dimensions
        - "text-embedding-3-small" 512 and 1536
        - "text-embedding-3-large" 256, 1024 and 3072
    - Huggingface models have a very wide range of dimensions

# Embeddings

- Embedding models are unique

- Each dimension has a different meaning, individual to the model
- vectors from different models are incompatible with each other

- Some embedding models are multi-language, but not all

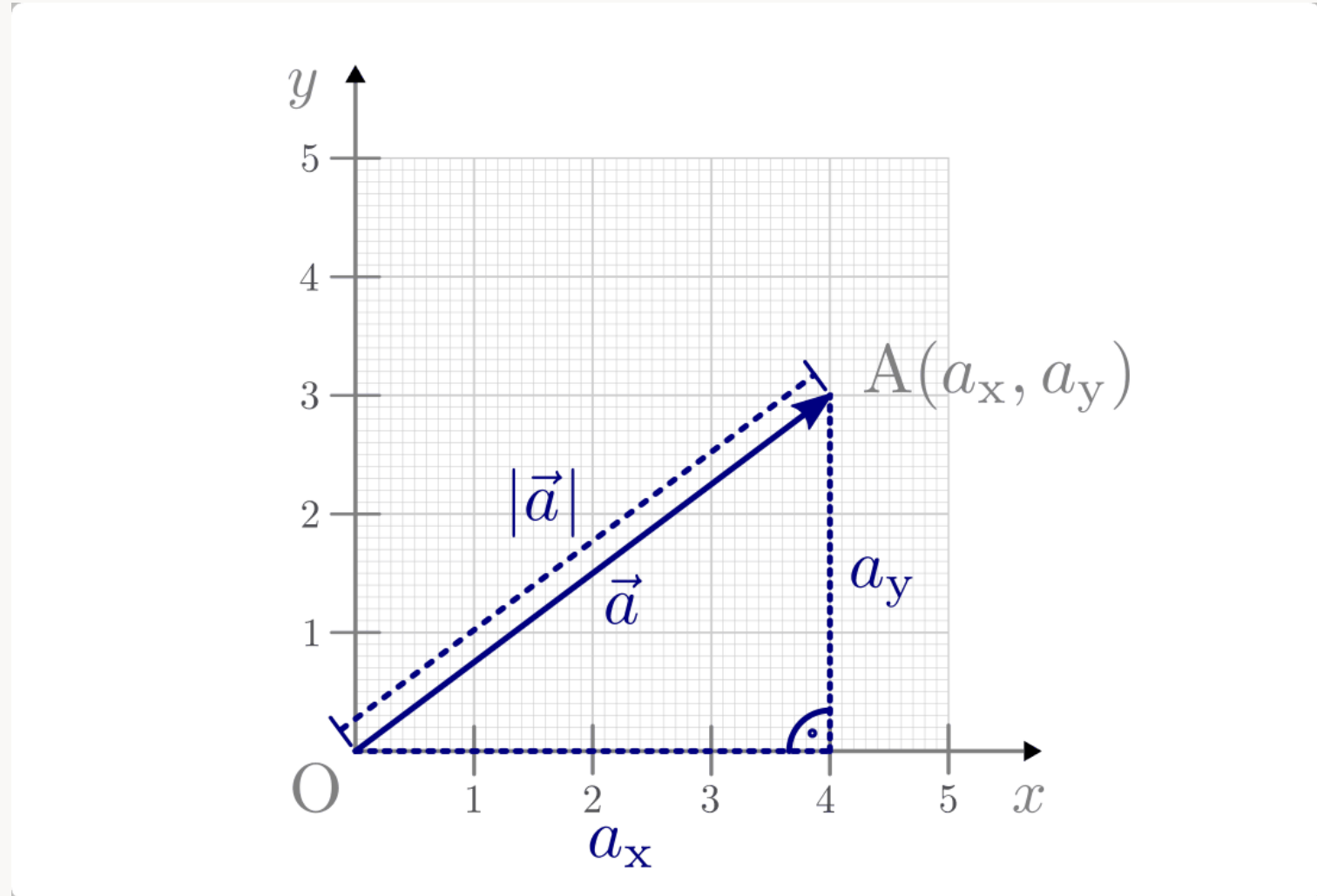- In an LLM, also the first step is to embed the input into a lower dimensional space

# What is a vector?

- Mathematical quantity with a direction and length

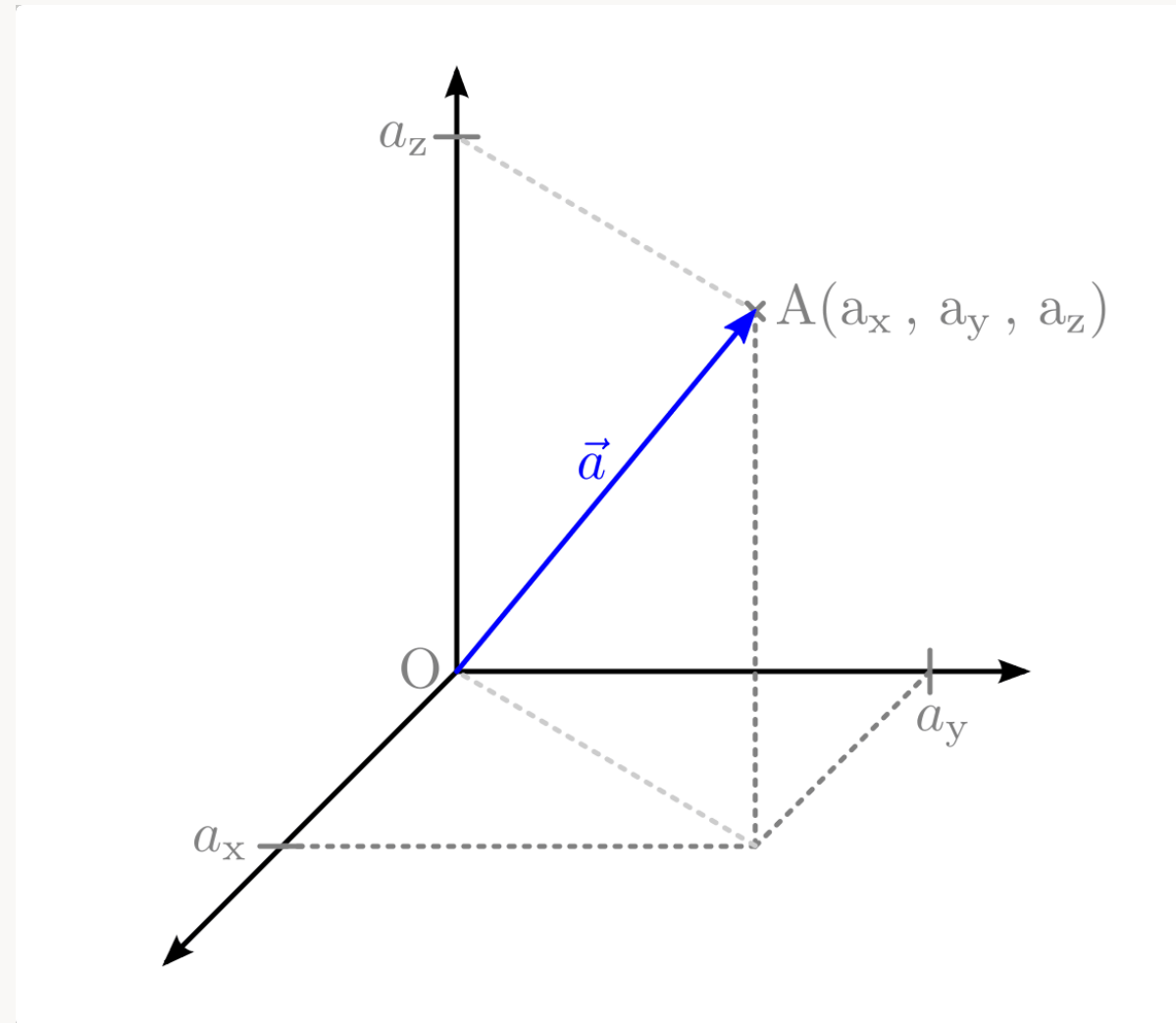- $\vec{a} = \begin{pmatrix} a_x \\ a_y \end{pmatrix}$

# Vectors in 2D

$$\vec{a} = \begin{pmatrix} a_x \\ a_y \end{pmatrix}$$
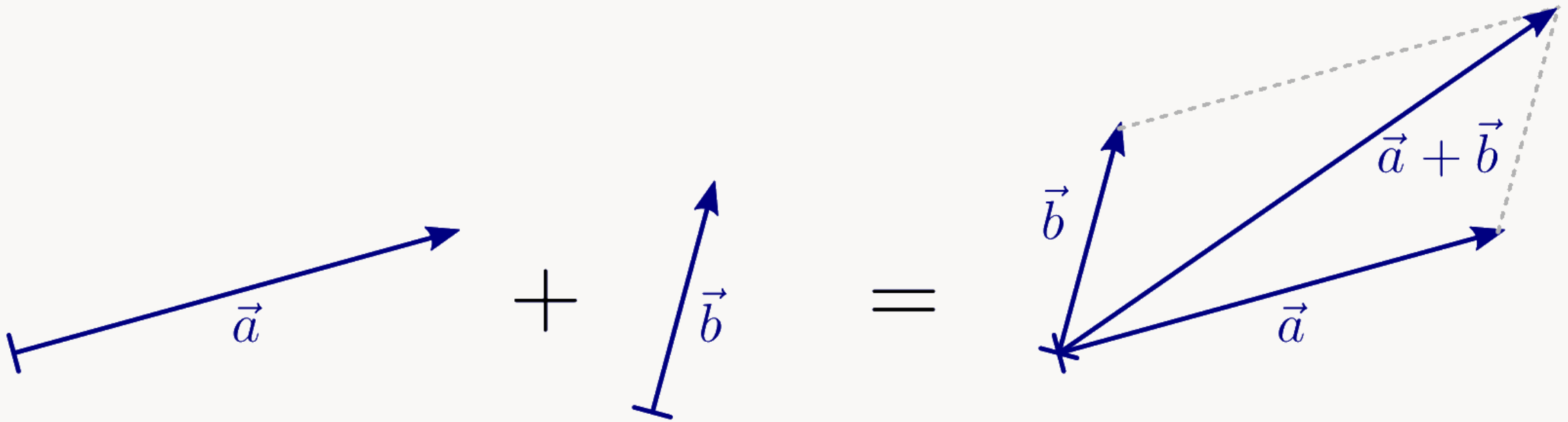
# Vectors in 3D



$$\vec{a} = \begin{pmatrix} a_x \\ a_y \\ a_z \end{pmatrix}$$

# Vectors in multidimensional space

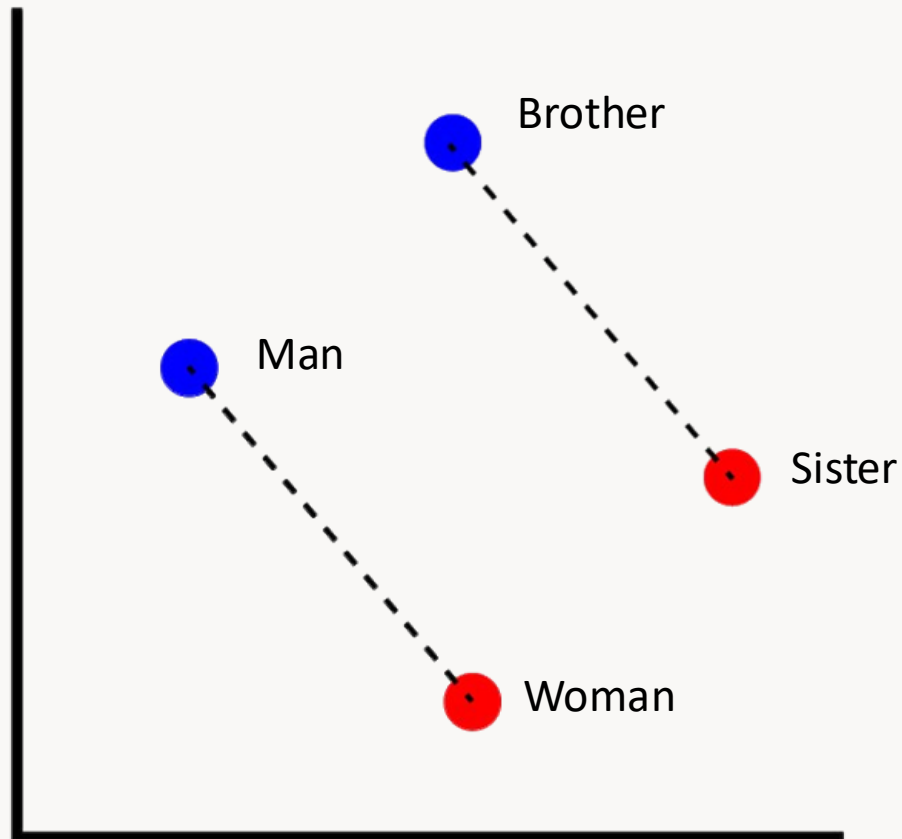$$\vec{a} = \begin{pmatrix} a_u \\ a_v \\ a_w \\ a_x \\ a_y \\ a_z \end{pmatrix}$$

# Calculation with vectors
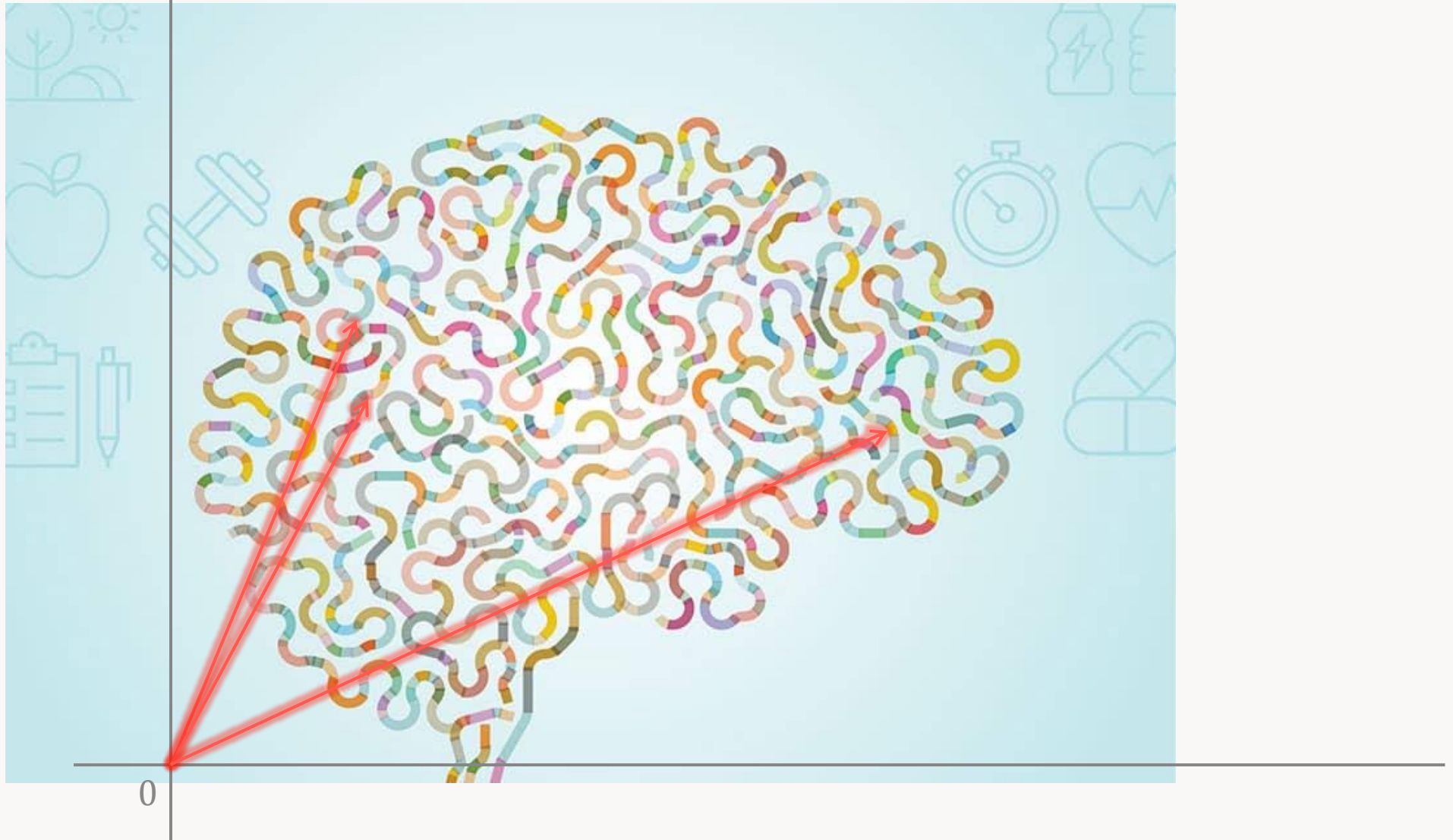
# Word2Vec
**Mikolov et al., Google, 2013**

$$Brother - Man + Woman \approx Sister$$

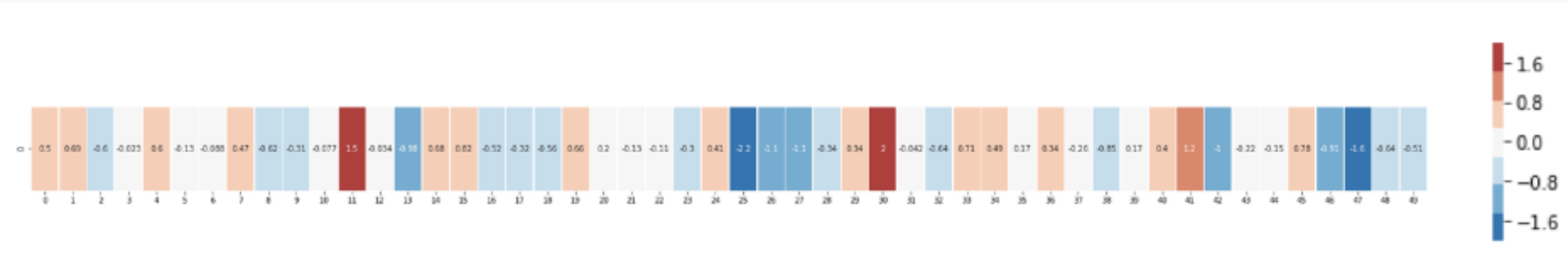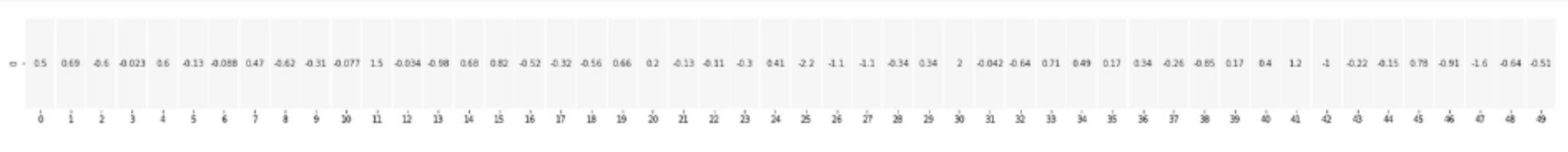# Embedding-Model

- Task: Create a vector from an input
  - Extract meaning / semantics

- Embedding models usually are very shallow & fast
  Word2Vec is only two layers

- Similar to the first step of an LLM
  - Convert text to values for input layer

- This comparison is very simplified, but one could say:
  - The embedding model 'maps' the meaning into the model's 'brain'

# Embedding-Model

# Embedding-Model

[ 0.50451 , 0.68607 , -0.59517 , -0.022801, 0.60046 , -0.13498 , -0.08813 , 0.47377 , -0.61798 , -0.31012 , -0.076666, 1.493 , -0.034189, -0.98173 , 0.68229 , 0.81722 , -0.51874 , -0.31503 , -0.55809 , 0.66421 , 0.1961 , -0.13495 , -0.11476 , -0.30344 , 0.41177 , -2.223 , -1.0756 , -1.0783 , -0.34354 , 0.33505 , 1.9927 , -0.04234 , -0.64319 , 0.71125 , 0.49159 , 0.16754 , 0.34344 , -0.25663 , -0.8523 , 0.1661 , 0.40102 , 1.1685 , -1.0137 , -0.21585 , -0.15155 , 0.78321 , -0.91241 , -1.6106 , -0.64426 , -0.51042 ]





http://jalammar.github.io/illustrated-word2vec/

# Embedding-Model



http://jalammar.github.io/illustrated-word2vec/

# Recap Embeddings

- Embedding model: "Analog to digital converter for text"

- Embeds the high-dimensional natural language meaning into a lower dimensional-space (the model's 'brain')

- No magic, just applied mathematics

- Math. representation: Vector of n dimensions

- Technical representation: array of floating point numbers

# DEMO

## Embeddings
Sentence Transformers, local embedding model

# Vector-Databases

| Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG |

# Vector-Databases

- Mostly document-based

- Index: Embedding (vector)

- Document (content)

- Metadata

- Query functionalities

# Vector-Databases

- Pinecone

- Milvus

- Chroma

- Weaviate

- Deep Lakee

- Qdrant

- Elasticsearch

- Vespa

- Vald

- ScaNN

- Pgvector
(PostgreSQL Extension)

- Faiss

- ...

- ... (probably) coming to a relational database near you soon(ish)
  SQL Server Example: https://learn.microsoft.com/en-us/samples/azure-samples/azure-sql-db-openai/azure-sql-db-openai/

# Vector-Databases

- (Search-)Algorithms

  - Cosine Similarity  $S_{C(a,b)} = \dfrac{a \cdot b}{\|a\| \times \|b\|}$

  - Manhatten Distance (L1 norm, taxicab)

  - Euclidean Distance (L2 norm)

  - Minkowski Distance (~ generalization of L1 and L2 norms)

  - L∞ ( L-Infinity), Chebyshev Distance

  - Jaccard index / similarity coefficient (Tanimoto index)

  - Nearest Neighbour

  - Bregman divergence

  - …

# DEMO

Vector database
LangChain, Chroma, local embedding model

# Indexing

| Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG |

# Indexing

- Loading

- Clean-up

- Splitting

- Embedding

- Storing

# Loading

- Import documents from different sources, in different formats

- LangChain has very strong support for loading data

- Support for cleanup

- Support for splitting



**Document loaders**

| | |
|---|---|
| 📄 **mhtml** <br> MHTML is a is used both for emails but also for archived webpag… | 📄 **Microsoft Excel** <br> The UnstructuredExcelLoader is used to load Microsoft Excel files. |
| 📄 **Microsoft OneDrive** <br> Microsoft OneDrive (formerly | 📄 **Microsoft OneNote** <br> This notebook covers how to load documents from OneNote. |
| 📄 **Microsoft PowerPoint** <br> [Microsoft | 📄 **Microsoft SharePoint** <br> Microsoft SharePoint is a |
| 📄 **Microsoft Word** <br> Microsoft Word | 📄 **Modern Treasury** <br> Modern Treasury simplifies complex |

# Clean-up

- HTML Tags
- Formatting information
- Normalization
  - lowercasing
  - stemming, lemmatization
  - remove punctuation & stop words
- Enrichment
  - tagging
  - keywords, categories
  - metadata

# Splitting (Text Segmentation)

- Document is too large / too much content / not concise enough



- by size (text length)
- by character (\n\n)
- by paragraph, sentence, words (until small enough)
- by size (tokens)
- overlapping chunks (token-wise)

# Vector-Databases

- Indexing



Document      Splitted (smaller) parts      Embedding-Model      Embedding      Vector-Database

Metadata: Reference to original document

$$\begin{pmatrix} a \\ b \\ c \\ ... \end{pmatrix}$$

# Retrieval (Search)

| Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG |

# Retrieval



"What is the name
of the teacher?"

Query

Embedding-
Model

$$\begin{pmatrix} a \\ b \\ c \\ \dots \end{pmatrix}$$

Embedding

Vector-
Database

Doc. 1:  0.86
Doc. 2:  0.84
Doc. 3:  0.79

Weighted result

... (Answer generation)

# DEMO

Store and retrieval
LangChain, Chroma, local embedding model, OpenAI GPT

# Indexing II
# Not good enough?

| Introduction | Embeddings | Vector-DBs | Indexing | Retrieval | Indexing II | RAG |

# Not good enough?

# Not good enough?

- Semantic search still only uses your index


- It's just as good as your embeddings
  - All chunks need to be


- Sh*t in, sh*t out

# HyDE (Hypothetical Document Embedddings)

- Search for a hypothetical Document



"What should I do, if I missed the last train?"
Query

Write a company policy that contains all information which will answer the given question: {QUERY}

LLM, e.g. GPT-3.5-turbo

Hypothetical Document

$$\begin{pmatrix} a \\ b \\ c \\ ... \end{pmatrix}$$

Embedding-Model

Embedding

Vector-Database

Doc. 3:  0.86
Doc. 2:  0.81
Doc. 1:  0.81

Weighted result

# What else?

- Downside of HyDE:
    - Each request needs to be transformed through an LLM
      (slow & expensive)
    - A lot of requests will probably be very similar to each other
    - Each time a different hyp. document is generated,
      even for an extremely similar request
        - Leads to very different results each time
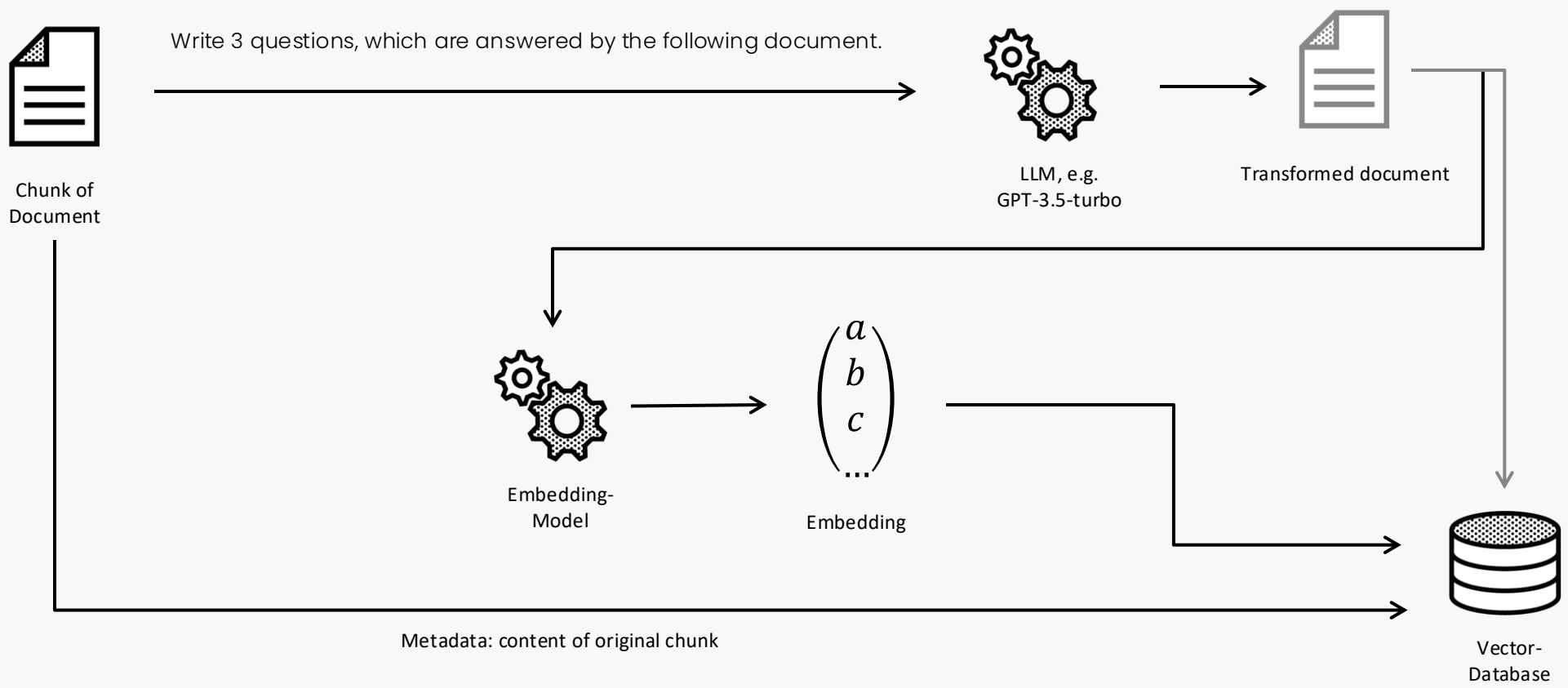

- Idea: Alternative indexing
    - Transform the document, not the query

# Alternative Indexing

## HyQE: Hypothetical Question Embedding



Chunk of
Document

Write 3 questions, which are answered by the following document.

LLM, e.g.
GPT-3.5-turbo

Transformed document

$$\begin{pmatrix} a \\ b \\ c \\ \dots \end{pmatrix}$$

Embedding-
Model

Embedding

Metadata: content of original chunk

Vector-
Database

# Alternative Indexing

- Retrieval



"What should I do, if I missed the last train?"

Query

Embedding-Model

$$\begin{pmatrix} a \\ b \\ c \\ ... \end{pmatrix}$$

Embedding

Vector-Database

Doc. 3: 0.89
Doc. 1: 0.86
Doc. 2: 0.76

Weighted result

Original document from metadata

# DEMO

Compare embeddings
LangChain, Qdrant, OpenAI GPT
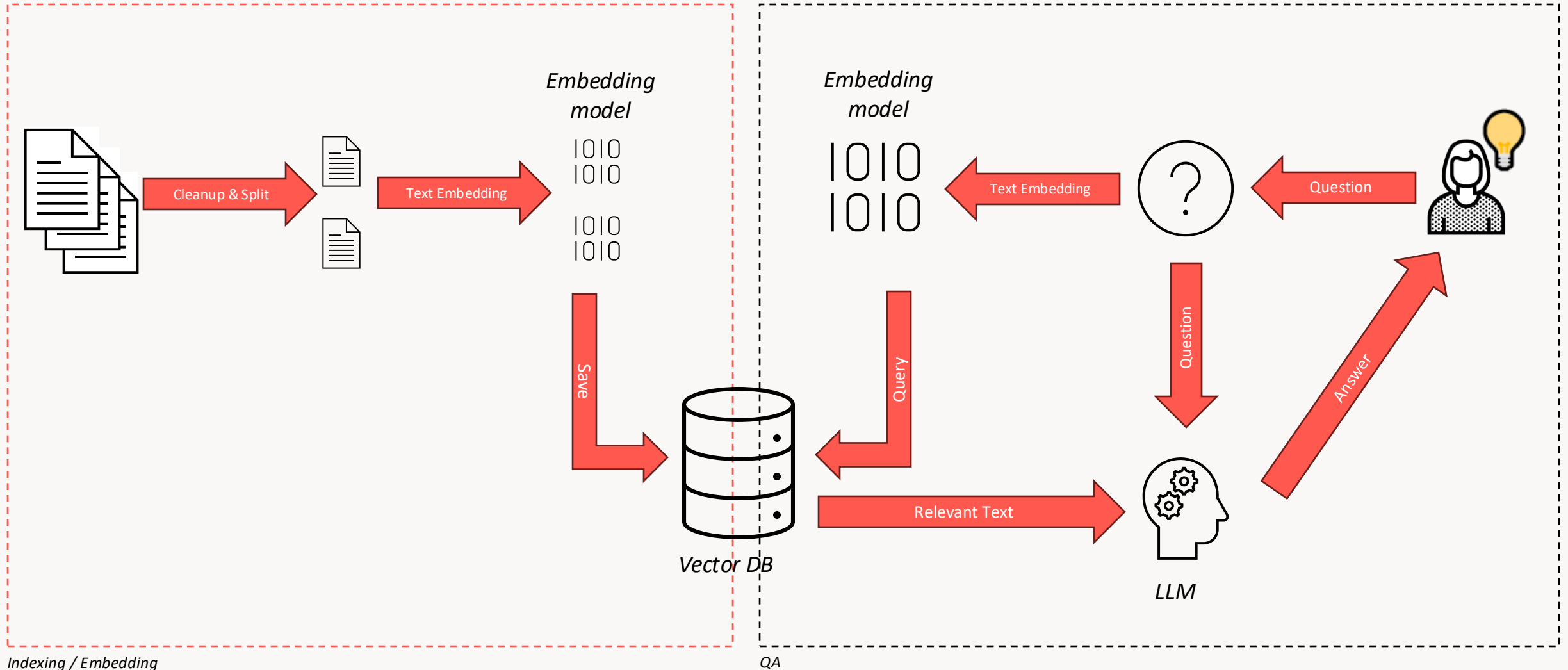
# Conclusion

# Retrieval-augmented generation (RAG)
## Indexing & (Semantic) search



*Embedding model*

Cleanup & Split

Text Embedding

Save

*Embedding model*

Text Embedding

Question

Query

Question

Answer

Relevant Text

*Vector DB*
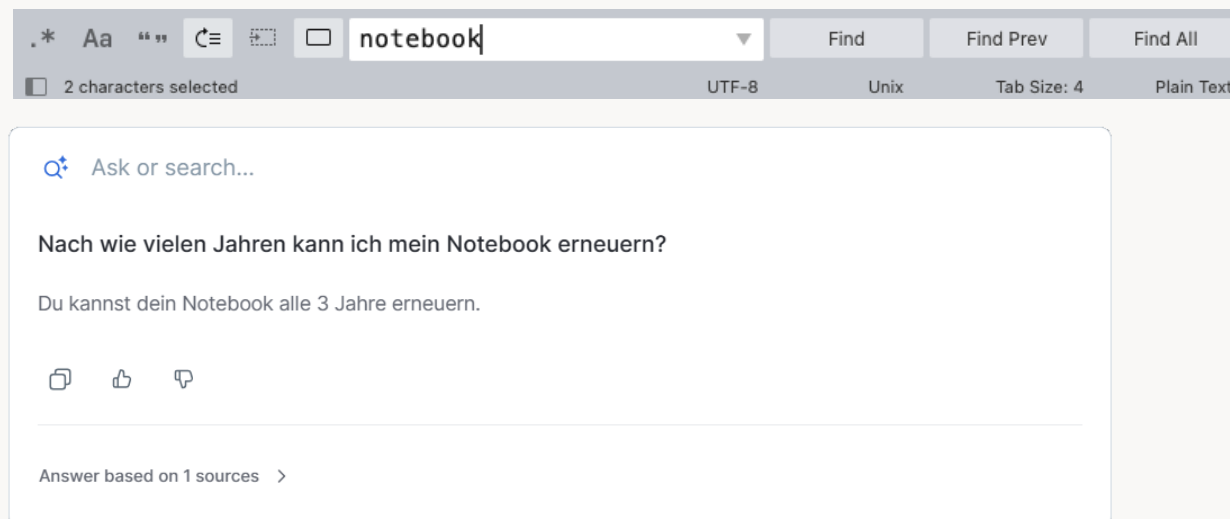
*LLM*

*Indexing / Embedding*

*QA*

# Recap: Not good enough?

- Tune text cleanup, segmentation, splitting

- HyDE or HyQE or alternative indexing
    - How many questions?
    - With or without summary

- Other approaches
    - Only generate summary
    - Extract "Intent" from user input and search by that
    - Transform document and query to a common search embedding
    - HyKSS: Hybrid Keyword and Semantic Search
      https://www.deg.byu.edu/papers/HyKSS.pdf

- Always evaluate approaches with your own data & queries

- The actual / final approach is more involved
  as it seems on the first glance

55

# Conclusion

- Semantic search is a first and fast Generative AI business use-case

- Quality of results depend heavily on data quality
  and preparation pipeline

- RAG pattern can produce breathtaking good results
  without the need for user training

# Thank you!

think
tecture

Demos:

https://github.com/thinktecture-labs/basta-2024-advanced-rag

Sebastian Gingter

https://thinktecture.com/sebastian-gingter

# Real-World RAG:
# Eigene Daten & Dokumente mit semantischer Suche & LLMs erschließen

think
tecture

## Slides & Code

https://www.thinktecture.com/de/sebastian-gingter

Sebastian Gingter

sebastian.gingter@thinktecture.com

Developer Consultant