



# Vector DBs: a deep dive with Weaviate



**JP Hwang**

Senior developer educator



# Study guide

Key topics covered:

- What vectors are
- Differences between search types
- Role of embedding models
- Role of generative models
- Overview of HNSW
- Agentic AI

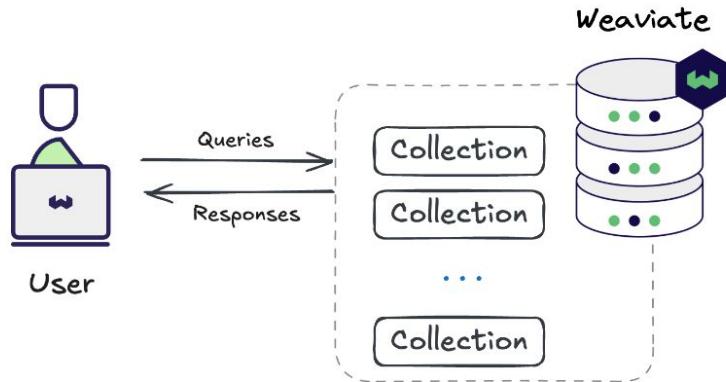
Example questions:

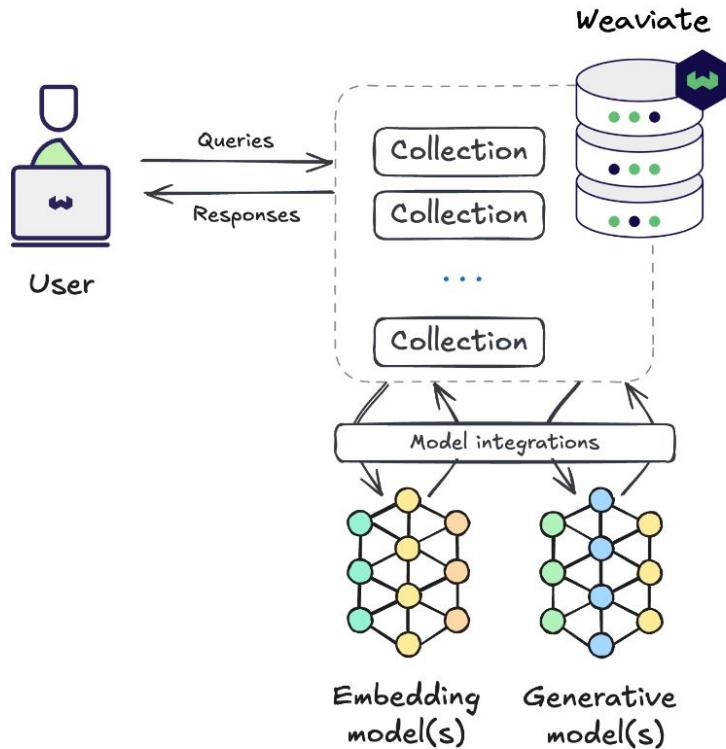
- What does a vector embedding capture?
- Which of hybrid, vector or keyword searches are based on a comparison of meanings?
- What does an embedding model output?
- What does a generative AI model output?
- In the G step of RAG, which type (embedding or generative) of model is used?
- In an HNSW index, does search begin on a higher (sparse) or a lower (dense) layer?
- What makes an application “agentic”?

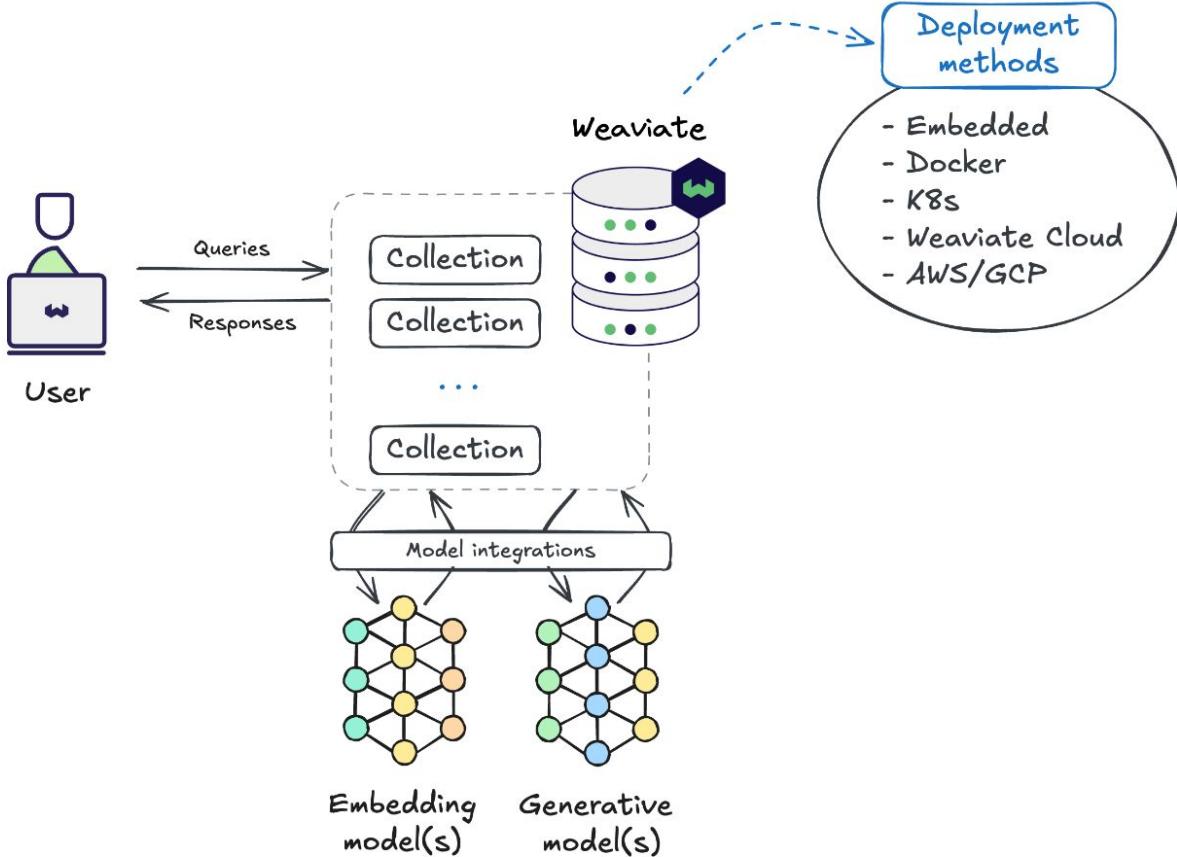


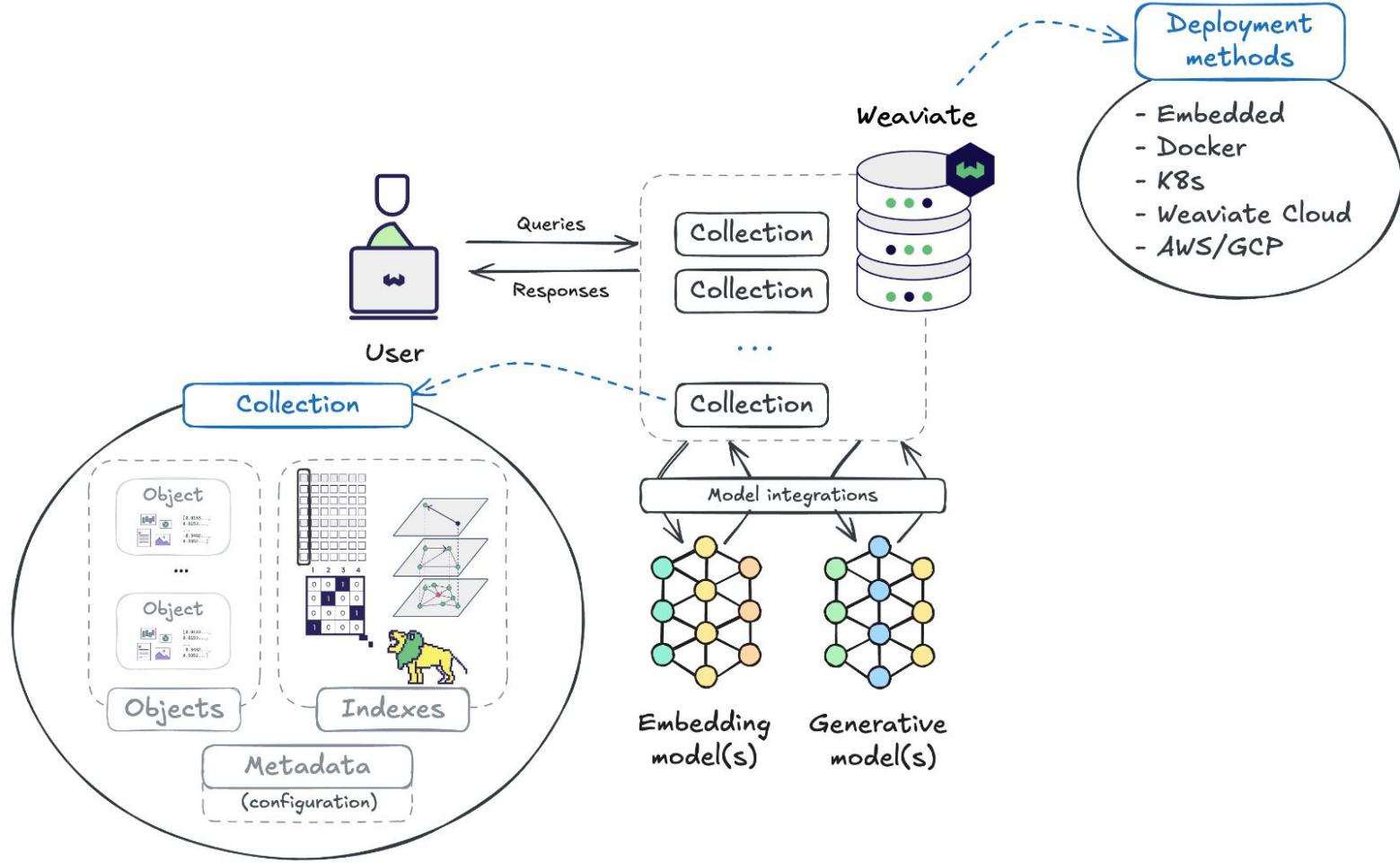
# What is Weaviate ?

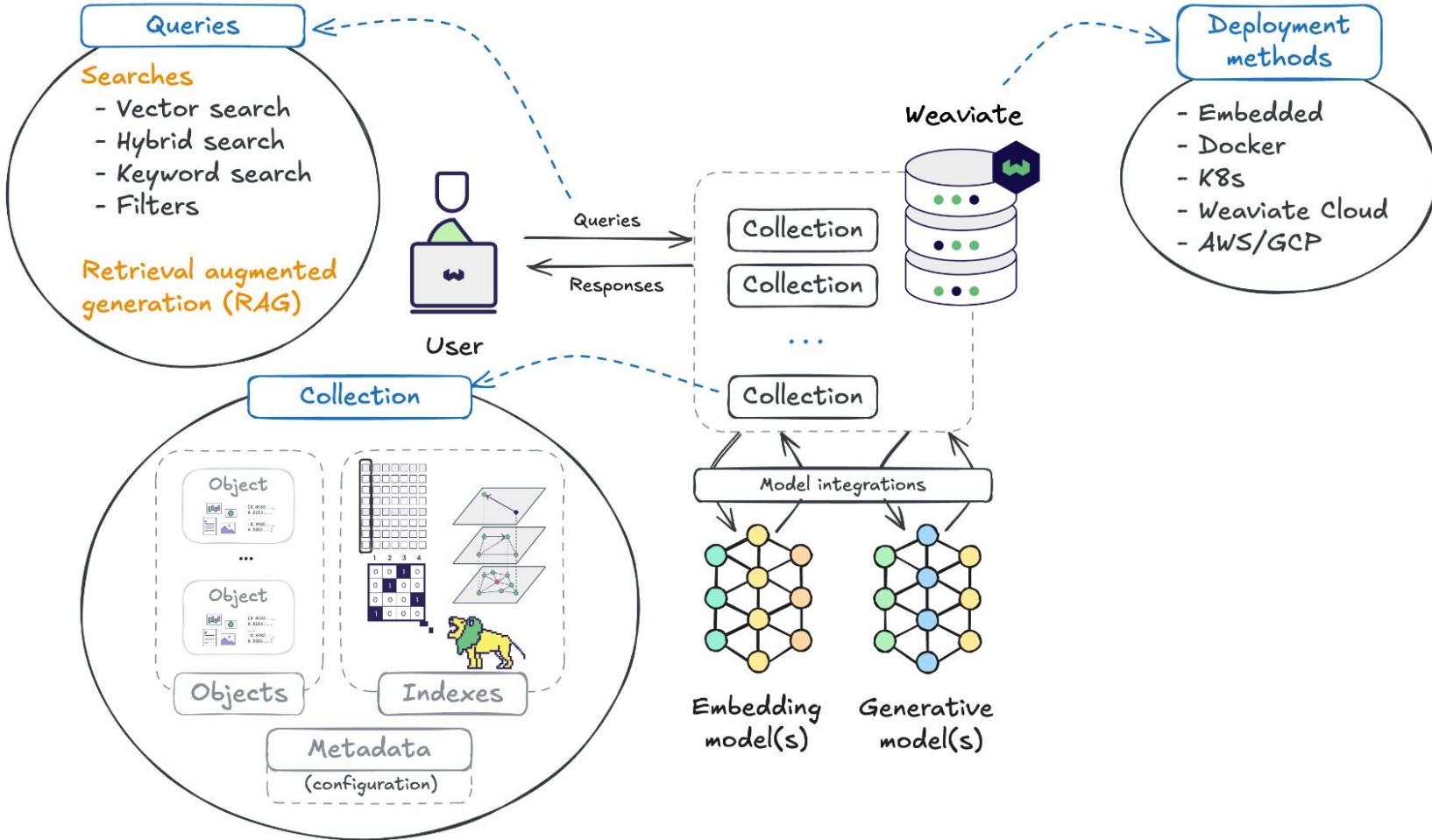
An introduction to vectors and vector DBs

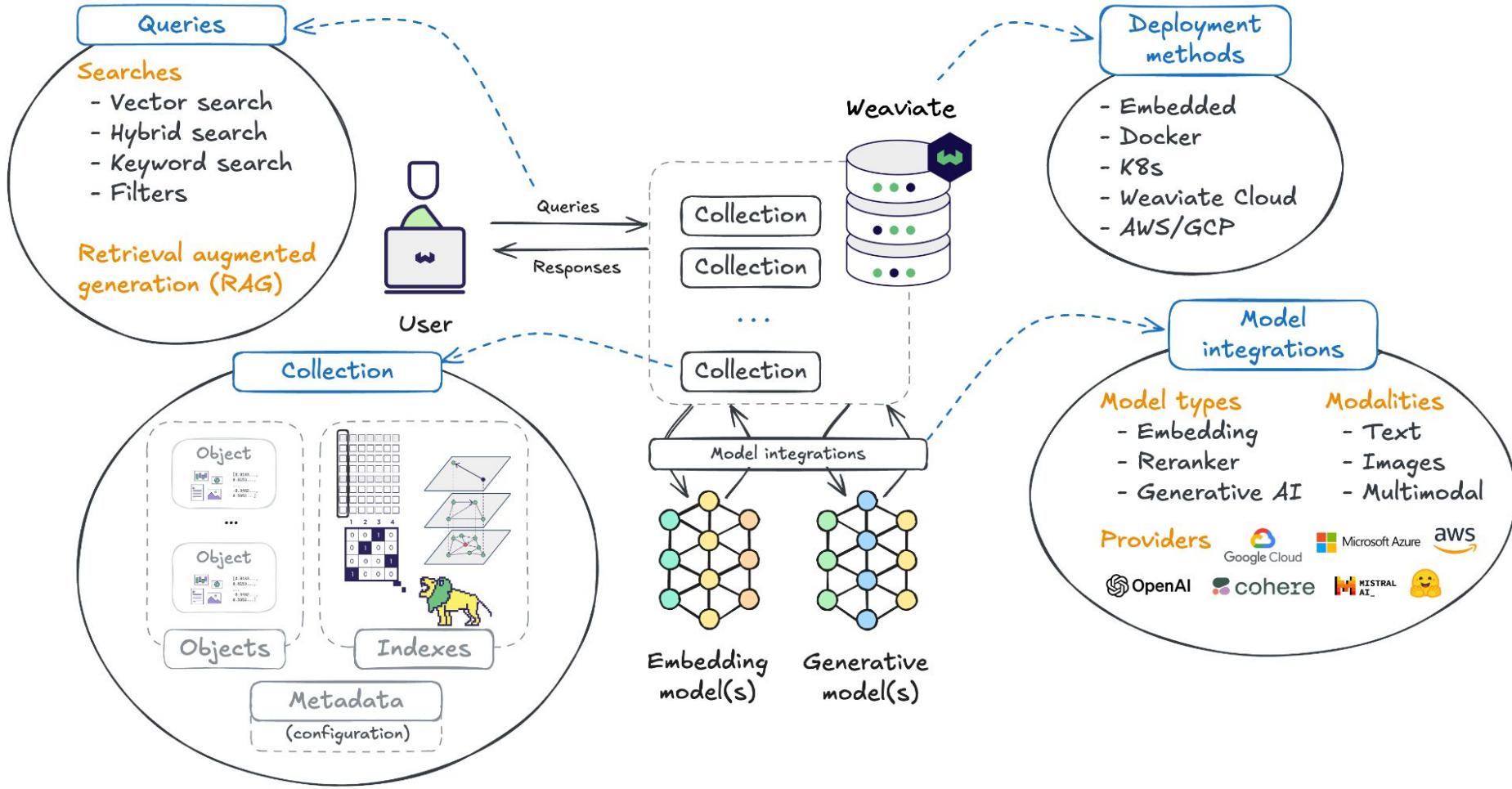














# A quick demo

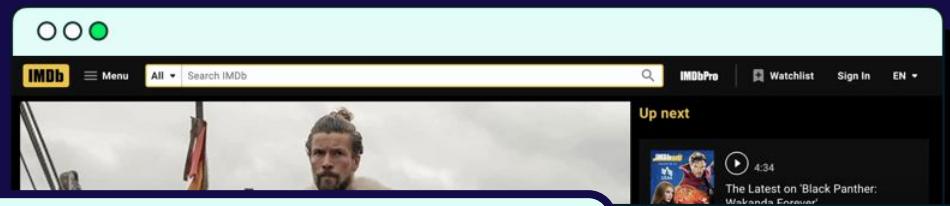
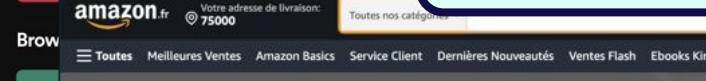
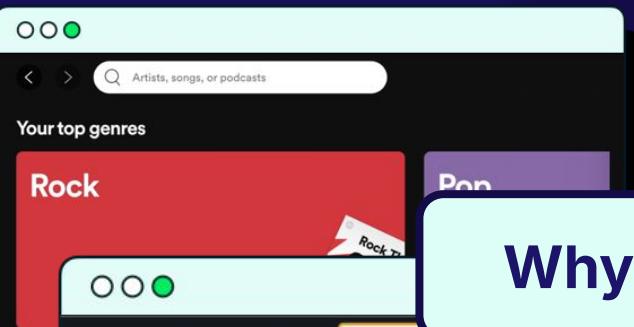
Where we are headed



**Everything we do  
starts with search**

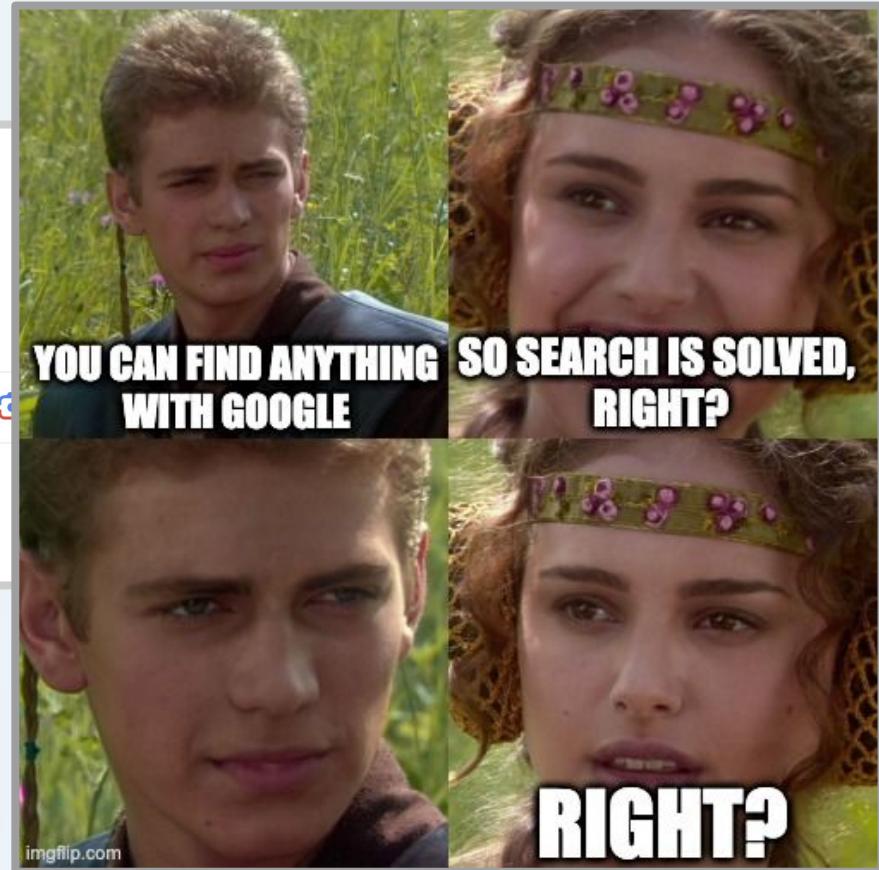
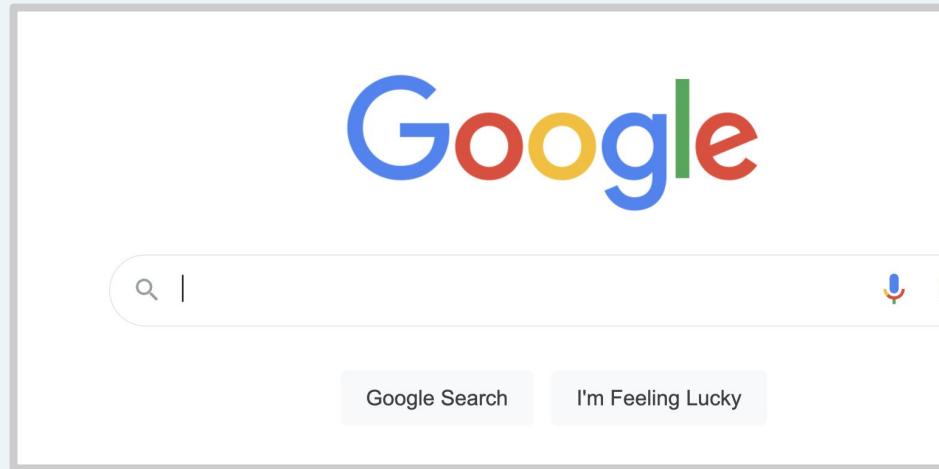


# Why is search so ubiquitous?





# We are spoiled

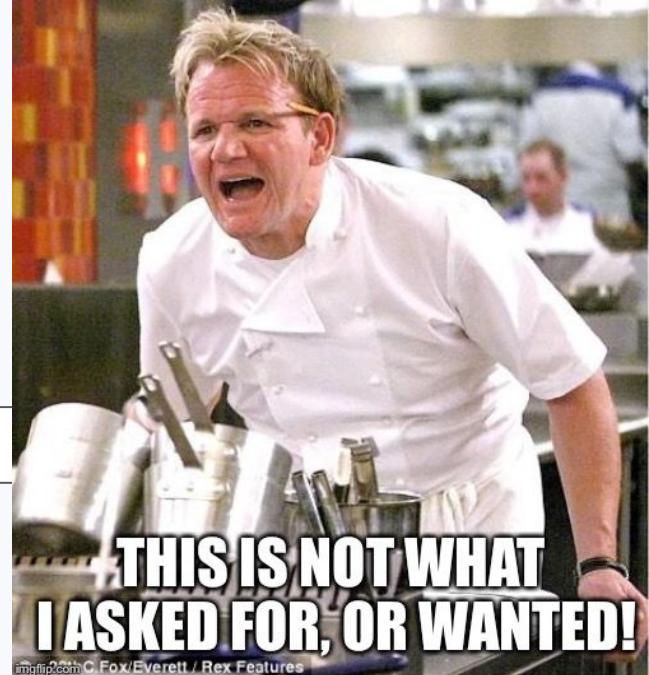




# The problem of effective search

# Traditional search has shortcomings

How do airplanes fly?



[Why you should fly with ExpensiveAir's airplanes!](#)

...fly with ExpensiveAir's modern fleet of airplanes to  
over 4000 destinations. How to book: ...

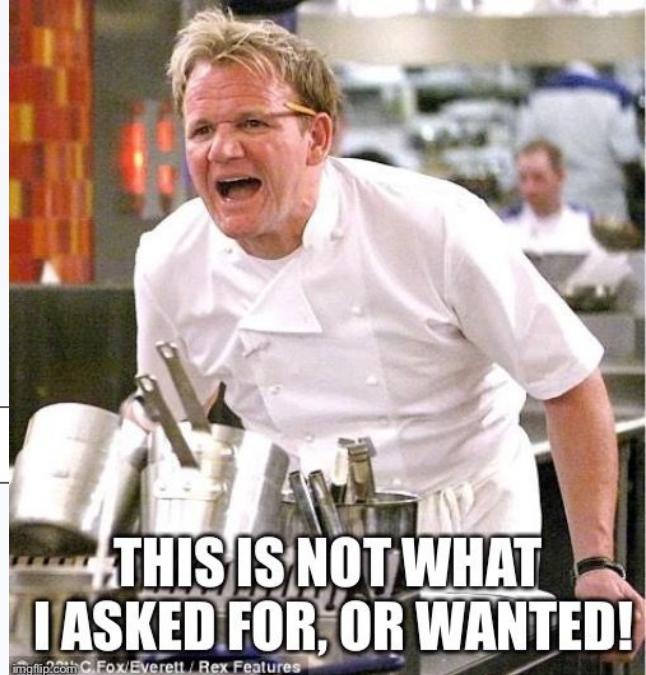
# Traditional search has shortcomings

How do airplanes fly?



[Why you should fly with ExpensiveAir's airplanes!](#)

...[fly](#) with ExpensiveAir's modern fleet of [airplanes](#) to  
over 4000 destinations. [How](#) to book: ...



THANK YOU

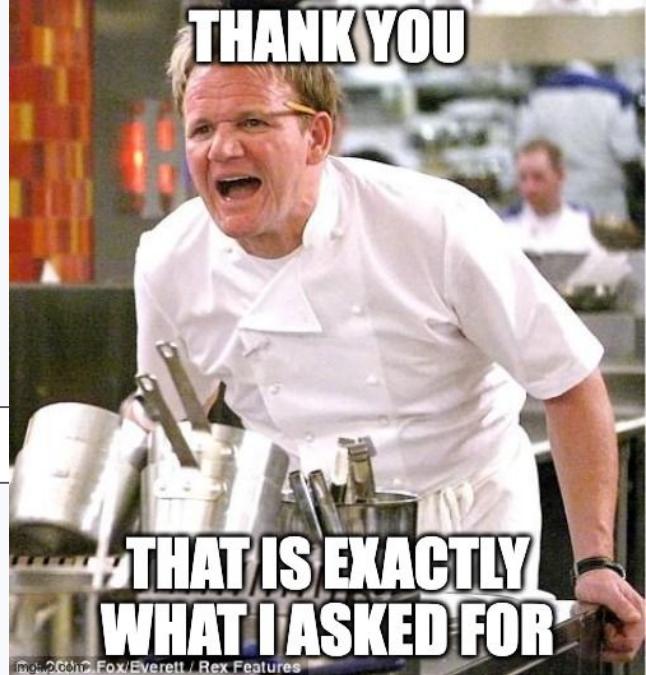
# Semantic search to the rescue!

How do airplanes fly?



[Dynamics of Flight - NASA](#)

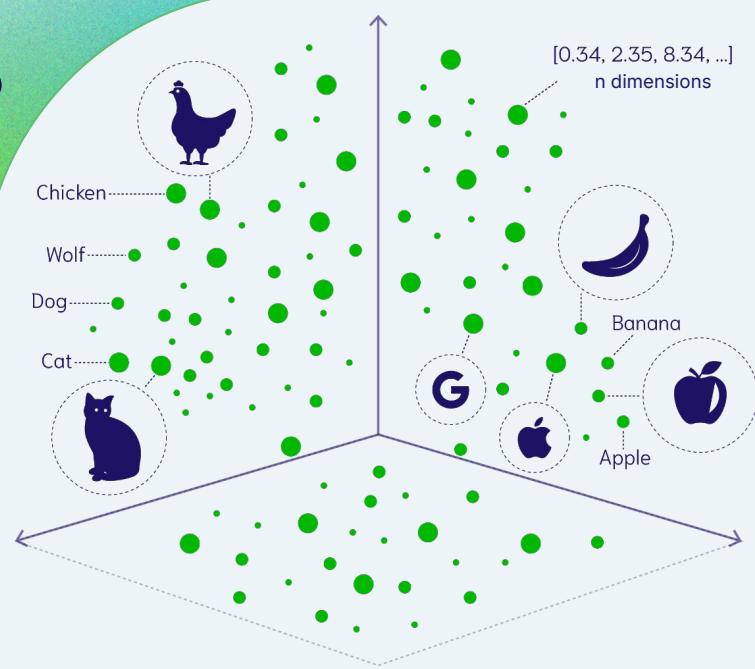
Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases.





# How does this work?

Vectors!





# What is a vector?

A vector is a set of numbers  
That represents some  
meaning

Like

[1, 0]

or

[0.513, 0.155, 0.983, 0.001,  
0.932]

or

[0.0009420722, 0.020158706,  
-0.03939992, 0.018441677, ...,  
-0.025270471, -0.056622636, ...]

# Numbers represent meaning?

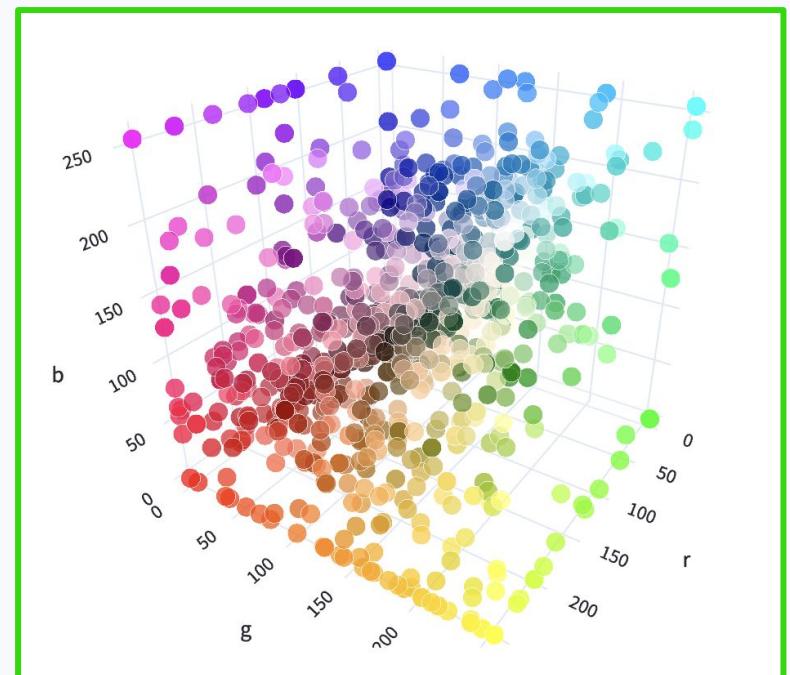
Here's an example:

RGB numbers represent colors, like:

(255, 0, 0) = red

(80, 200, 120) = emerald

Each number is a dial  
for (red, green, blue) ness.

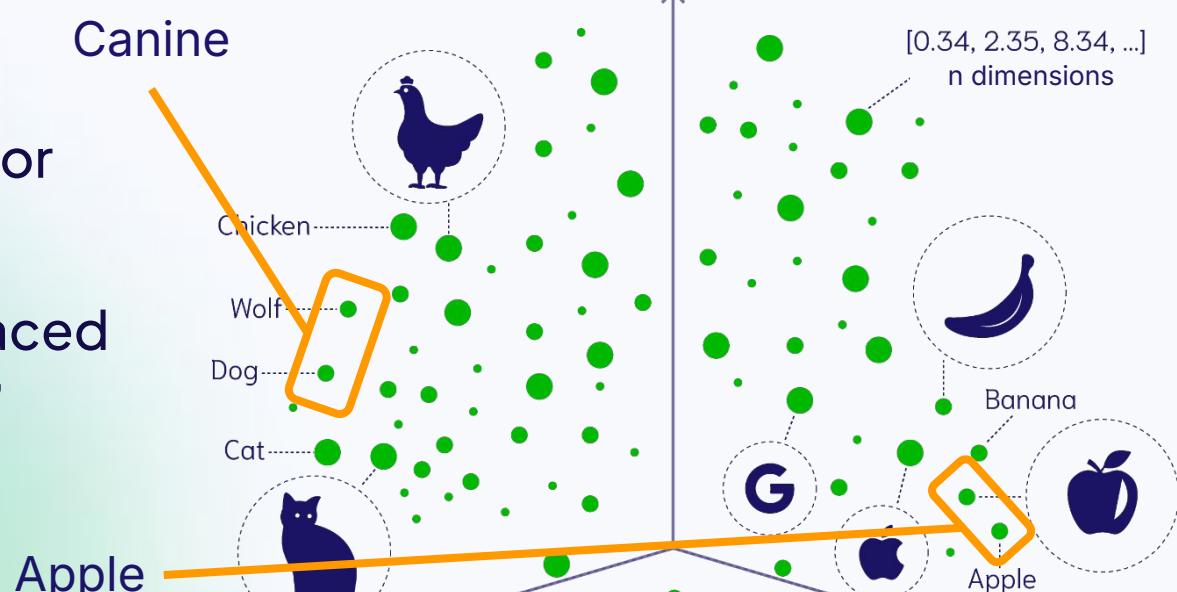


# Numbers represent meaning

Any meaning can be represented with a vector

Similar meanings are placed closer in “vector space”

Aerodynamics  
of flight

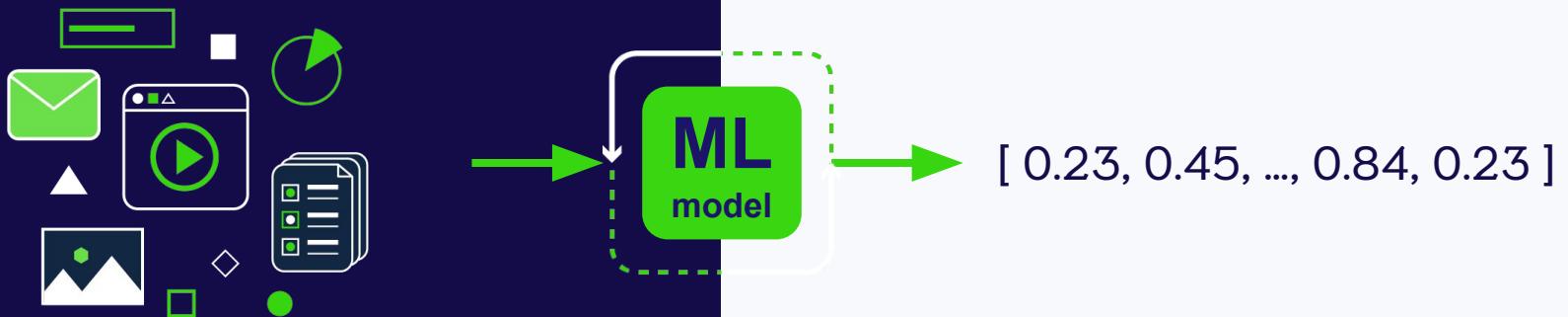




# Where do vectors come from?



# Machine Learning models





# Vector Embeddings

“Fly with ExpensiveAir’s modern fleet of airplanes to over 4000 destinations. How to book...”

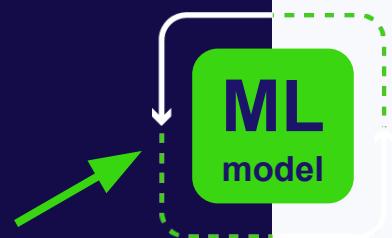




# Vector Embeddings

“Fly with ExpensiveAir’s modern fleet of airplanes to over 4000 destinations. How to book...”

Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases.

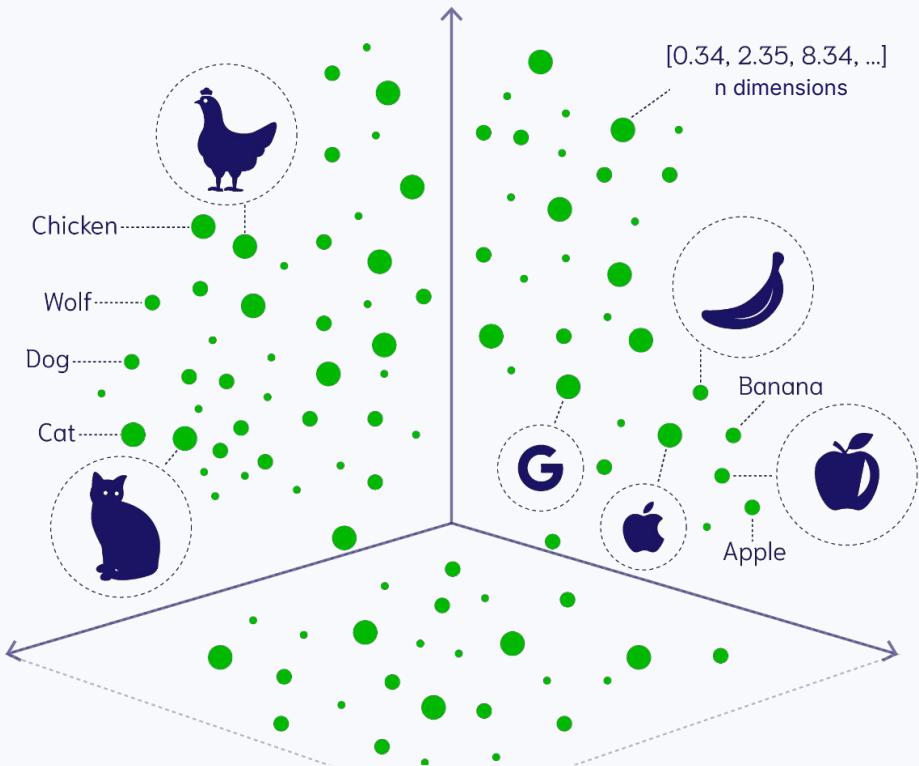


[ 0.23, 0.45, ..., 0.84, 0.23 ]

[ 0.26, 0.31, ..., 0.12, 0.44 ]

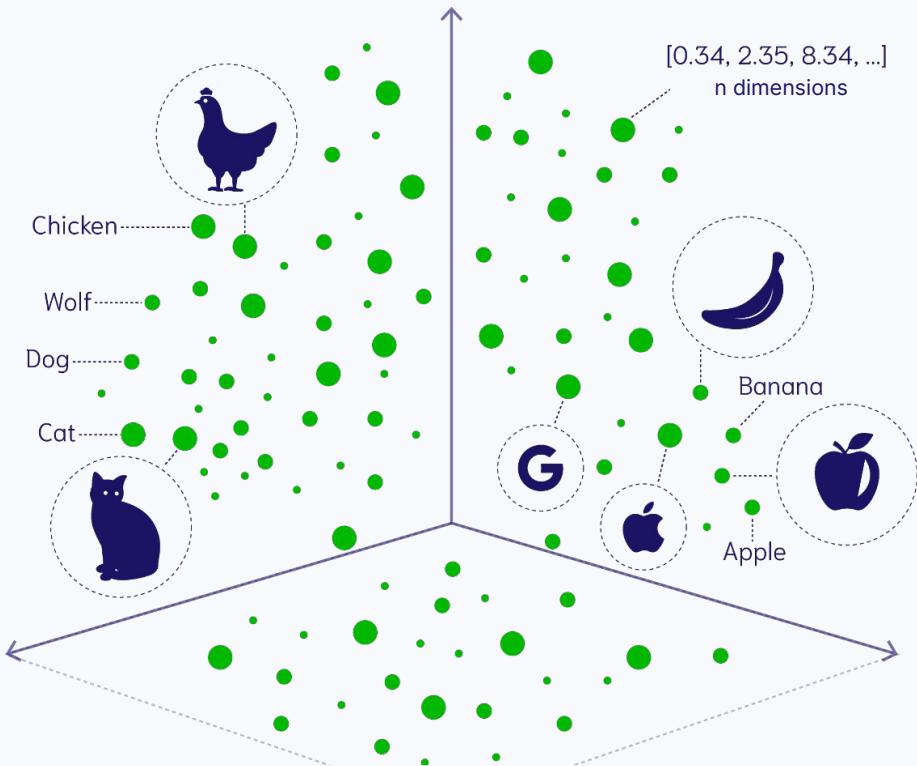
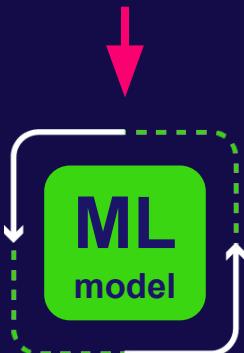


# Vector Space



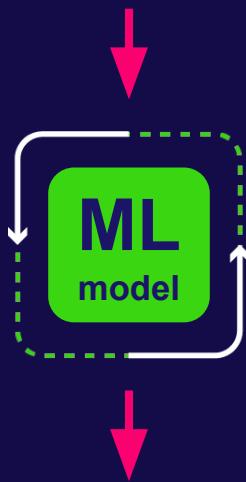
# Vector Search

How do airplanes fly?

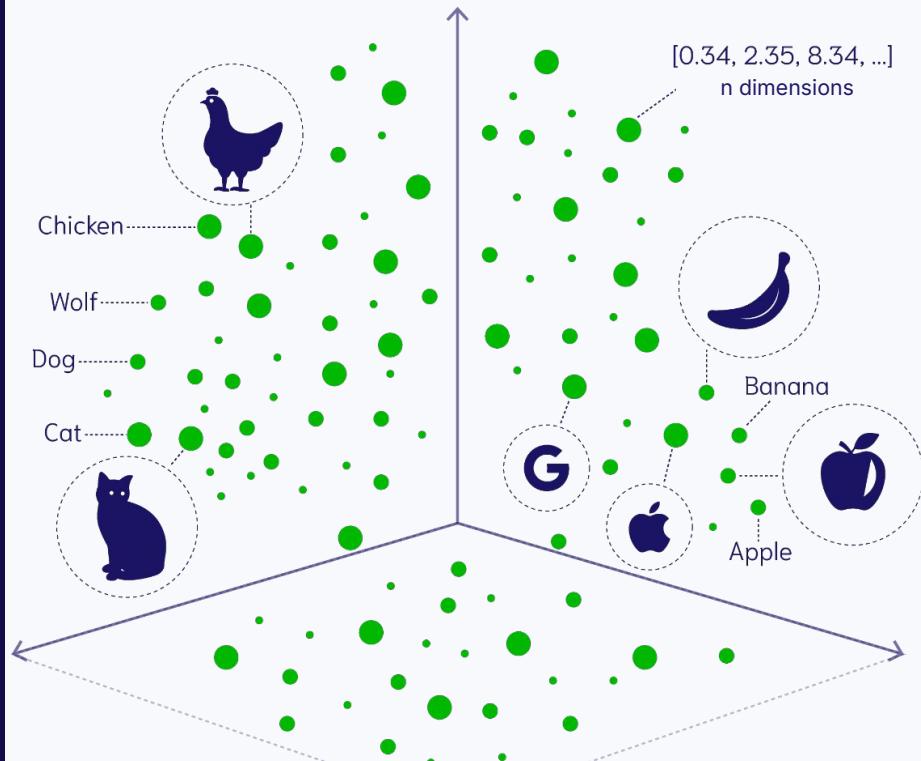


# Vector Search

How do airplanes fly?



[ 0.24, 0.36, ..., 0.16, 0.46 ]

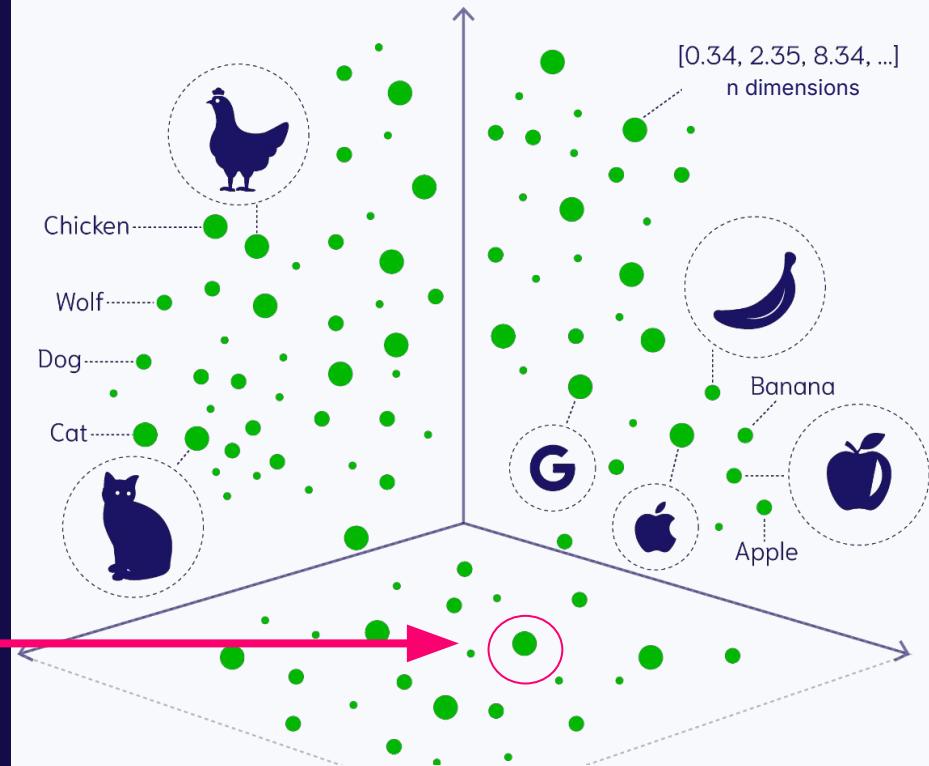


# Vector Search

How do airplanes fly?

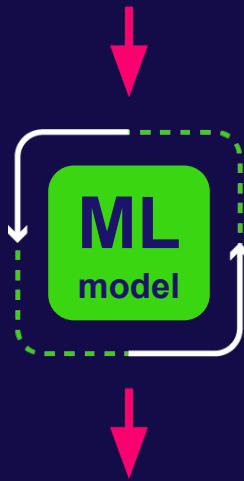


[ 0.24, 0.36, ..., 0.16, 0.46 ]



# Vector Search

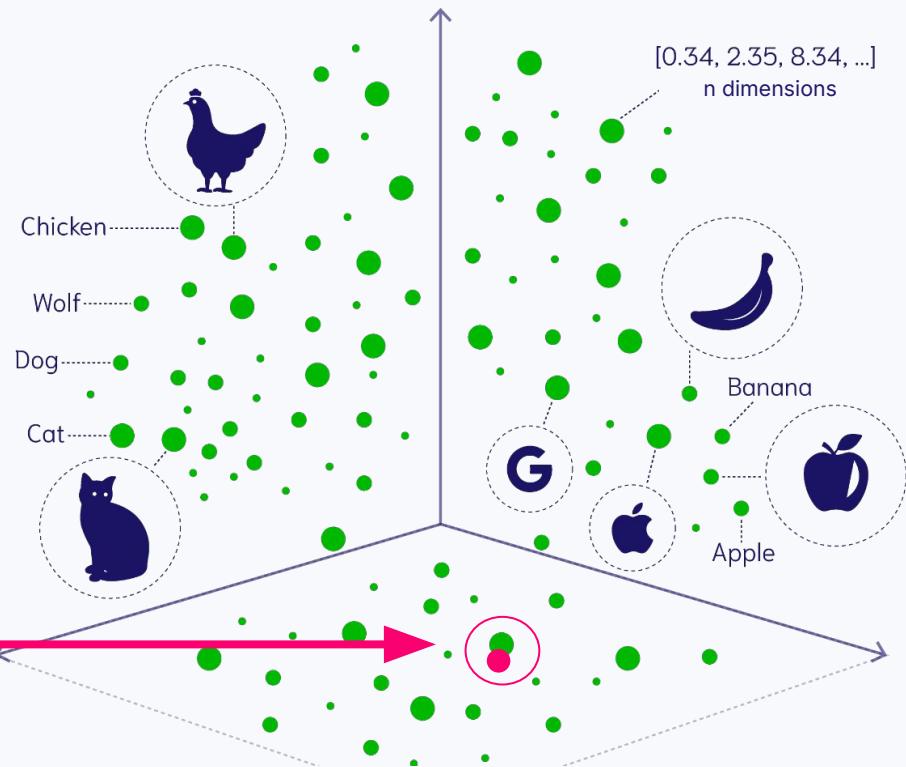
How do airplanes fly?



[ 0.24, 0.36, ..., 0.16, 0.46 ]

[ 0.26, 0.31, ..., 0.12, 0.44 ]

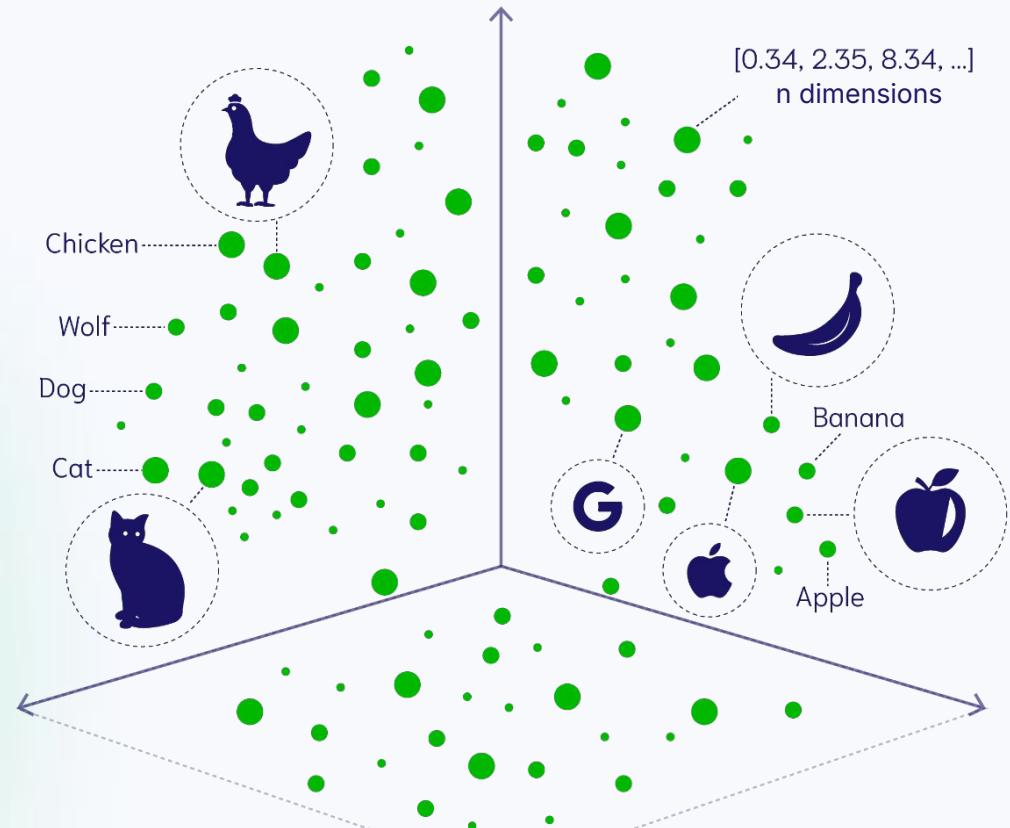
The NASA article



# What is a vector database?

A database whose primary data organization is by meaning (vectors)

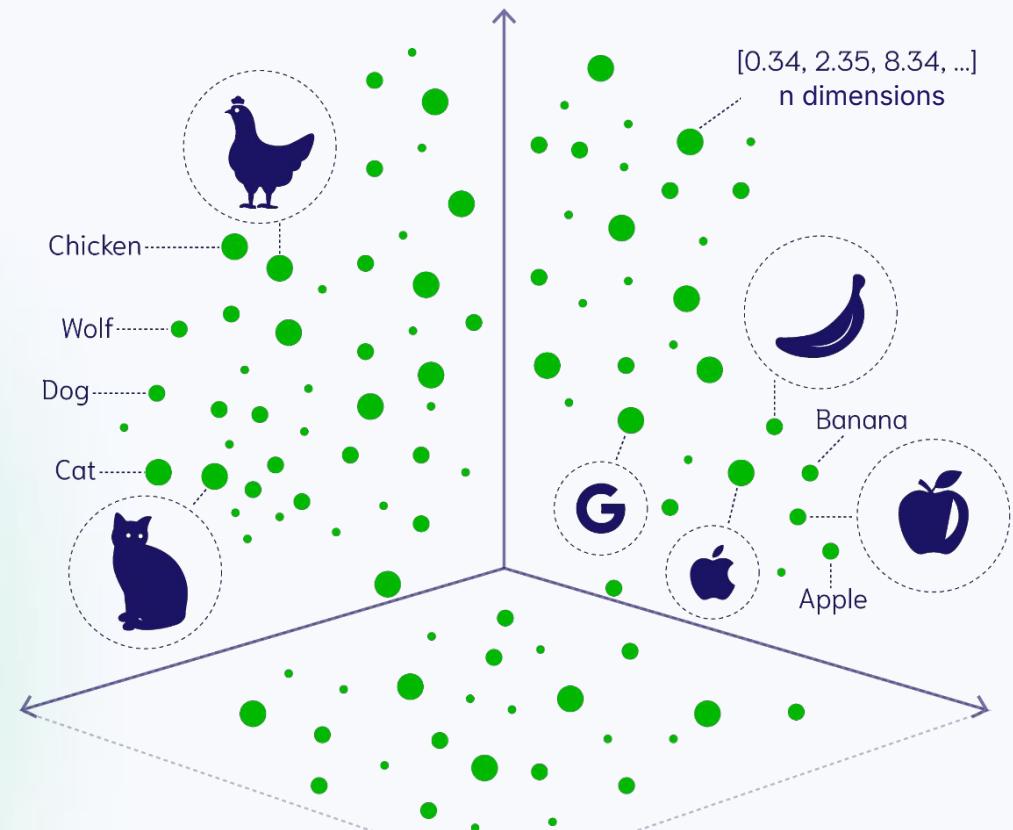
And optimized for data management and queries for using vectors **at scale.** (billions)



# What is an AI-native database?

## Vectors and AI-native workflow as first-class citizens

- AI model integrations
- Hybrid search
- Scalability (multi-tenancy, quantization)
- Developer experience





# Better search

Use semantic, keyword and hybrid searches  
to get the right data, faster



# Keyword, Vector or Hybrid?

Differences between search types



# Search: example

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

Ad for an airline

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

Explains aerodynamics

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

Advantages of flight

## Query:

"How do airplanes fly?"

# Search: example

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

## Ad for an airline

## Explains aerodynamics

## Advantages of flight

## Query:

"How do airplanes fly?"

how are they compared / searched?

# Keyword search (BM25)

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

Ad for an airline

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

Explains aerodynamics

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances.

Advantages of flight

## Query:

"How do airplanes fly?"

## **Keyword search:**

match keywords and calculate a score based on counts & typical frequency (TF-IDF)

# Keyword search (BM25)

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

*Ad for an airline*

*Explains aerodynamics*

*Advantages of flight*

## Query:

"How do airplanes fly?"

score: 4.8  
(for example)

## *Keyword search:*

match keywords and calculate a score based on counts & typical frequency (TF-IDF)



# Keyword search results

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

Ad for an airline

Rank: 1

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

Explains aerodynamics

Not in results

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances.

Advantages of flight

Rank: 2

## Query:

"How do airplanes fly?"

finds documents with the most important keyword matches

## Keyword search:

match keywords and calculate a score based on counts & typical frequency (TF-IDF)



# When to use...

Keyword search

When exact matches important.

Examples:

- Domain-specific data (jargon)
- Product / part numbers
- URL matches



# Keyword search syntax



```
1 products = client.collections.get("Product")
2
3 response = products.query.bm25(
4     "How do airplanes fly?",
5     limit=2
6 )
```



# Vector search

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

Ad for an airline

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

Explains aerodynamics

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances.

Advantages of flight

## Query:

"How do airplanes fly?"

## **Vector search:**

Convert each object into a vector and compare their numeric similarity



# Vector search

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

Ad for an airline

[ 0.26, 0.31, ..., 0.12, 0.44 ]

Explains aerodynamics

Advantages of flight

## Query:

"How do airplanes fly?"

[ 0.23, 0.45, ..., 0.84, 0.23 ]

similarity: 0.4  
(for example)

## Vector search:

Convert each object into a vector and compare their numeric similarity



# Vector search results

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

Ad for an airline

Rank: 3

Explains aerodynamics

Rank: 1

Advantages of flight

Rank: 2

## Query:

"How do airplanes fly?"

finds documents with the highest similarity in meaning

## Vector search:

Convert each object into a vector and compare their numeric similarity



# When to use...

Vector search

When robustness is important, or cross-lingual or cross-modal search is desired.

Examples:

- Descriptions
- Image to text / image to image
- Product reviews



# Vector search syntax

```
1 products = client.collections.get("Product")
2
3 response = products.query.near_text(
4     "How do airplanes fly?",
5     limit=2
6 )
```



# Hybrid search

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

Ad for an airline

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

Explains aerodynamics

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances.

Advantages of flight

## Query:

"How do airplanes fly?"

## Hybrid search:

Perform both a vector search and keyword search & combine the results for a final ranking



# Hybrid search

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

Ad for an airline

vector search similarity: 0.4  
keyword search score: 4.8

Explains aerodynamics

Advantages of flight

## Query:

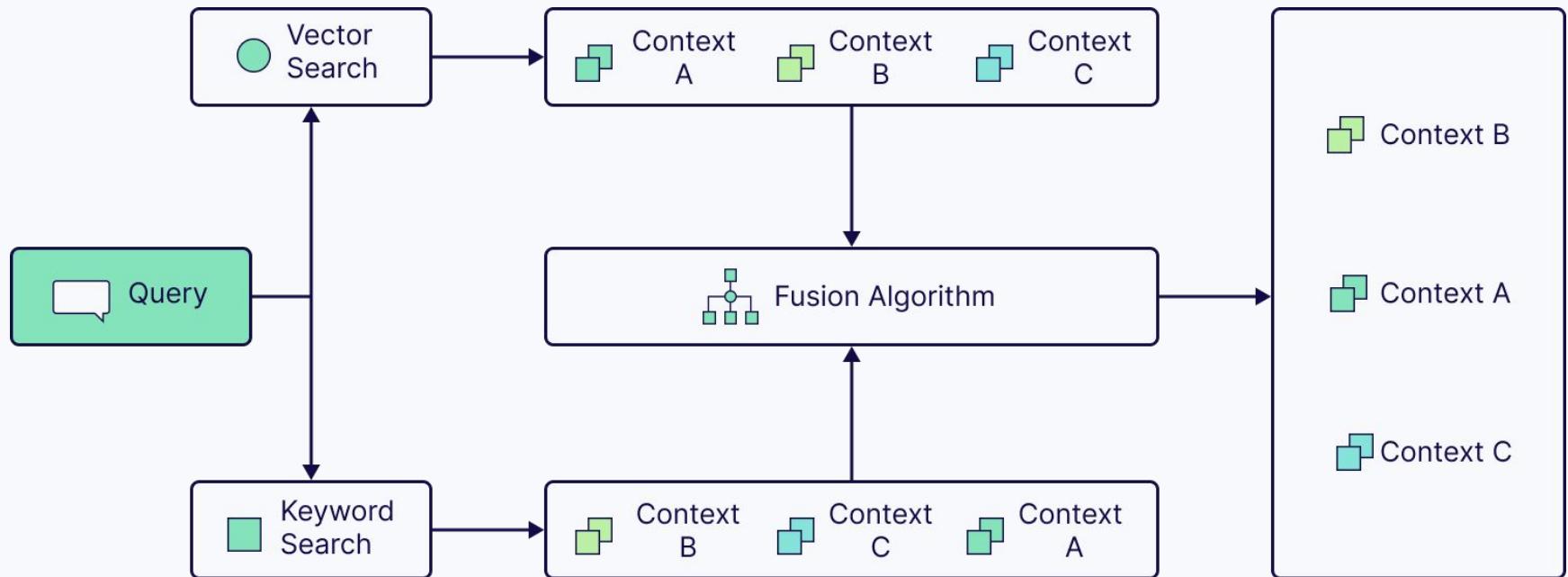
"How do airplanes fly?"

Calculate overall score & rank against others

## Hybrid search:

Perform both a vector search and keyword search & combine the results for a final ranking

# Hybrid Search





# Hybrid search results

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

*Ad for an airline*

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

*Explains aerodynamics*

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances.

*Advantages of flight*

Rank: 3

## Query:

"How do airplanes fly?"

mediocre results in both are not rewarded

## Hybrid search:

Perform both a vector search and keyword search & combine the results for a final ranking

# Hybrid search results

## Source documents

"Come and fly with ExpensiveAir's modern fleet of airplanes to over 4000 destinations. How to book a flight with ExpensiveAir?..."

"Airplane wings are shaped to make air move faster over the top of the wing. When air moves faster, the pressure of the air decreases."

"Flying is a great way to travel long distances quickly. Airplanes are able to carry large numbers of passengers and cargo over long distances."

Ad for an airline

Rank: 1/2

Explains aerodynamics

Rank: 1/2

Advantages of flight

Rank: 3

## Query:

"How do airplanes fly?"

hybrid results reward high rankings in at least one of keyword/vector

mediocre results in both are not rewarded

## Hybrid search:

Perform both a vector search and keyword search & combine the results for a final ranking



# When to use...

Hybrid search

Best-of-both-worlds search. Surfaces top results from both keyword & vector searches

Good results reported “in the wild” from users due to resilience.



# Hybrid search syntax

```
1 products = client.collections.get("Product")
2
3 response = products.query.hybrid(
4     "How do airplanes fly?",
5     limit=2,
6     alpha=0.5
7 )
```



# Quick demos



# Vector search: in depth

A deeper dive

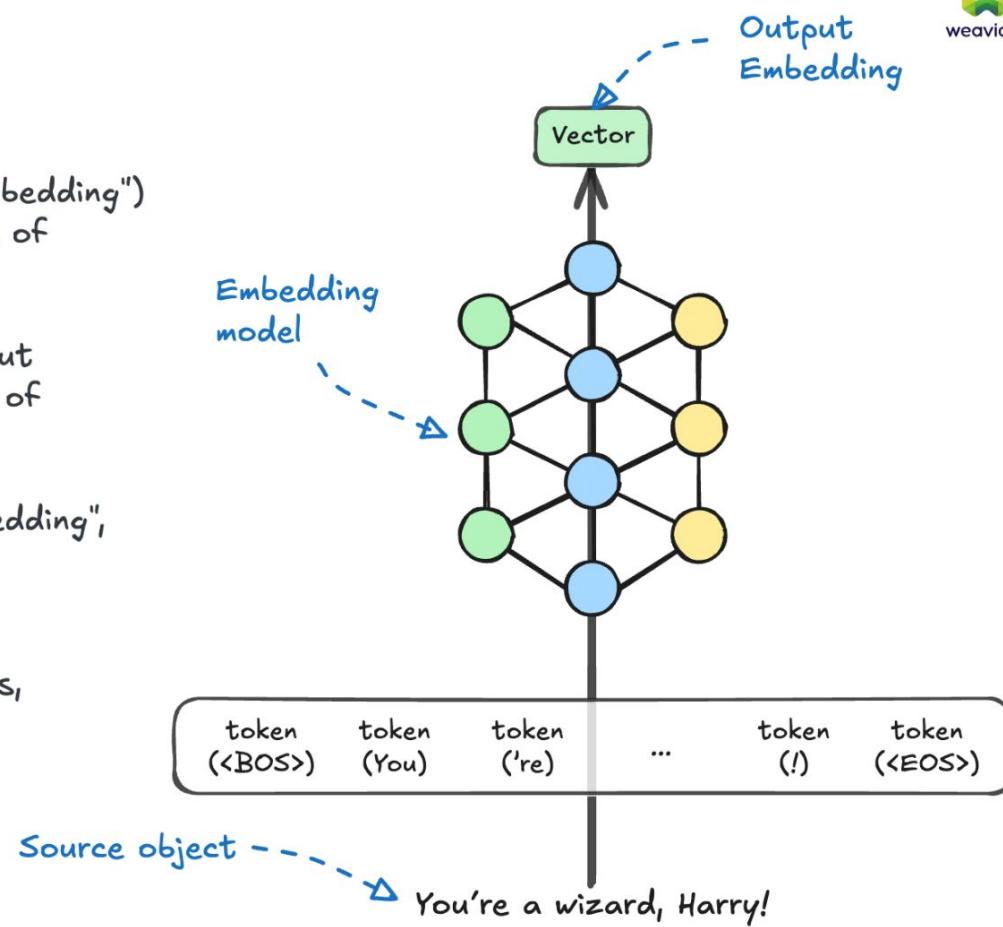
# How embeddings are created

A vector embedding (a.k.a. "vector" or "embedding") captures meaning of an object as a series of numbers.

An embedding is created by passing an input into an "embedding model", a specific form of a deep learning or AI model.

The model then produces an output "embedding", which looks like this:  
[0.0134, 0.8723, -0.4532, ..., 0.5892]

The more similar the meaning of the inputs, the more similar these embeddings will be.



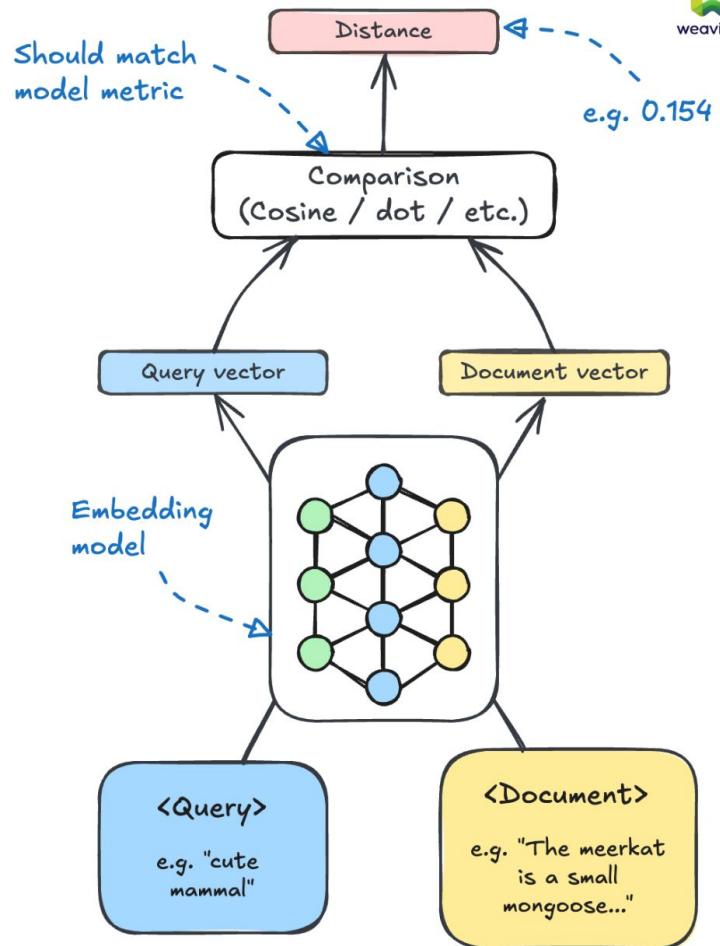
## How embeddings are compared

Vector embeddings capture meanings numerically. As a result, two embeddings can be compared to produce a measure of similarity.

The result is often measured as a "distance", measured using the same distance metric that the embedding model is trained on.

A "distance" indicates dissimilarity in meaning.

Two objects with a larger distance are more dissimilar than two objects with a smaller distance.



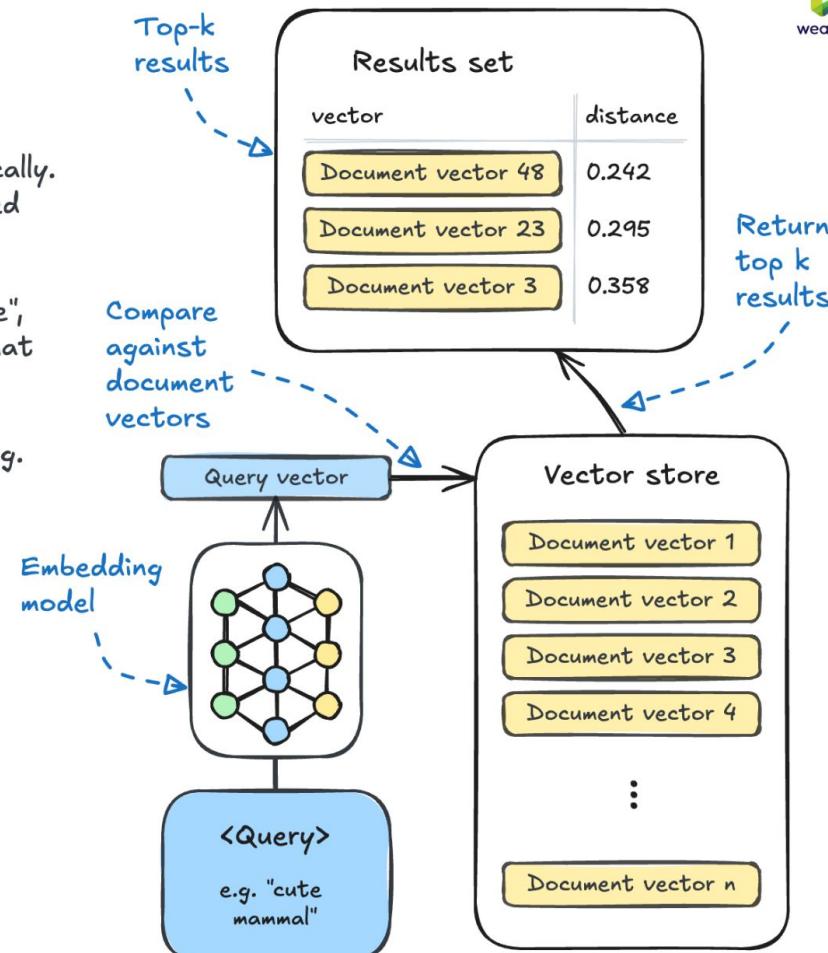
## How vector search works

Vector embeddings capture meanings numerically. As a result, two embeddings can be compared to produce a measure of similarity.

The result is often measured as a "distance", measured using the same distance metric that the embedding model is trained on.

A "distance" indicates dissimilarity in meaning.

Two objects with a larger distance are more dissimilar than two objects with a smaller distance.





# How do we search at scale?

What happens as our dataset grows larger?

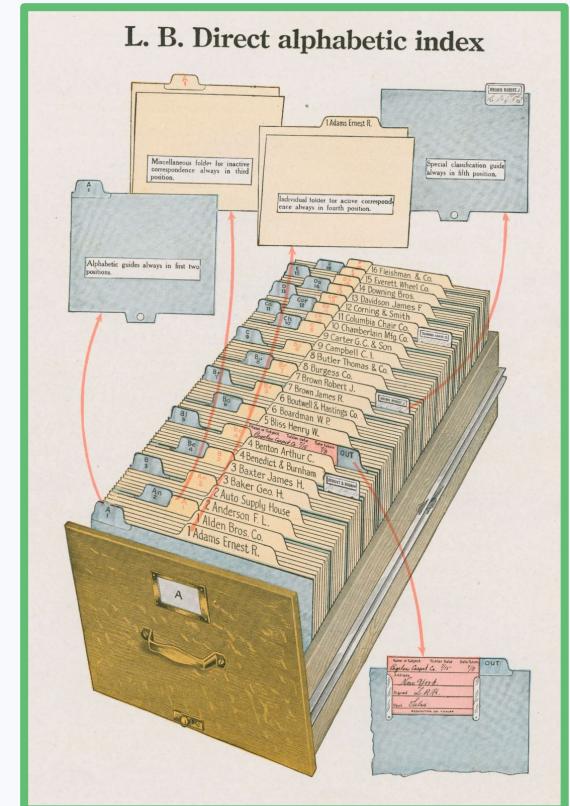
- How long does vector search take at:
  - 10,000 objects?
  - 10,000,000 objects?
  - 10,000,000,000 objects?
- And what are the typical bottlenecks?



# Indexes

Used by databases for speed

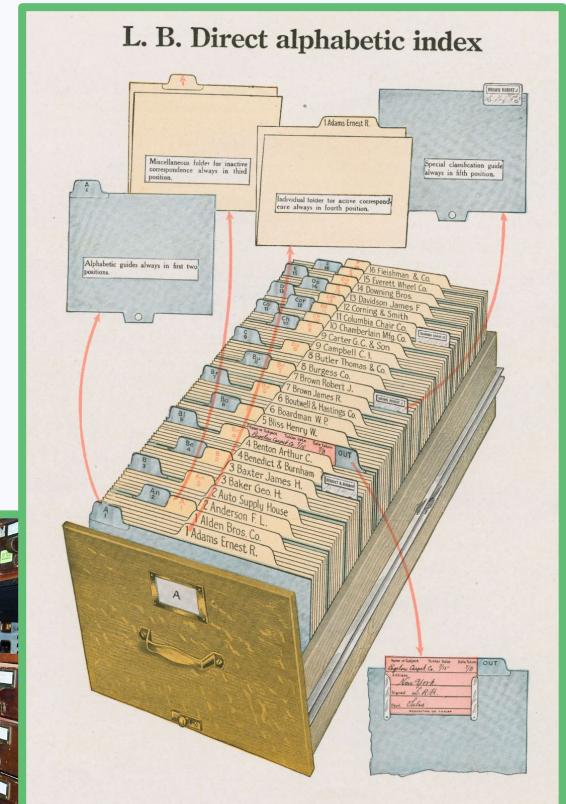
- e.g. Library catalogues



# Indexes

In databases:

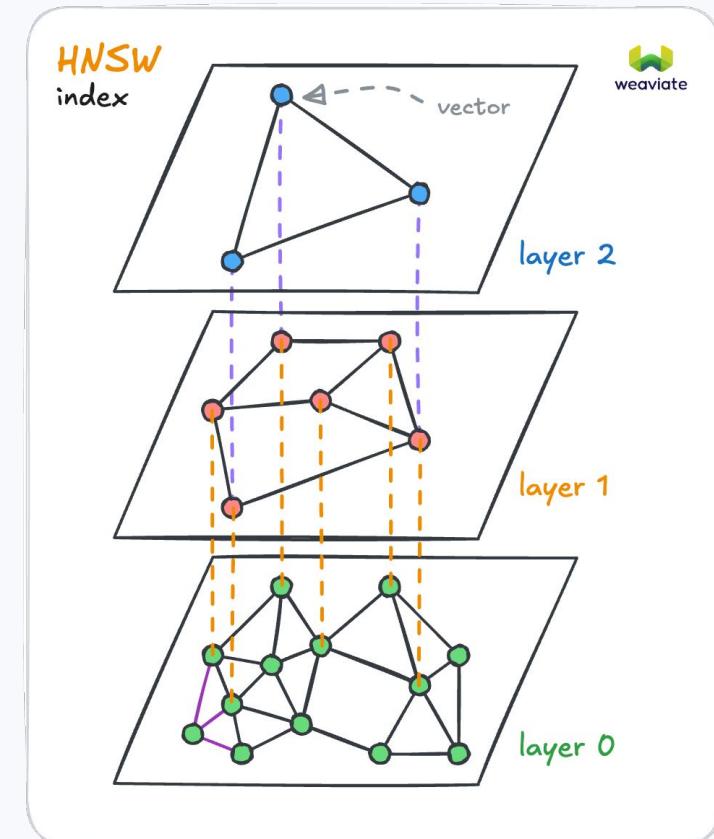
- Catalogue data
- Speed up search & filtering



# Vector indexing (ANN)

A way to scale search up to billions of vectors.

Most common index type:

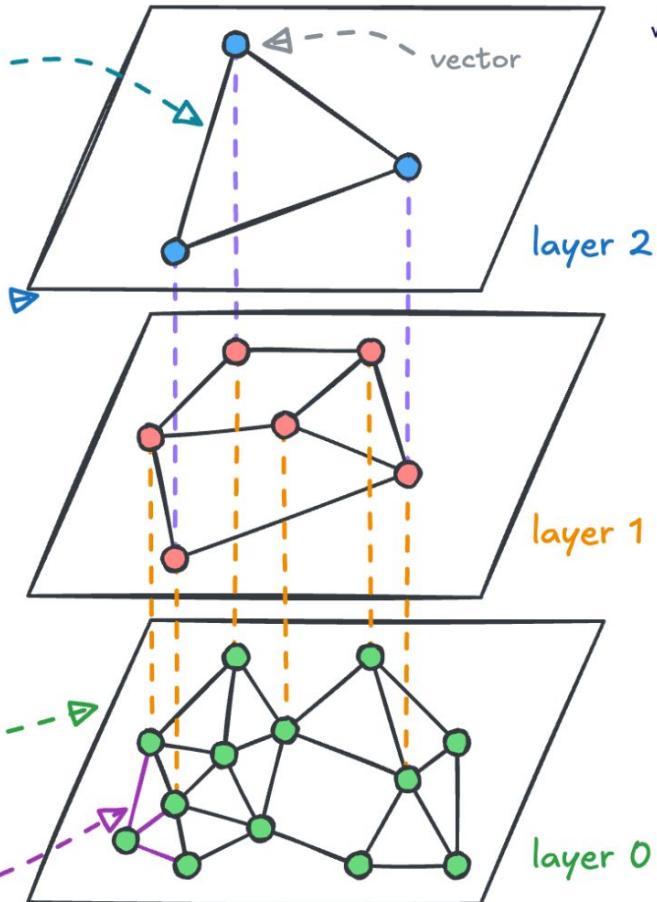


# Vector indexes explained: HNSW

An HNSW (hierarchical navigable small world) index enables **fast vector searches** even at large sizes.

The HNSW algorithm builds a **multi-layer graph** that balances speed & search quality.

higher layers:  
exponentially  
fewer nodes



Upper layers enable fast traversal,  
and lower layers enable detailed searching.

(Bi)directional connections,  
enabling "neighborhood"-based  
searches.

# Vector indexes explained: HNSW search

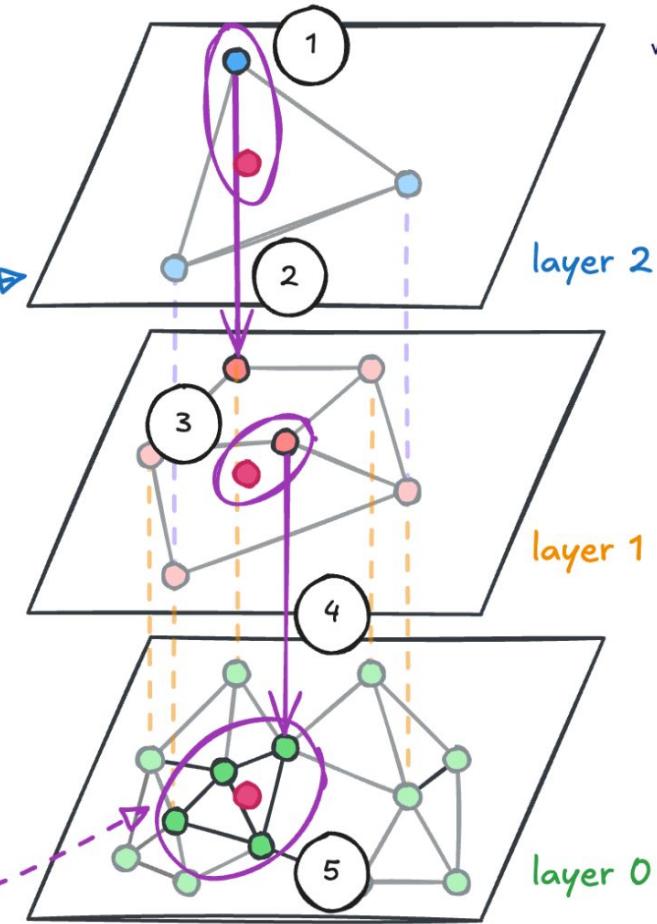
Vector search in HNSW proceeds by starting at the top layer, and proceeding down.

How would HNSW search for this vector: ● work?

Starting from the top layer,

- 1: find the most similar node,
- 2: travel down
- 3-4: repeat until bottom layer
- 5: then, find the closest k nodes

Search results - - -





# Where embeddings come from

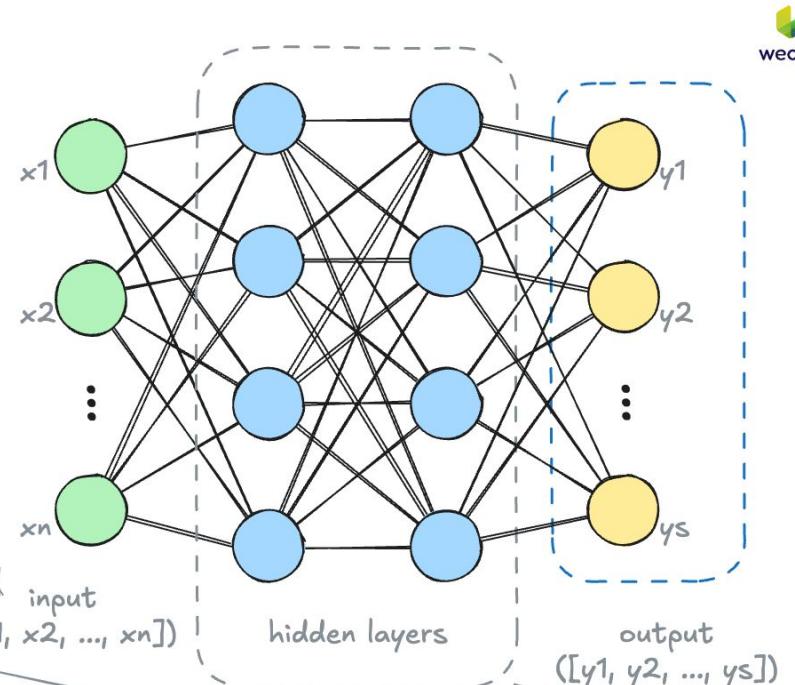
# Vectoriser: input → vectors

## Embedding models

Embedding models take many (e.g. thousands of) tokens as an input

And output a fixed shape ( $s$ ) of numbers.

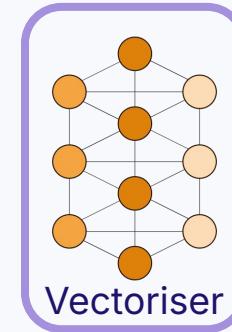
This is called an "**embedding**"



# Vectoriser: input ➔ vectors

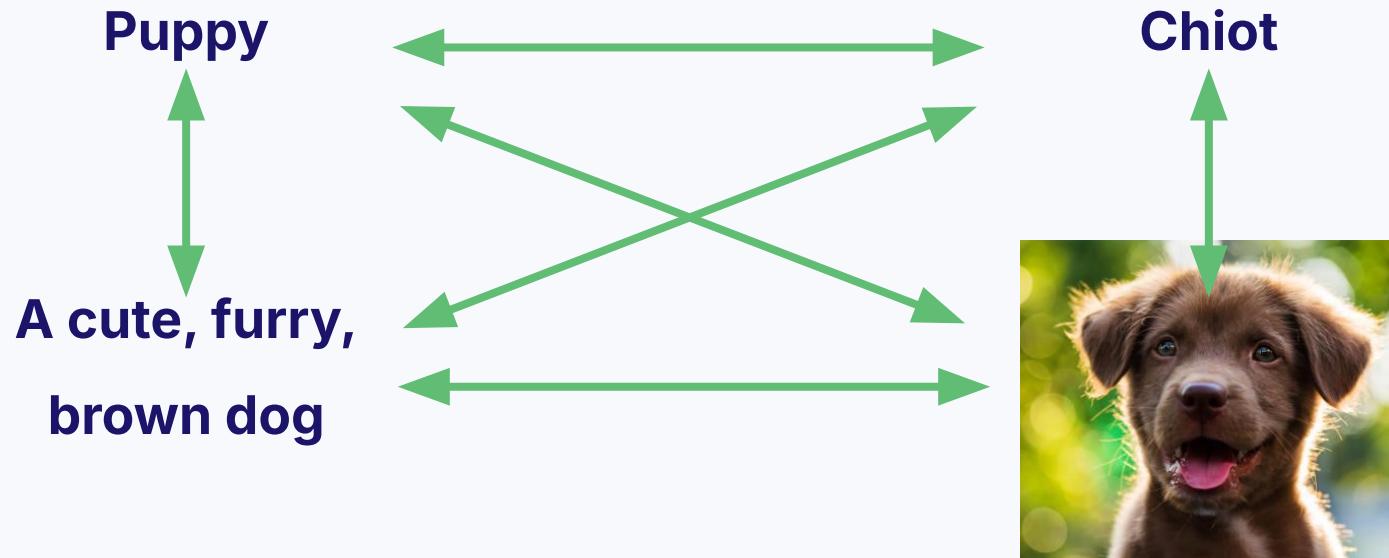
Hundreds of models are available:

- Proprietary models
  - Via Cohere, OpenAI, Google, AWS, etc.
- Open-source models
  - Via Hugging Face
- Why so many?



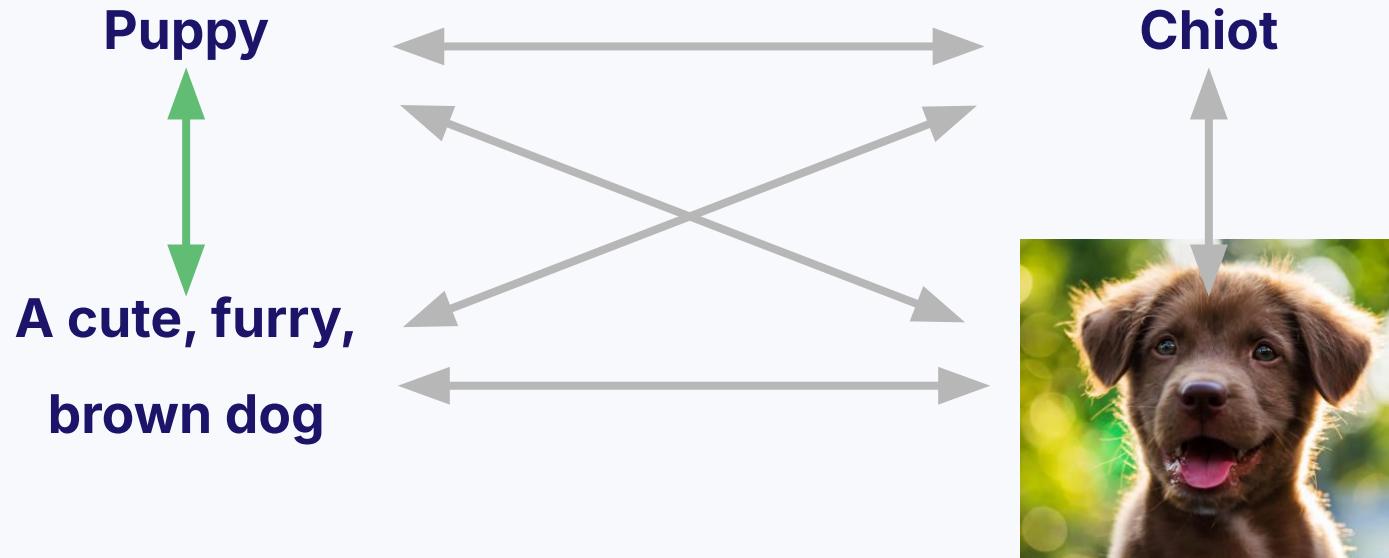


# Which two are the \*most\* similar?



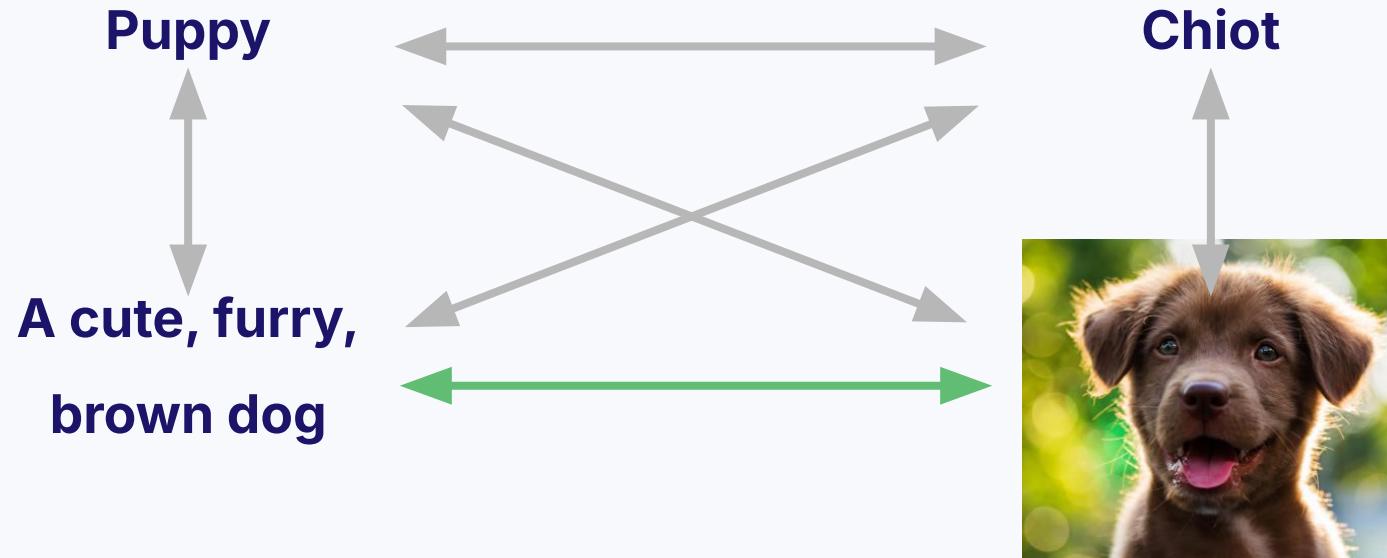


# The only English descriptions



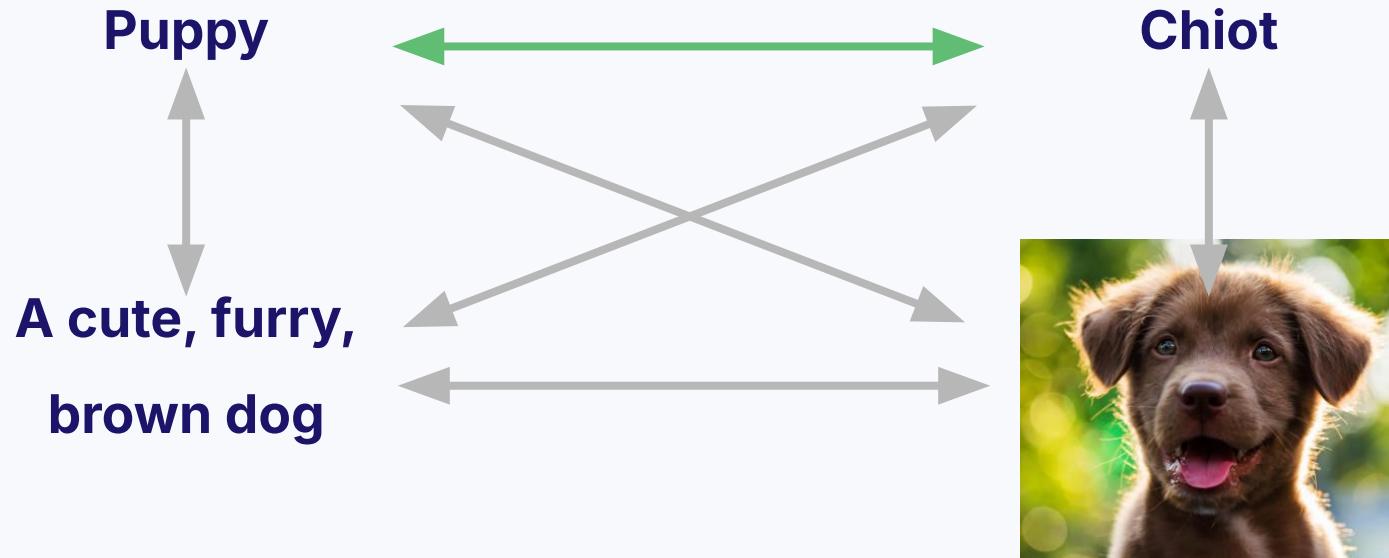


# Best matching image to text





# If you speak English & French





# Which two are the **\*most\*** similar?

Puppy

Chiot

Depends on your perspective & goals  
A cute, fury,  
brown dog





# Which two are the **\*most\*** similar?

Puppy

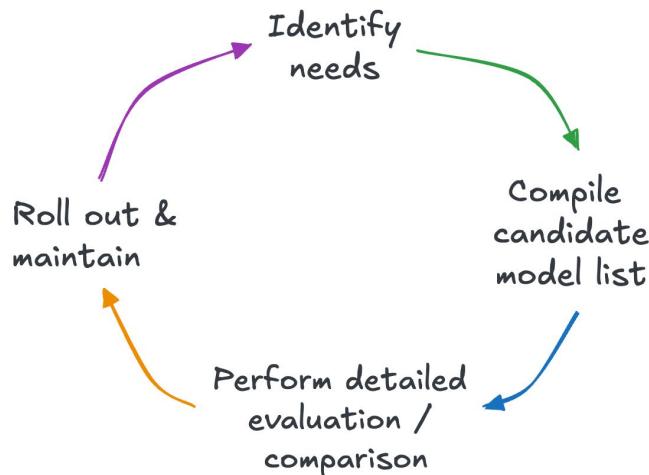
Chiot

Is determined by the **vectoriser model**  
A cute, furry,  
brown dog



# Vectorizer model selection

## Embedding models: Selection workflow



# Vectorizer model selection

Embedding model selection: identify your needs



## Data characteristics

- Modality
- Language
- Length
- Domain
- Asymmetry



## Performance needs

- Accuracy
- Latency
- Throughput
- Volume
- Task type



## Operational Factors

- Hardware limitations
- Org capabilities
- API Rate limits
- Deployment



## Business Requirements

- Hosting requirements
- Licensing
- Privacy
- Budget
- Compliance



# Examples of model cards

**model name**

**model provider**

**license info**

**model size**

**benchmarks**

**input / output details**

The diagram illustrates four examples of model cards:

- microsoft/phi-4**: Shows the model card for the Microsoft phi-4 model. Annotations point to the "model name" (phi-4), "model provider" (Microsoft Research), and "license info" (MIT license).
- Cohere/Coher-embed-english-v3.0**: Shows the model card for the Cohere embed-english-v3.0 model. Annotations point to the "model size" (14.7B params) and "input / output details" (Usage Cohere API section).
- llama3.3**: Shows the model card for the llama3.3 70B model. Annotations point to the "model size" (70B parameters) and "input / output details" (Usage Cohere API section).
- benchmarks**: Shows a collection of benchmarks for the Arctic-embed model. Annotations point to the "benchmarks" section.

**Model Summary** (from Microsoft phi-4 card):

**Developers**: Microsoft Research

**Description**: phi-4 is a state-of-the-art open model built upon a blend of synthetic datasets, data from filtered public domain websites, and acquired academic books and Q&A datasets. The goal of this approach was to ensure that small capable models were trained with data focused on high quality and advanced reasoning.  
phi-4 underwent a rigorous enhancement and alignment process, incorporating both supervised fine-tuning and direct preference optimization

**Tools**: 70b  
Downloads last month: 508,640

**Community**: 41

**Model card**, **Files and versions**, **Community**

**Model Summary** (from Cohere embed-english-v3.0 card):

**Transformers**: mteb, **Eval Results**, **Inference Endpoints**

**Model card**, **Files and versions**, **Community**

This repository contains the tokenizer for the Cohere embed-english-v3.0 model. See our blogpost [Cohere Embed V3](#) for more details on this model.

You can use the embedding model either via the Cohere API, AWS SageMaker or in your private deployments.

**Usage Cohere API**

The following code snippet shows the usage of the Cohere API. Install the cohere SDK via: pip install cohere

**Inputs and outputs**

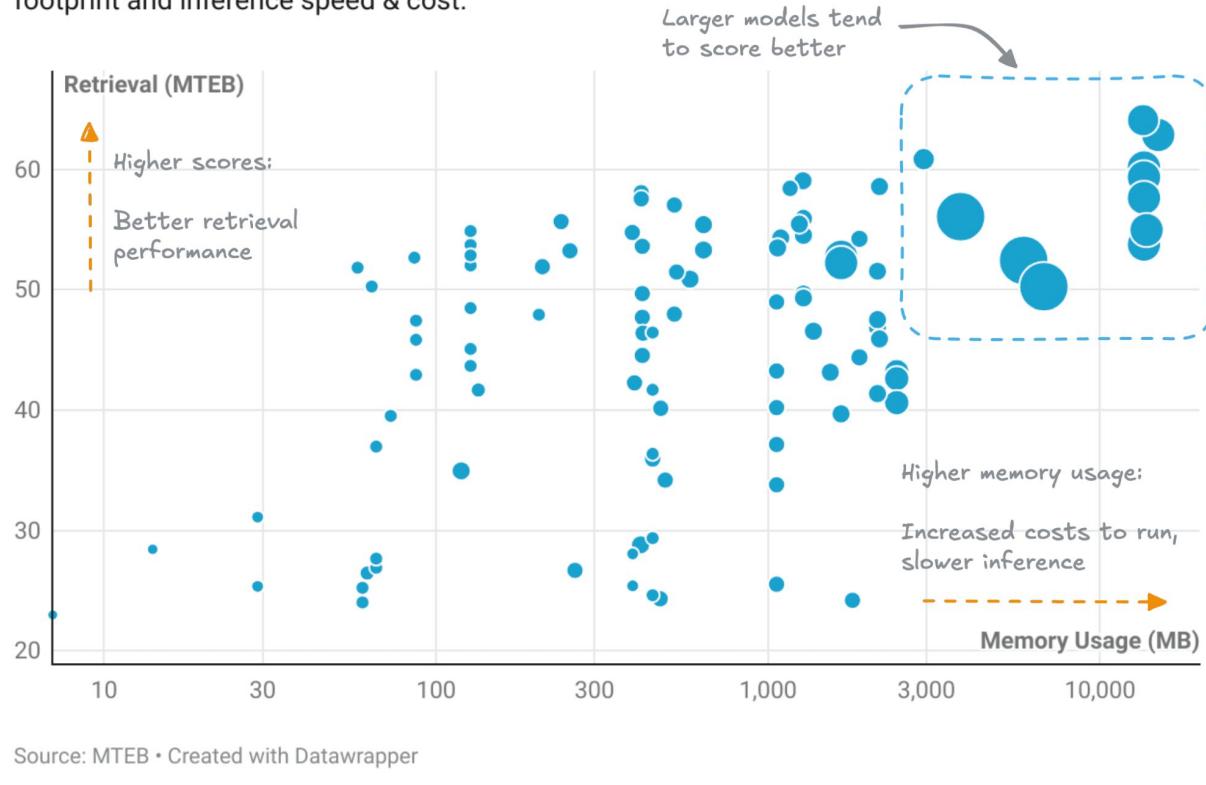
- Input:**
  - Text string, such as a question, a prompt, or a document to be summarized
  - Images, normalized to 896 x 896 resolution and encoded to 256 tokens each
  - Total input context of 128K tokens for the 4B, 12B, and 27B sizes, and 32K tokens for the 1B size
- Output:**
  - Generated text in response to the input, such as an answer to a question, analysis of image content, or a summary of a document
  - Total output context of 8192 tokens

**On this page**

- Model Information
- Description
- Inputs and outputs
- Citation
- Model Data
- Training Dataset
- Data Preprocessing
- Implementation Information
- Hardware
- Software
- Evaluation

# Vectorizer model selection

Larger embedding models generally perform better at retrieval, at a cost of increased memory footprint and inference speed & cost.





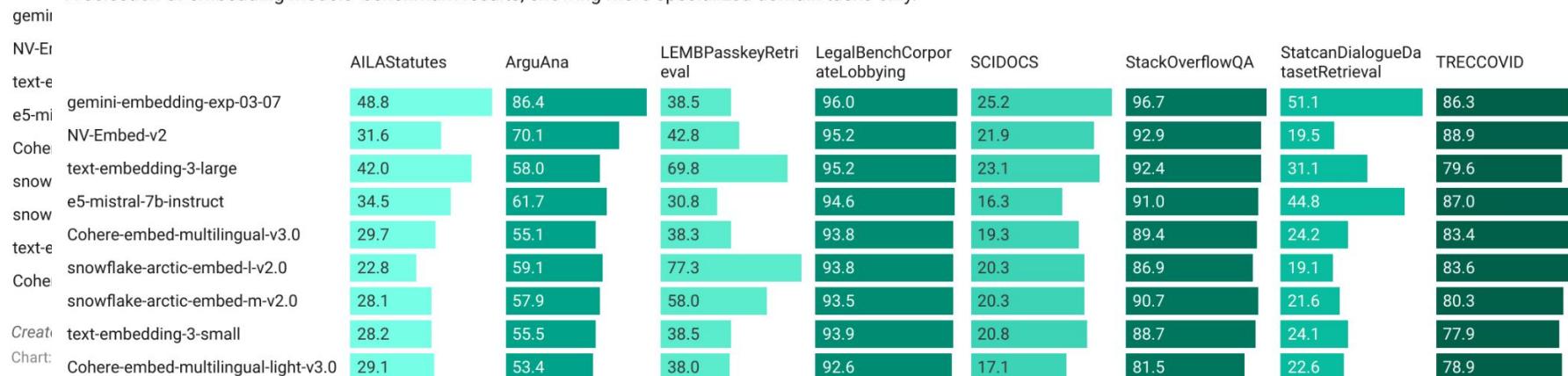
# Evaluate benchmarks

## MTEB Benchmarks by task (general)

A selection of embedding models' benchmark results, showing general domain tasks only.

## MTEB Benchmarks by task (specialized)

A selection of embedding models' benchmark results, showing more specialized domain tasks only.



Created for educational purposes to demonstrate embedding model selection.

Chart: JP Hwang @ Weaviate • Source: MTEB @ 2025.Apr.16 • Created with Datawrapper



# Detailed evaluation

Evaluate at system-level,  
based on your own

- Custom benchmarks
  - Curate/create own datasets
  - Define specific metrics
- Quantitative analysis
- Qualitative analysis
  - Apply domain expertise
  - Look for patterns
- End-to-end system evaluation
- Periodic re-evaluation



# AI-native workflows

About Weaviate's AI-native developer experience,  
integrating embedding and genAI models directly

## Model provider integrations

## Model provider integrations

	Model provider	Embeddings	Generative AI	Others
Anthropic	-	Text	-	
Anyscale	-	Text	-	
AWS	Text	Text		
Cohere	Text, Multimodal	Text	Reranker	
Databricks	Text	Text	-	
FriendliAI	-	Text	-	
Google	Text, Multimodal	Text	-	
Hugging Face (API)	Text	-	-	
Jina AI	Text, Multimodal	-	Reranker	
Mistral	Text	Text	-	
NVIDIA	Text	-	-	
OctoAI (Deprecated)	Text, Multimodal	-	Reranker	
OpenAI	Text	Text	-	
Azure OpenAI	Text, Multimodal	Text	Reranker	
Voyage AI	Text, Multimodal	-	Reranker	
Weaviate	Text	Text	-	
xAI	Text	Text	-	
GPT4All (locally hosted)	Text	Text	-	
KubeAI (locally hosted)	Text, Multimodal	-	Reranker	
Hugging Face (locally hosted)	Text, Multimodal	-	Reranker	
Meta ImageBind (locally hosted)	Text	-	-	
Ollama (locally hosted)	-	Text	-	

# Out-of-the-box integrations

## Locally hosted

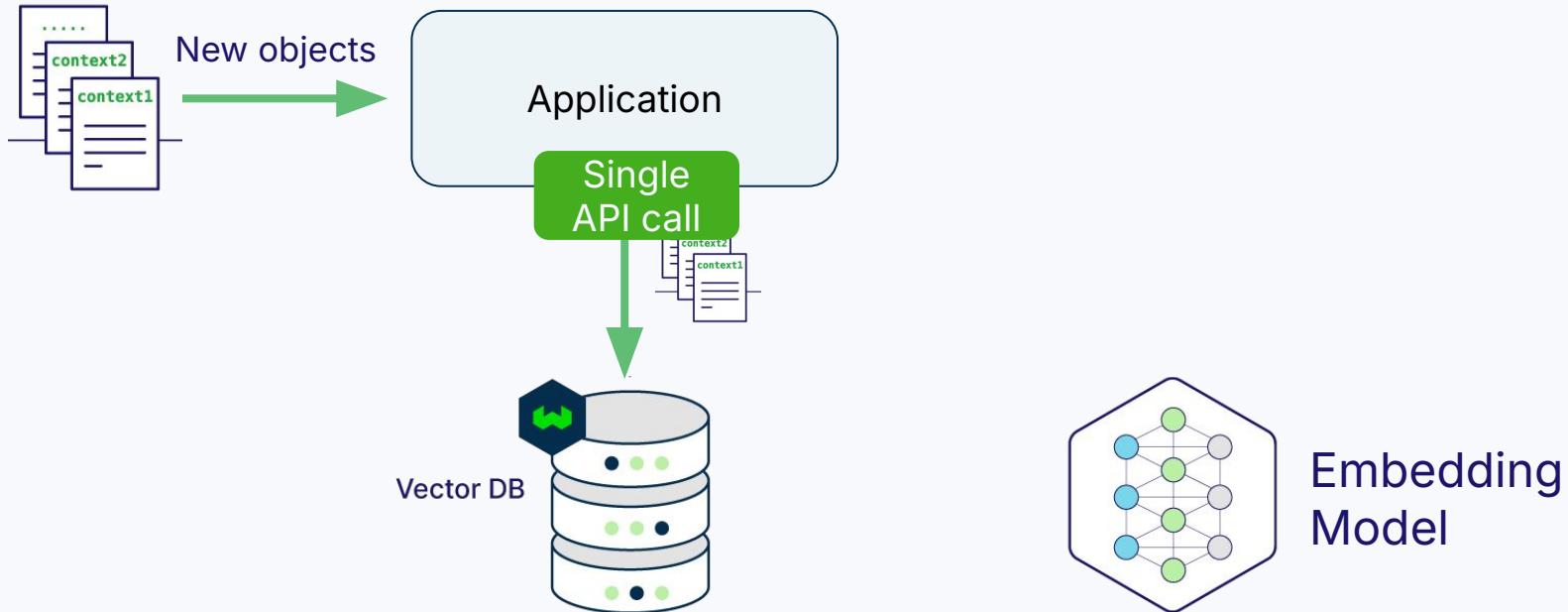
Model provider	Embeddings	Generative AI	Others
GPT4All	Text	-	-
Hugging Face	Text, Multimodal (CLIP)	-	Reranker
Meta ImageBind	Multimodal	-	-
Ollama	Text	Text	-



# Embedding model integration

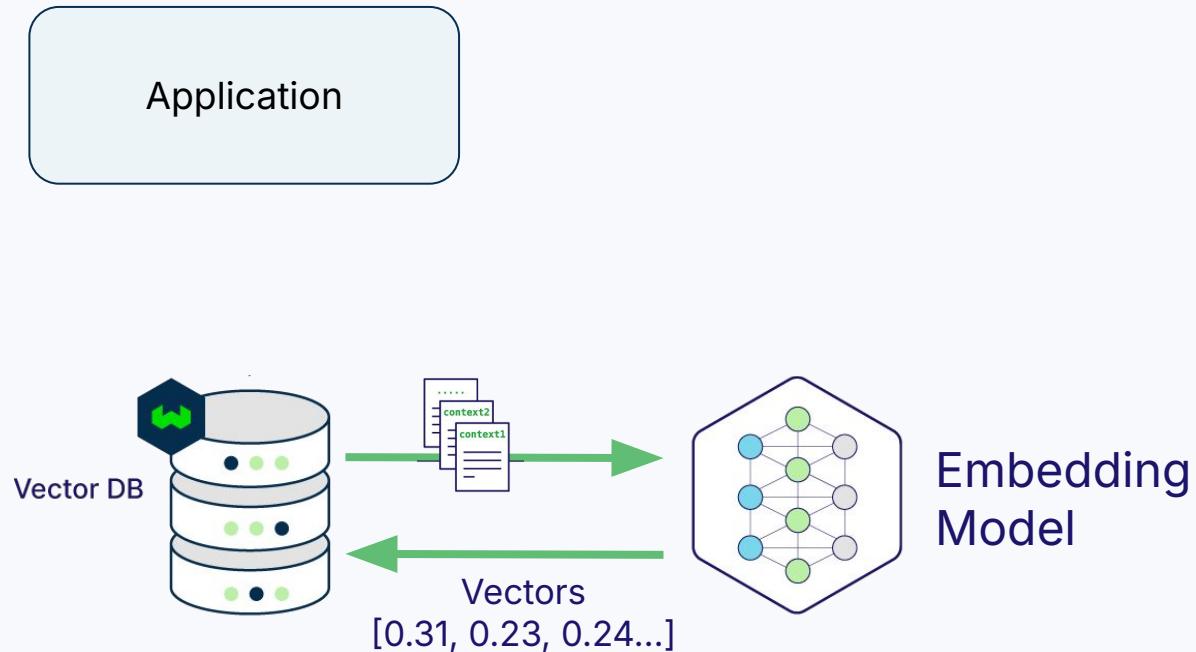


# Data operations (i.e. data import)



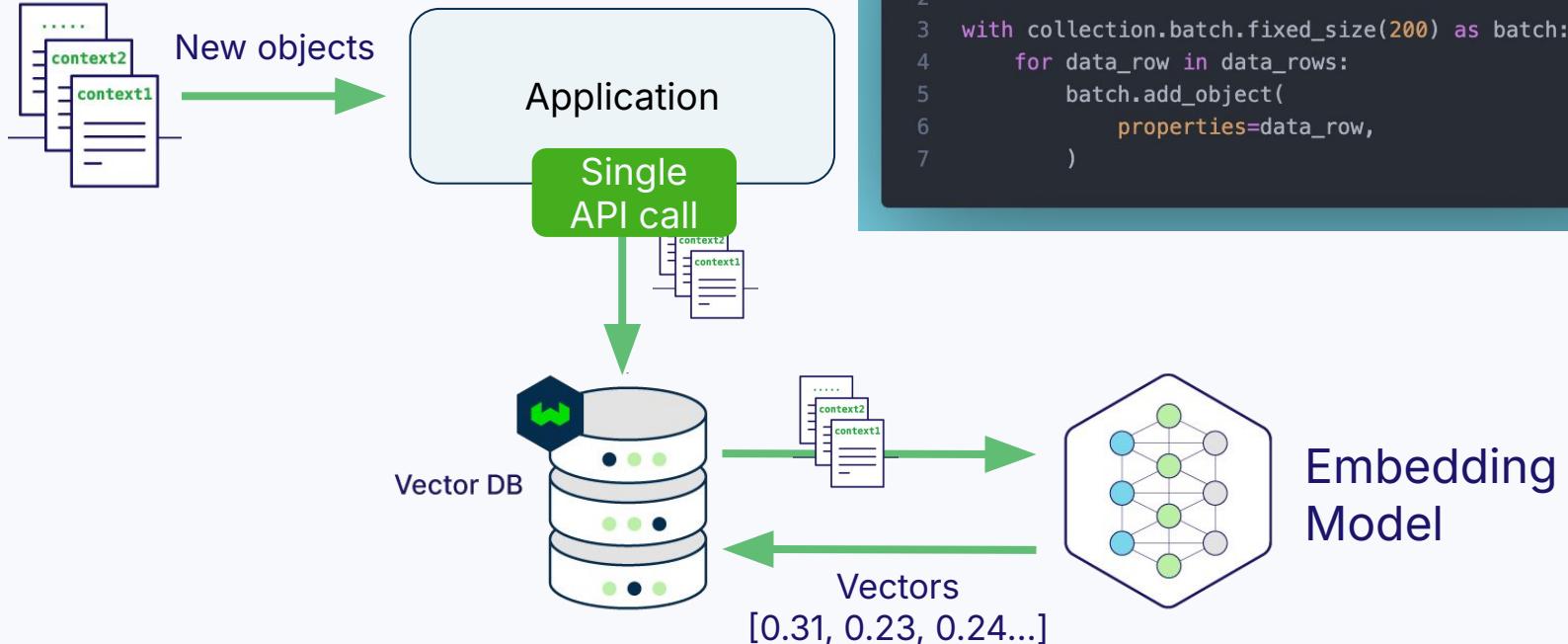


# Data operations (i.e. data import)

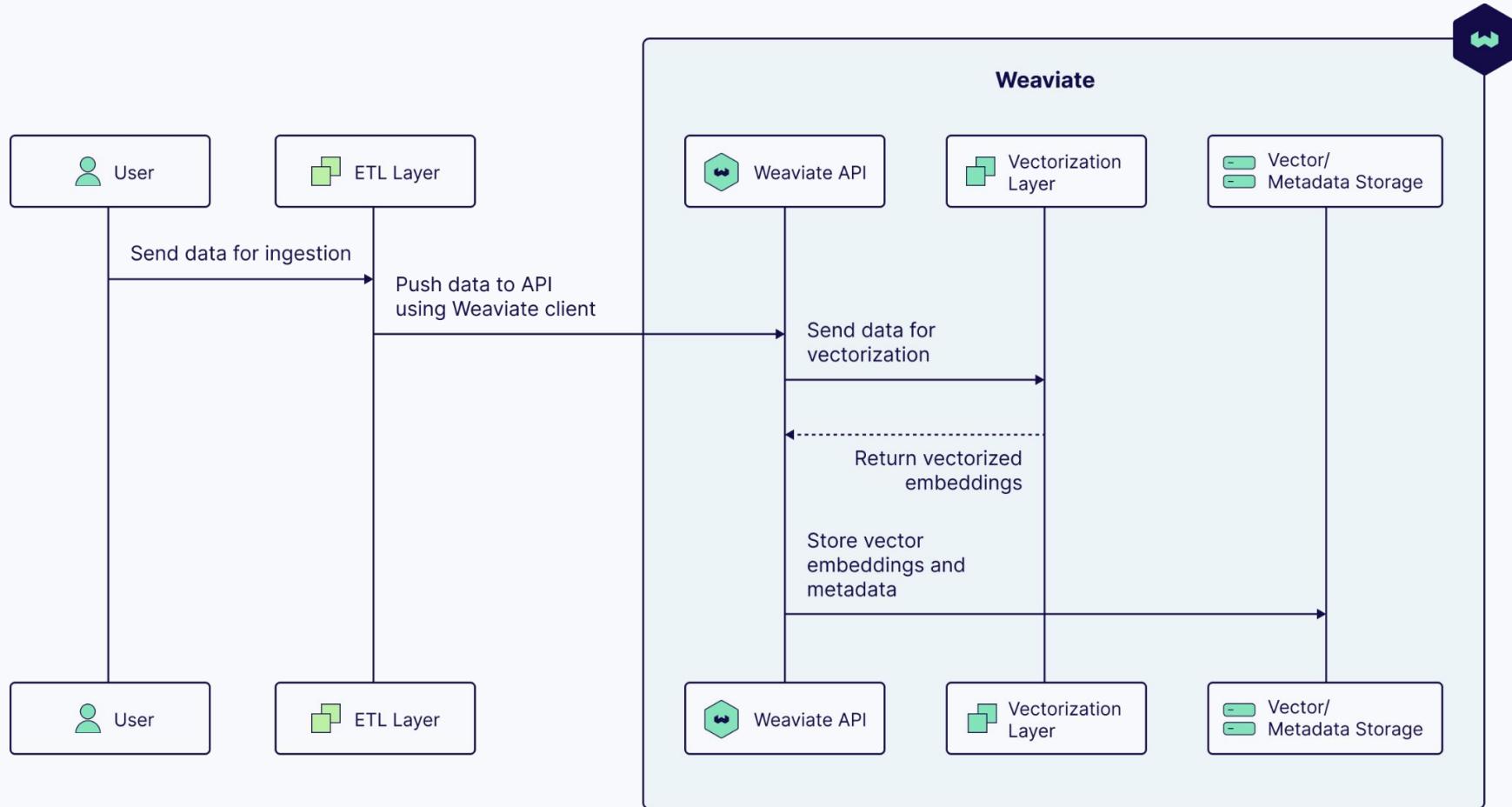




# Data operations (i.e. data import)

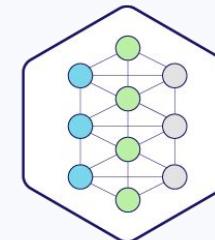
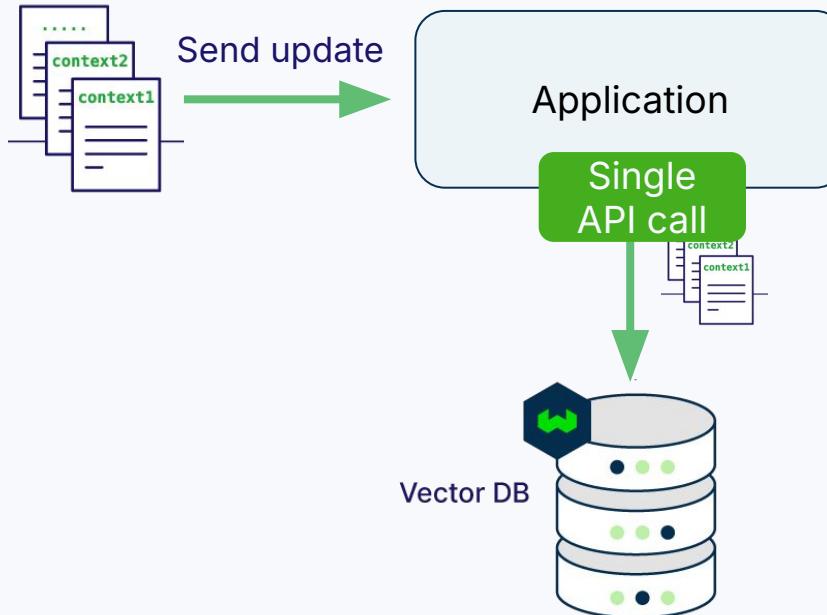


# Data Ingestion





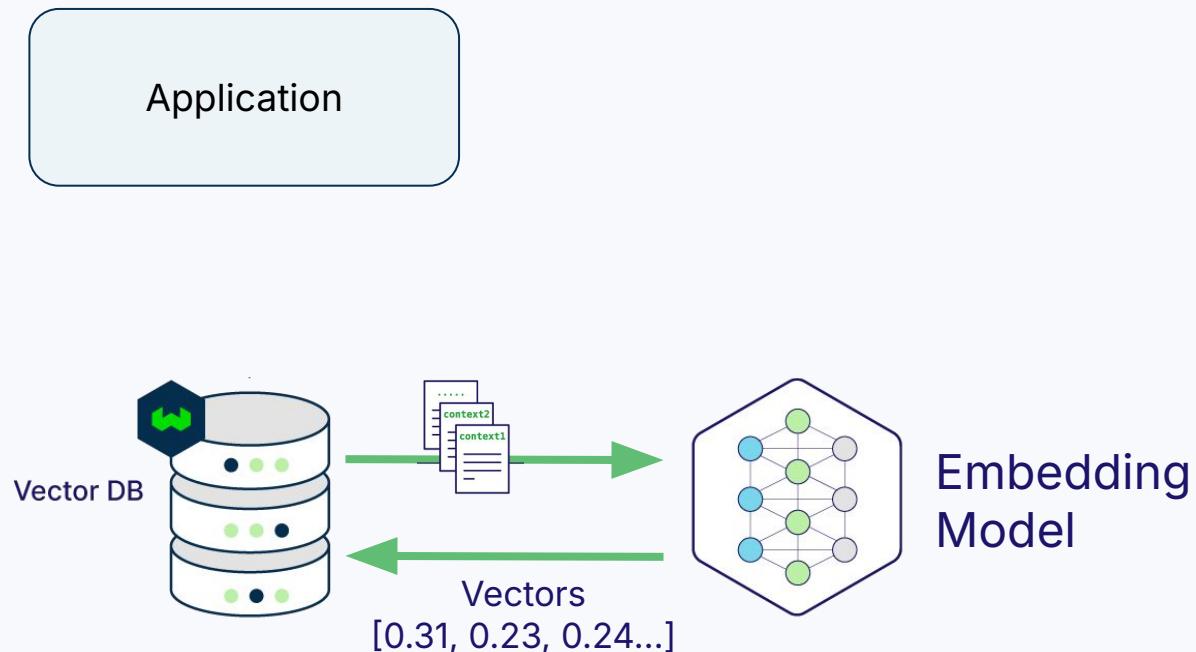
# Data operations (i.e. update)



Embedding  
Model

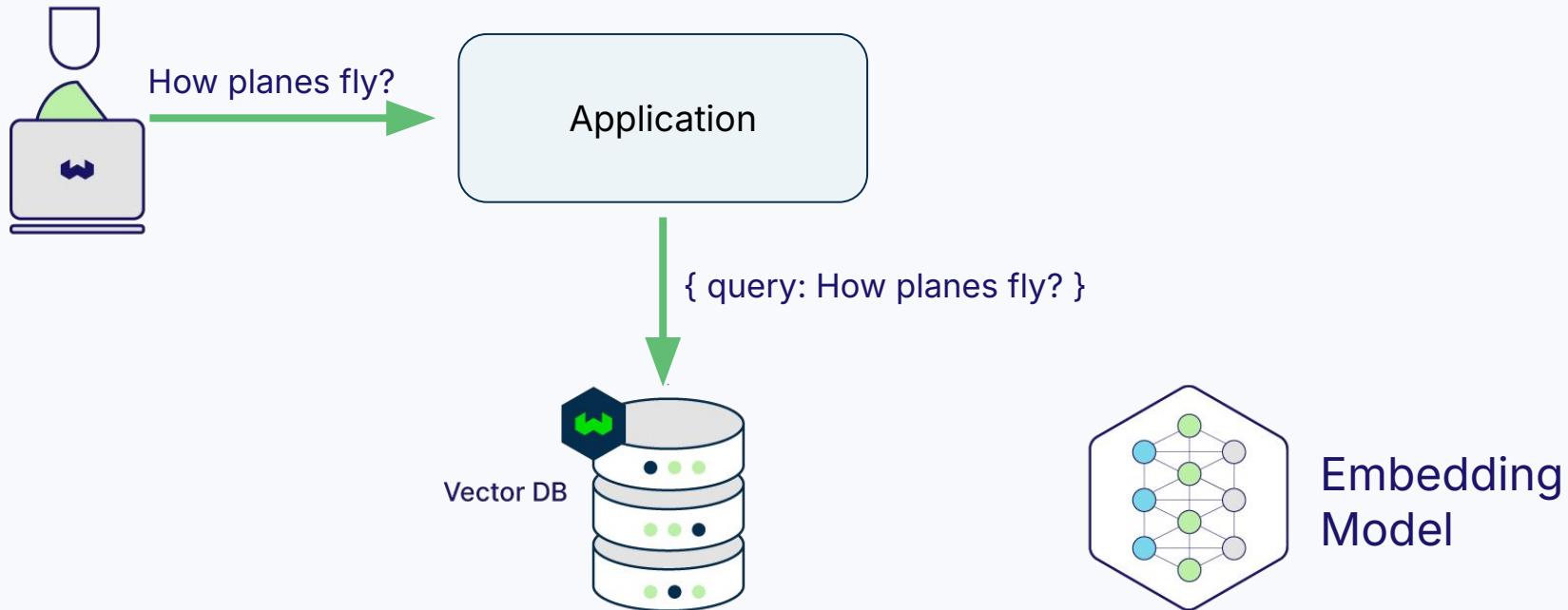


# Data operations (i.e. update)



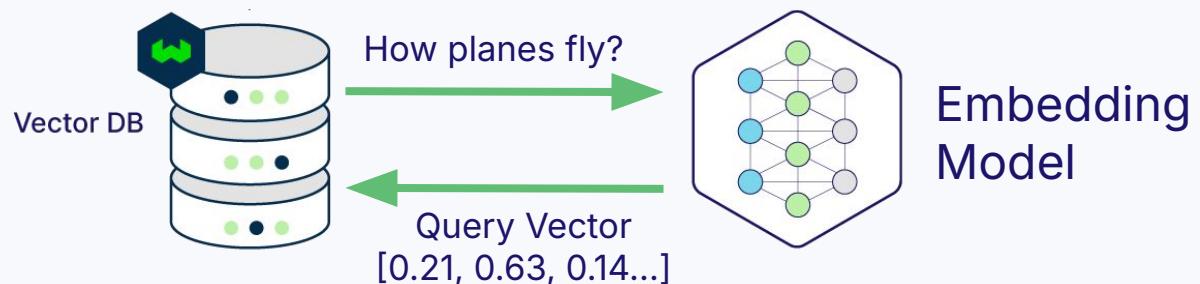
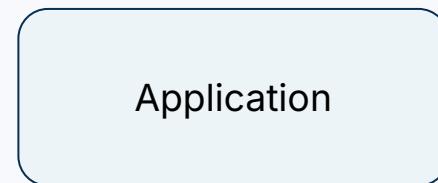


# Query



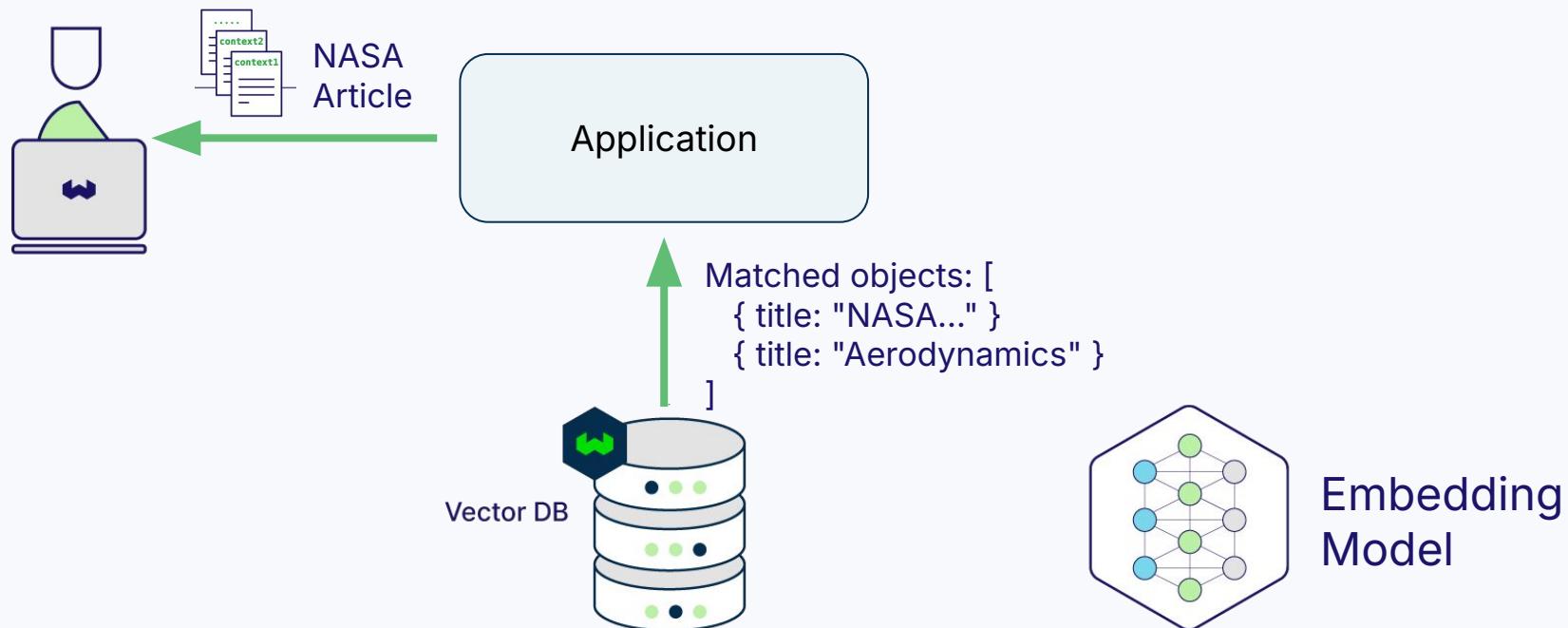


# Query



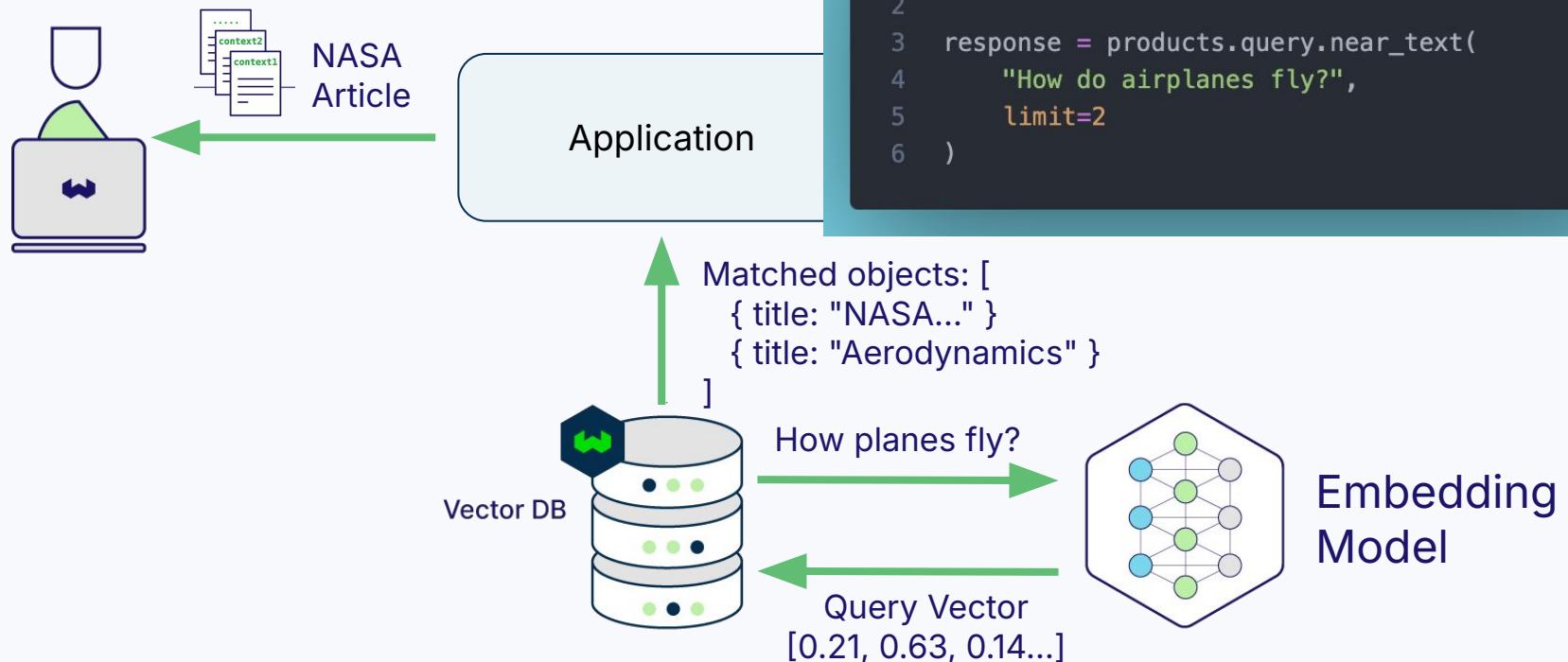


# Query





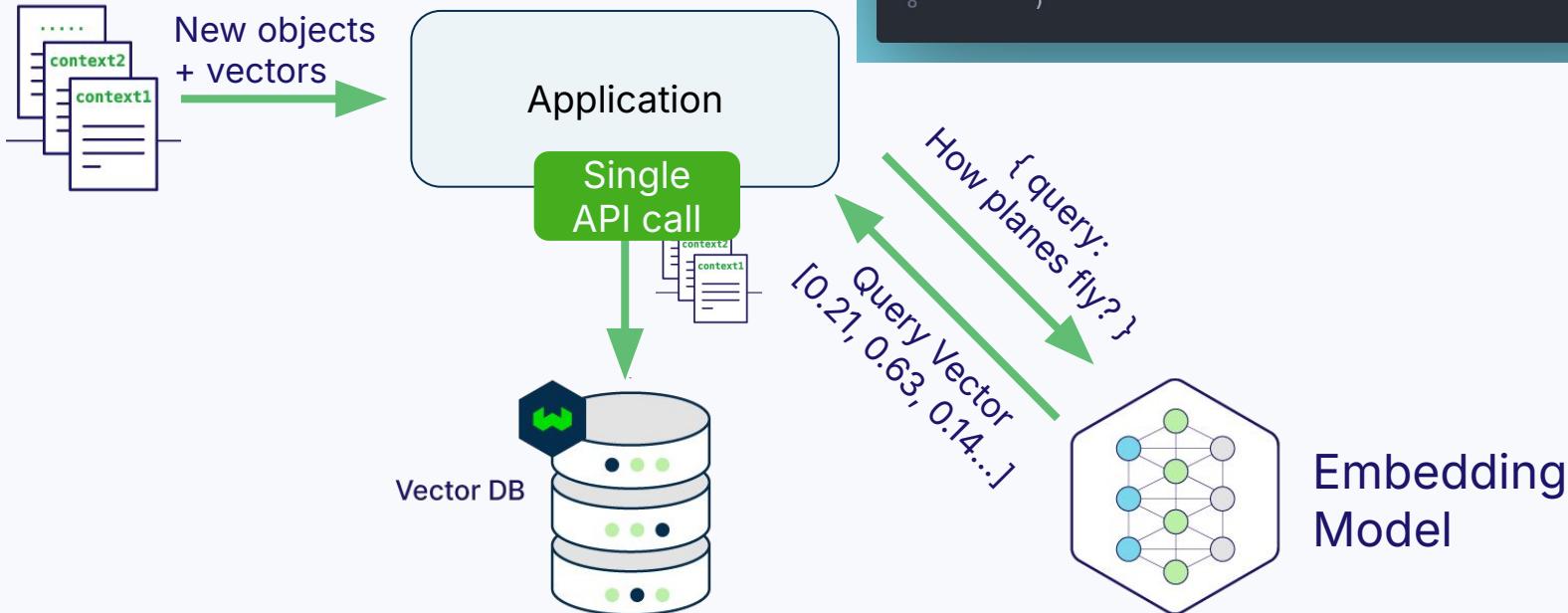
# Query





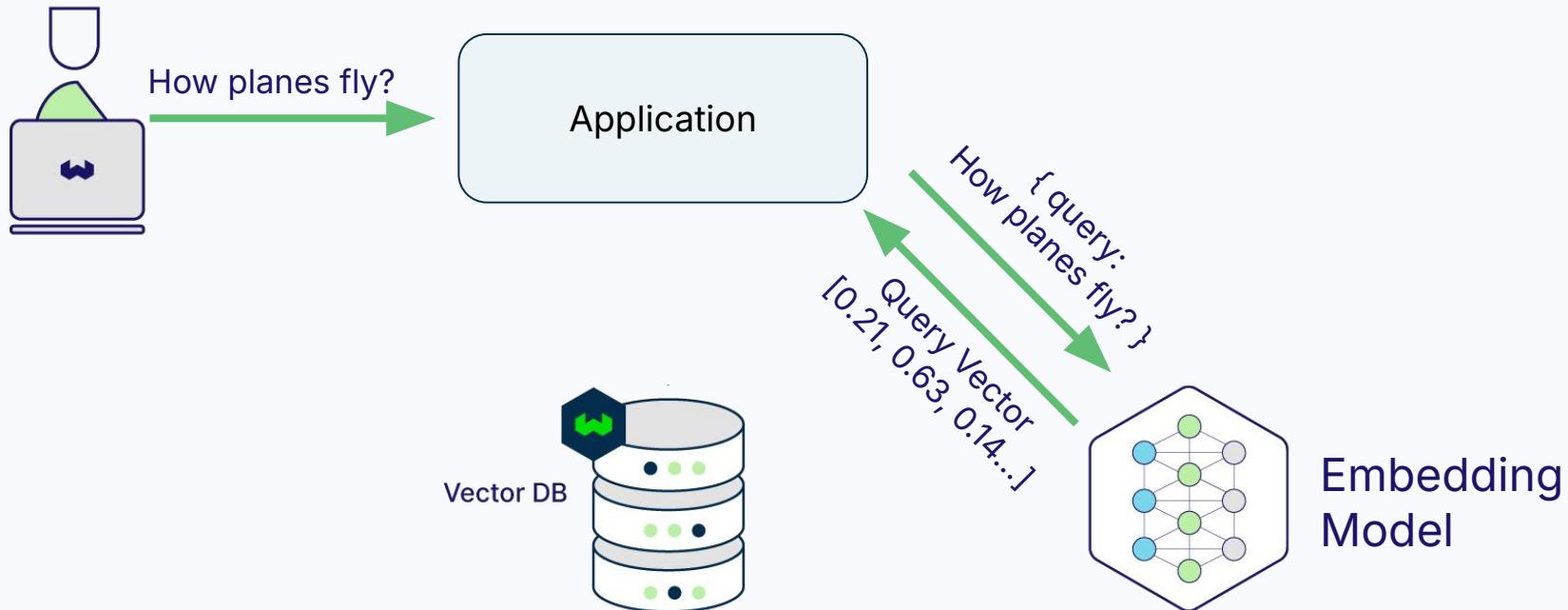
# Bring your own vectors

# Import Data with Vectors (no vectorizer)



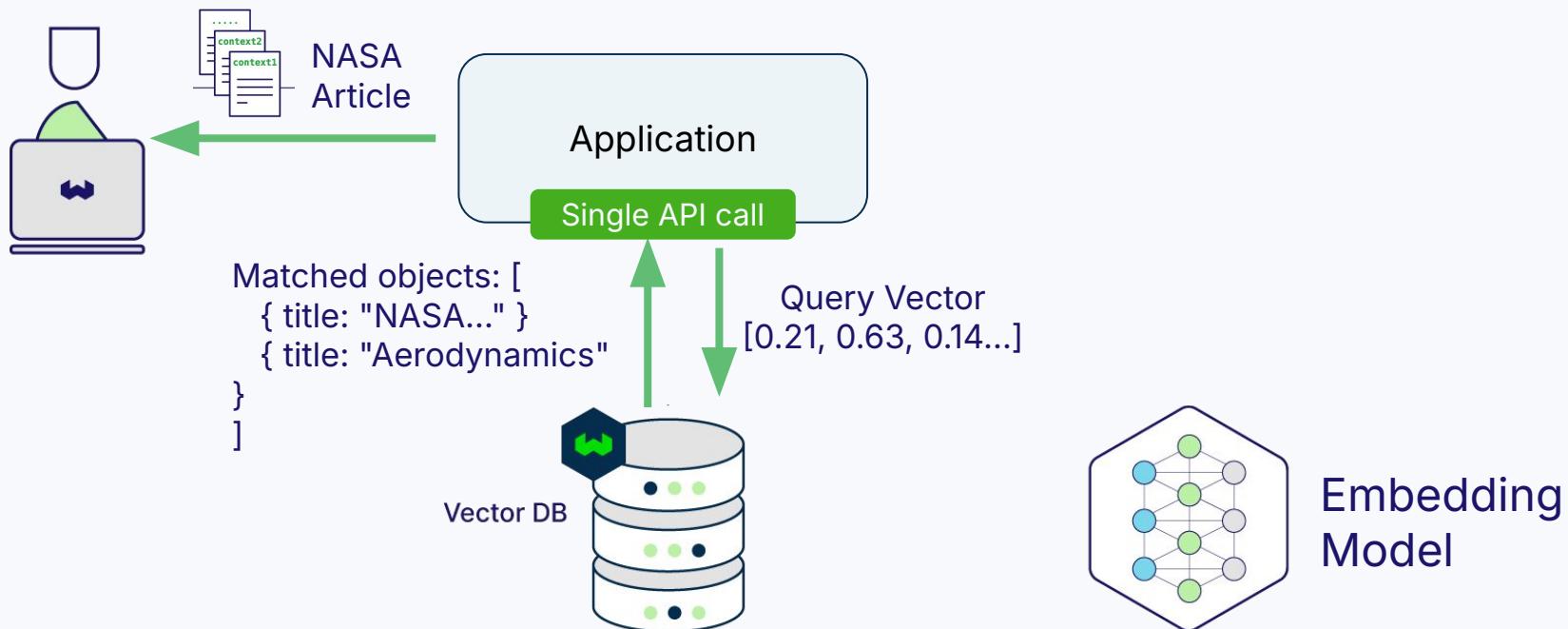


# Query (no vectorizer)



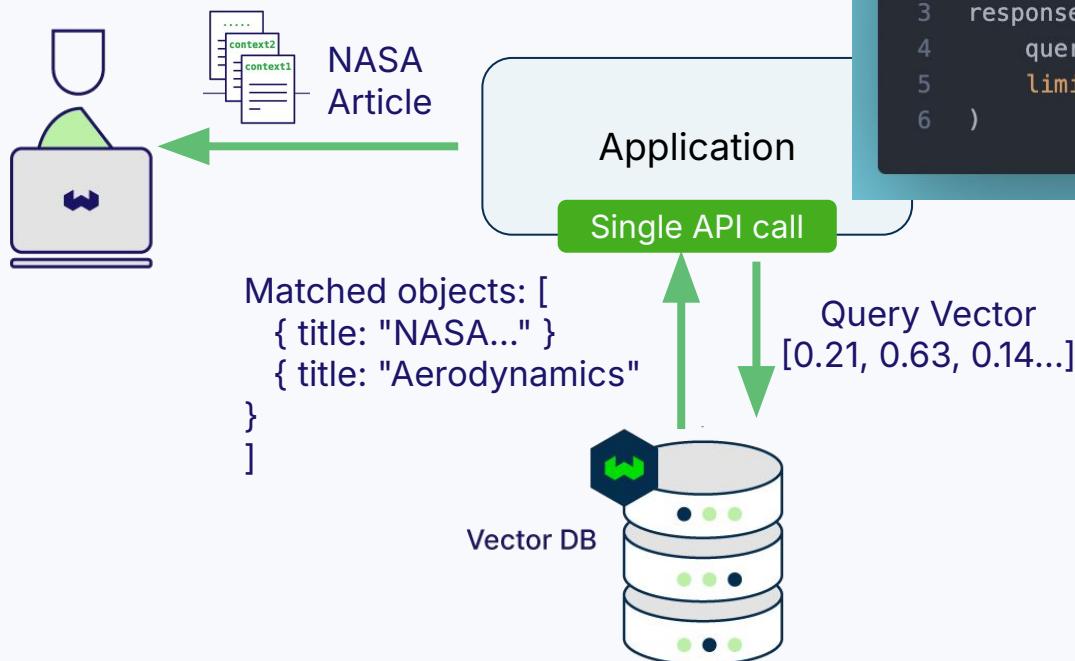


# Query (no vectorizer)



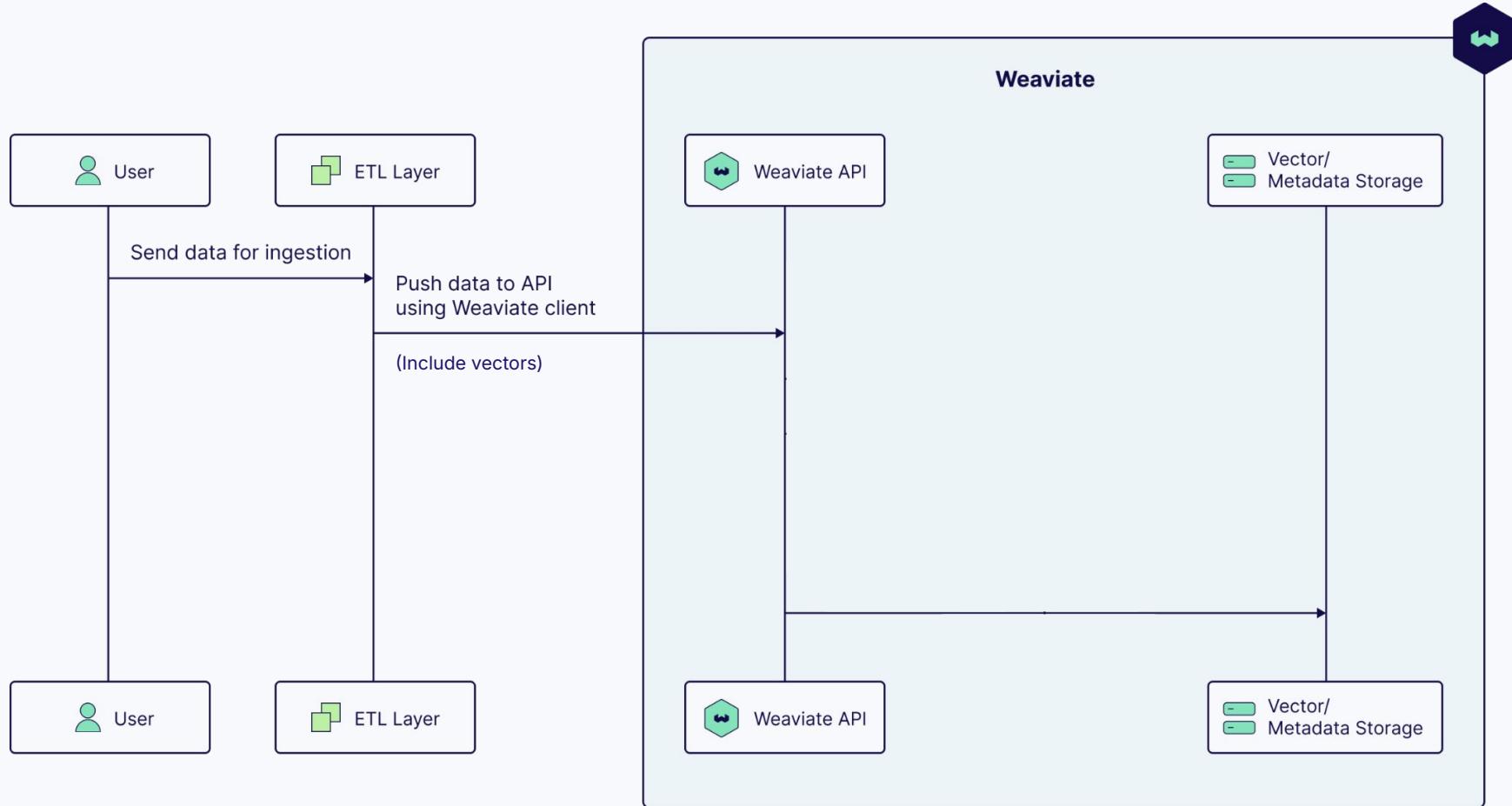


# Query (no vectorizer)



```
1 products = client.collections.get("Product")
2
3 response = products.query.near_vector(
4     query_vector, # Manually obtained query vector
5     limit=2
6 )
```

# Data Ingestion





# GenAI model integration

# Generative AI - basic usage

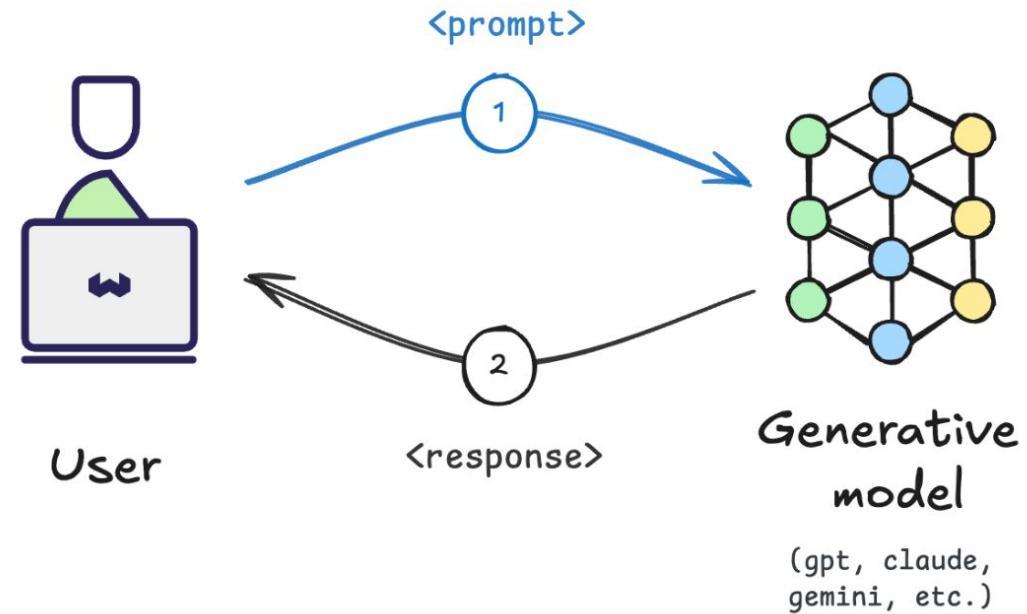
A generative AI model (e.g. an LLM) can generate a response based on a prompt.  
The prompt can be a single "question", or simulate a conversation.

Example prompt:

'What is the capital of France?'

Example response:

'The capital of France is Paris.'





Vector Database

Large Language Model

# One Key Difference





# Vector Database (Stateful)

Large Language Model  
(Stateless)



## Vector Database (Stateful)

- Insert data into the database, and generate vectors.

## Large Language Model (Stateless)

- The knowledge base is trained into the model.



## Vector Database (Stateful)

- Insert data into the database, and generate vectors.
- Data changes on the fly.

## Large Language Model (Stateless)

- The knowledge base is trained into the model.
- The model doesn't learn on the fly.



# Vector Database (Stateful)

- Insert data into the database, and generate vectors.
- Data changes on the fly.
- CRUD operations to update the state.

# Large Language Model (Stateless)

- The knowledge base is trained into the model.
- The model doesn't learn on the fly.
- Requires fine-tuning to update its worldview.



# This is **not** a story of VDBs vs LLMs





# This **is** a story of **VDBs + LLMs**





# This **is** a story of **VDBs + LLMs**

## **Focus on content:**

- Keep the data up to date
- Provide accurate context to LLMs



## **Focus on response:**

- Speak with your brand voice
- Boost understanding of the underlying context



## Retrieval augmented generation

In retrieval augmented generation (RAG), additional information, or "context" is passed to the generative AI model along with a task.

The context is retrieved from a trusted data source, so that the generated text can be based on accurate and up-to-date information.

Example query:

```
wiki.query.hybrid(<search_string>)
```

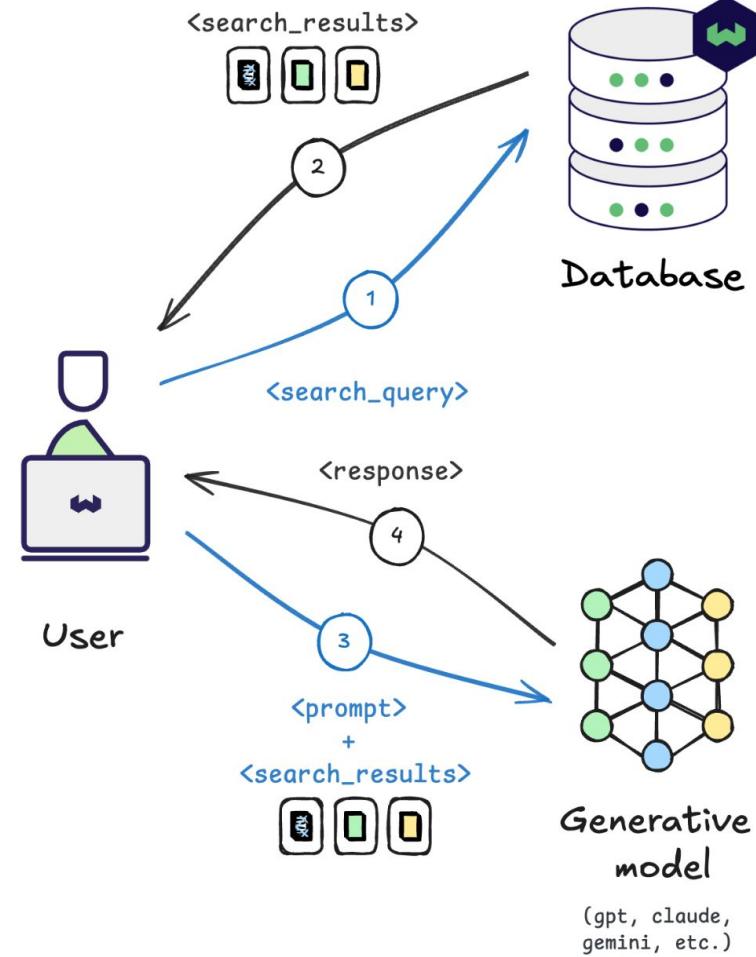
Example prompt:

```
Please answer <question>,
using the included information:
```

```
===== INFORMATION =====
<search_results>
=====
'
```

Example response:

```
'Based on the provided information, ...'
```



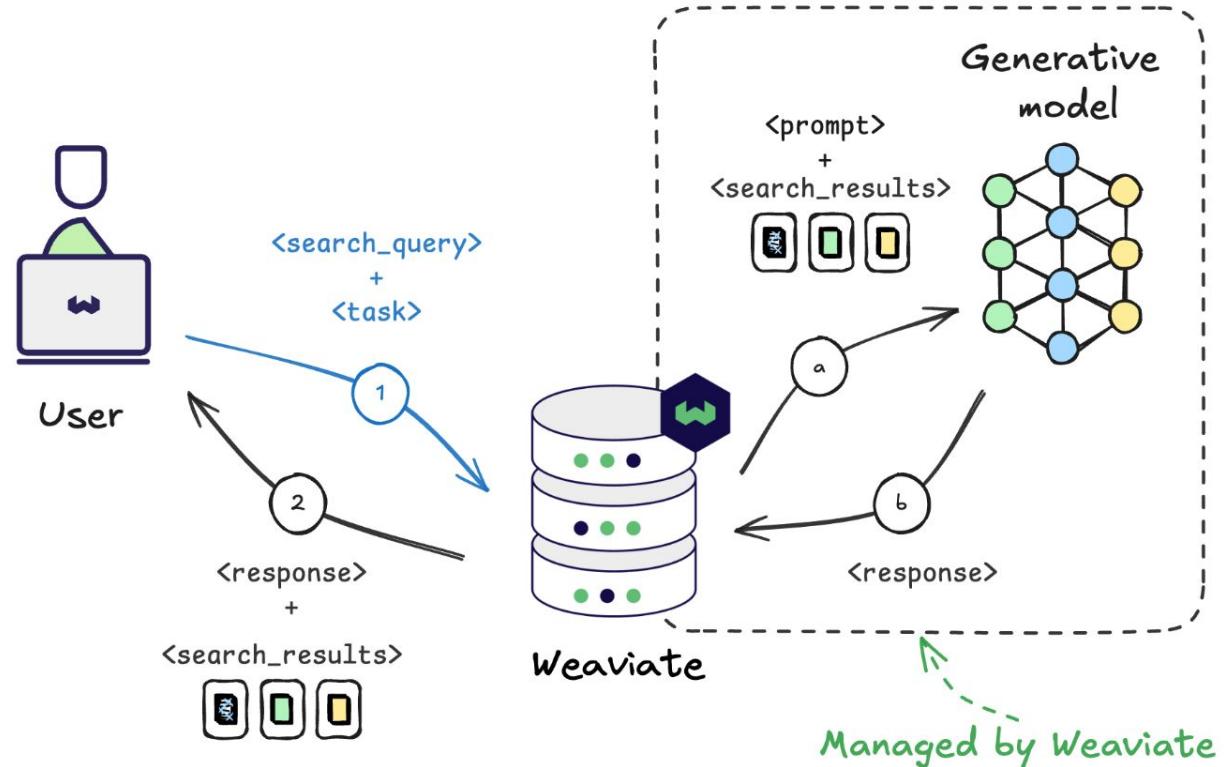
## Retrieval augmented generation with Weaviate

Weaviate simplifies the RAG workflow by combining the retrieval and generation steps into one.

Define and send the search query to Weaviate along with the task, and Weaviate manages the rest.

Example query:  
`wiki.generate.hybrid(  
 <search_string>,  
 grouped_task=<question>  
)`

Example response:  
`'Based on the provided information, ...'`





# RAG – grouped task

```
1 reviews = client.collections.get("WineReview")
2 response = reviews.generate.near_text(
3     query="a sweet German white wine",
4     limit=2,
5     target_vector="description",
6     grouped_task="Summarize these reviews",
7     generative_provider=GenerativeConfig.openai(),
8 )
9
```



# RAG – provider options

```
● ● ●  
1  response = collection.generate.near_text(  
2      query="A holiday film",  
3      limit=2,  
4      grouped_task="Write a tweet promoting these two movies",  
5      # highlight-start  
6      generative_provider=GenerativeConfig.anthropic(  
7          # # These parameters are optional  
8          # base_url="https://api.anthropic.com",  
9          # model="claude-3-opus-20240229",  
10         # max_tokens=512,  
11         # temperature=0.7,  
12         # stop_sequences=["\n\n"],  
13         # top_p=0.9,  
14         # top_k=5,  
15     ),  
16     # highlight-end  
17 )
```



# Multi- modality



# Why PDFs?



# PDFs are amazing... and complex

**Emptying your bin**

Empty your bin as soon as the dirt reaches the Max mark. You may need to clean the filter more often if you use your machine with a full bin.

Press the red wand release button and pull the wand away from the bin.  
Don't pull the trigger when emptying the bin.

**Washing your bin**

Don't put any part of your machine in a dishwasher or use detergents, polishes or air fresheners.

Empty your bin before cleaning it.

Press the red button on the bin runner to release the bin, then slide the bin off the runner.

Rinse your bin under a cold tap. Wipe with a dry lint-free cloth to remove any debris.

Make sure that all parts are completely dry.

Slide the bin onto the bin runner.

Push the bin back until the bin clicks into place. Make sure the bin is securely closed.

6

Patented May 8, 1934  
Des. 92,189  
UNITED STATES PATENT OFFICE

DESIGN FOR A DOUBLE-DECK TRANSPORT AIRPLANE  
Robert J. Minshall and Frank R. Clegg, West Coast Aeroplane Corporation, Seattle, Wash., a corporation of Washington

2023 marked a significant increase in the number of newly funded generative AI companies, with 99 new startups receiving funding, compared to 56 in 2022, and 31 in 2019 (Figure 4.3.6).

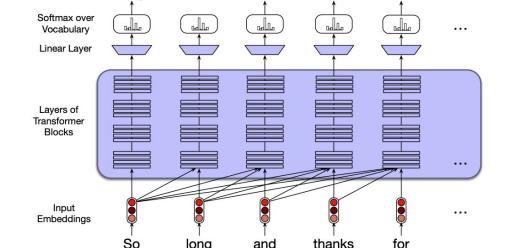
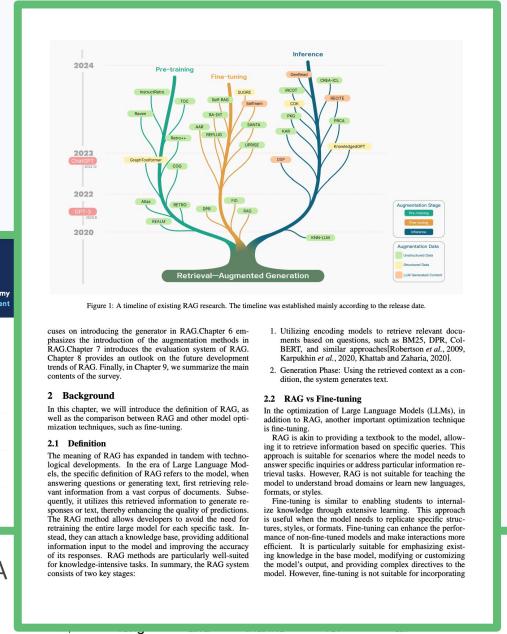
**Number of newly funded generative AI companies in the world, 2019–23**  
Source: Gartner, 2023 (Cited: 2024 AI Index report)

Year	Number of Companies
2019	~30
2020	~30
2021	~50
2022	~60
2023	99

**AI private investment events by size, 2022 vs. 2023**  
Source: CB Insights, 2024 (Cited: 2024 AI Index report)

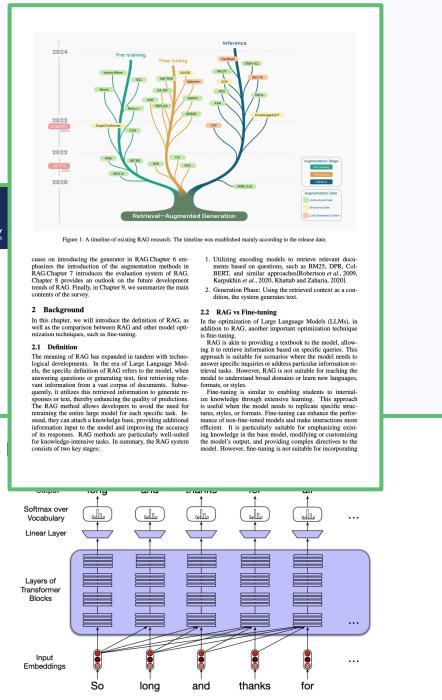
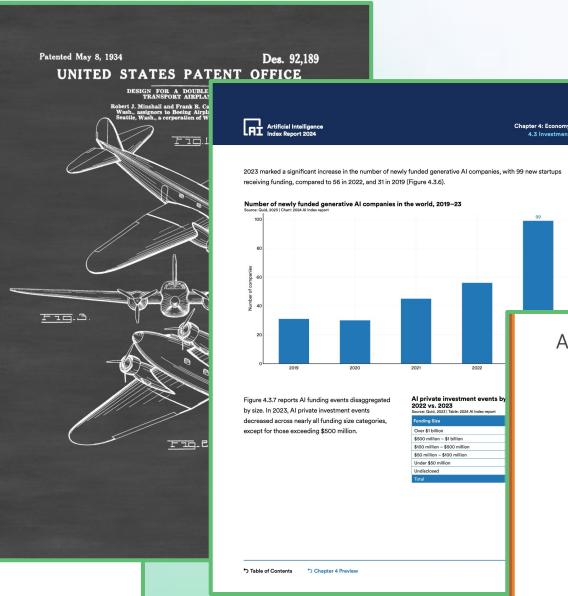
Funding Size	2022 Events	2023 Events
Over \$1 billion	~10	~10
\$500 million – \$1 billion	~10	~10
\$100 million – \$500 million	~10	~10
\$50 million – \$100 million	~10	~10
Under \$50 million	~10	~10
Undisclosed	~10	~10
Total	~50	~90

\* Table of Contents      \* Chapter 4 Preview





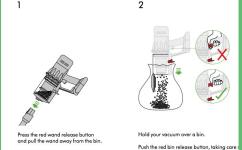
# How to go from here...





# How to go from here...

# to here?



**Patented May 8, 1934** **Des. 92,189**  
**UNITED STATES PATENT OFFICE**

DESIGN FOR A DOUBLE SEAT AIRPLANE

Robert J. Mishell and Frank R. Clegg, Bothell, Wash., a representation of W.

Figure 4.3 Artificial Intelligence Report 2024

Chapter 4: Economy 4.3 Investment

2023 marked a significant increase in the number of newly funded generative AI companies, with 99 new startups receiving funding, compared to 56 in 2022, and 31 in 2019 (Figure 4.3.6).

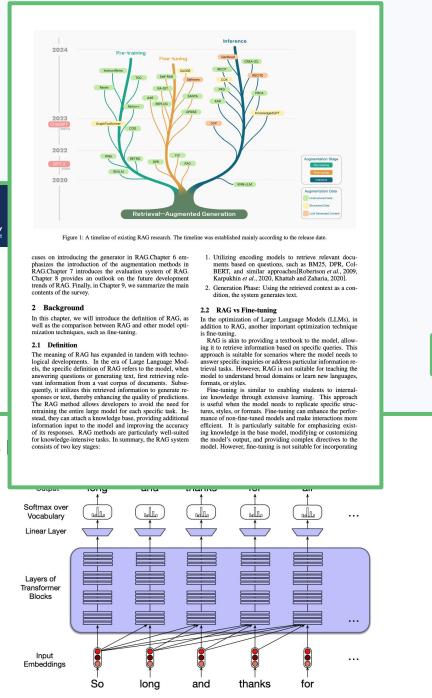
Number of newly funded generative AI companies in the world, 2019–23

Year	Number of companies
2019	~31
2020	~56
2021	~65
2022	~99
2023	~99

AI private investment events by 2022 vs. 2023

Category	2022	2023
Funding size	Under \$50 million	~\$50 million - \$1 billion
Others	~\$1 billion - \$500 million	~\$500 million - \$1 billion
Total	~\$500 million - \$1 billion	~\$500 million - \$1 billion

[Table of Contents](#) [Chapter 4 Preview](#)





# Solutions & challenges



# With existing tooling

- Text-only examples:
  - Simple layout
  - Easy to parse & convert

```
>>> p = re.compile(r'\W+')
>>> p.split('This is a test, short and sweet, of split().')
['This', 'is', 'a', 'test', 'short', 'and', 'sweet', 'of', 'split', '.']
>>> p.split('This is a test, short and sweet, of split().', 3)
['This', 'is', 'a', 'test, short and sweet, of split().']
```

Sometimes you're not only interested in what the text between delimiters was. If capturing parentheses are used in the RE, then their values are a following calls:

```
>>> p = re.compile(r'\W+')
>>> p2 = re.compile(r'(\W+)')
>>> p.split('This... is a test.')
['This', 'is', 'a', 'test', '.']
>>> p2.split('This... is a test.')
['This', '...', 'is', '.', 'a', ' ', 'test', '.', '']
```

The module-level function `re.split()` adds the RE to be used as the

```
>>> re.split(r'[\W]+', 'Words, words, words.')
['Words', 'words', 'words', '']
>>> re.split(r'([\W]+)', 'Words, words, words.')
['Words', ' ', 'words', ' ', 'words', ' ', '']
>>> re.split(r'[\W]+', 'Words, words, words.', 1)
['Words', 'words', 'words']
```

## 5.2 Search and Replace

Another common task is to find all the matches for a pattern, and `re.match` method takes a replacement value, which can be either a string or a function.

`.sub(replacement, string[, count=0])`

Returns the string obtained by replacing the leftmost non-overlapping replacement `replacement`. If the pattern isn't found, `string` is returned.

The optional argument `count` is the maximum number of pattern non-negative integer. The default value of 0 means to replace all

Here's a simple example of using the `sub()` method. It replaces colour

```
>>> p = re.compile('(blue|white|red)')
>>> p.sub('colour', 'blue socks and red shoes')
'colour socks and colour shoes'
>>> p.sub('colour', 'blue socks and red shoes', coun
'colour socks and red shoes'
```

The `subn()` method does the same work, but returns a 2-tuple containing replacements that were performed:

```
>>> p = re.compile('(blue|white|red)')
>>> p.subn('colour', 'blue socks and red shoes')
('colour socks and colour shoes', 2)
>>> p.subn('colour', 'no colours at all')
('no colours at all', 0)
```

Empty matches are replaced only when they're not adjacent to a previous

Published as a conference paper at ICLR 2025

## ColPali: EFFICIENT DOCUMENT RETRIEVAL WITH VISION LANGUAGE MODELS

Manuel Fayolle<sup>1,3</sup> Hugues Sibille<sup>1,\*4</sup> Tony Wu<sup>1</sup> Bilel Omrani<sup>1</sup>

Gautier Viard<sup>1</sup> Céline Hudelot<sup>1</sup> Pierre Colombo<sup>2,3</sup>

<sup>1</sup>Illiun Technology <sup>2</sup>EqualLai <sup>3</sup>CentraleSupélec, Paris-Saclay <sup>4</sup>ETH Zurich

manuel.fayolle@centralesupelec.fr

## ABSTRACT

Documents are visually rich structures that convey information through text, but also figures, page layouts, tables, or even fonts. Since modern retrieval systems mainly rely on the textual information they extract from document pages to index documents—often through lengthy and brittle processes—they struggle to exploit key visual cues efficiently. This limits their capabilities in many practical document retrieval applications such as Retrieval Augmented Generation (RAG). To benchmark current systems on visually rich document retrieval, we introduce the Visual Document Retrieval Benchmark *ViDoRe*, composed of various page-level retrieval tasks spanning multiple domains, languages, and practical settings. The inherent complexity and performance shortcomings of current systems motivate a new approach to document retrieval based on visual embeddings and document pages. We release *ColPali*, a Vision Language Model trained to produce high-quality multi-vector embeddings from images of document pages. Combined with a late interaction matching mechanism, *ColPali* largely outperforms modern document retrieval pipelines while being drastically simpler, faster and end-to-end trainable. We release models, data, code and benchmarks under open licenses at <https://hf.co/vidore>.

## 1 INTRODUCTION

Document Retrieval consists of matching a user query to relevant documents in a given corpus. It is central to many widespread industrial applications, either as a standalone ranking system (search engines) or as part of more complex information extraction or Retrieval Augmented Generation (RAG) pipelines.

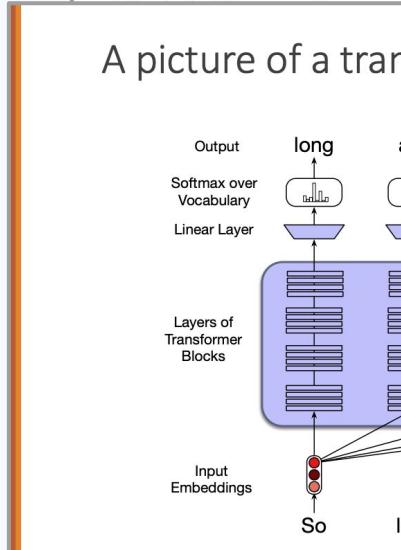
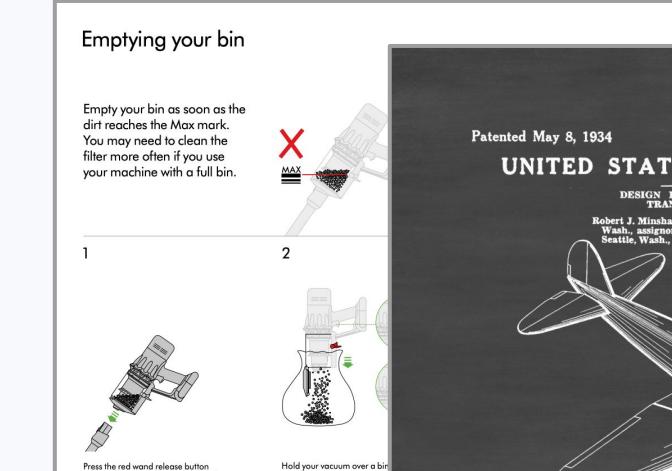
Over recent years, pretrained language models have enabled large improvements in text embedding models. In practical industrial settings, however, the primary performance bottleneck for efficient document retrieval stems not from embedding model performance but from the prior data ingestion pipeline. Indexing a standard PDF document involves several steps. First, PDF parsers or Optical Character Recognition (OCR) systems are used to extract words from the pages. Document layout detection models can then be run to segment paragraphs, titles, and other page objects such as tables, figures, and headings. A matching strategy is then defined to group passages with some semantic coherence and compute retrieval scores. Finally, an embedding function needs to describe visually rich elements in a natural language form, most suitable for embedding models. In our experiments (Table 2), we typically find that optimizing the ingestion pipeline yields much better performance on visually rich document retrieval than optimizing the text embedding model.

**Contribution 1: ViDoRe.** In this work, we argue that document retrieval systems should not be evaluated solely on the capabilities of their embedding models (Bajaj et al., 2016; Thakur et al., 2021; Mousighi et al., 2022), but should also consider the context visual elements of the documents to be retrieved. To this end, we create and openly release *ViDoRe*, a comprehensive benchmark to evaluate systems on page-level document retrieval with a wide coverage of domains, visual elements, and languages. *ViDoRe* addresses practical document retrieval scenarios, where

<sup>\*</sup>Equal Contribution

# With existing tooling

- Mixed media, different layouts
  - More challenging
    - How to decide boundaries?
    - What direction does information flow?
    - Do you know what the user will be looking for?





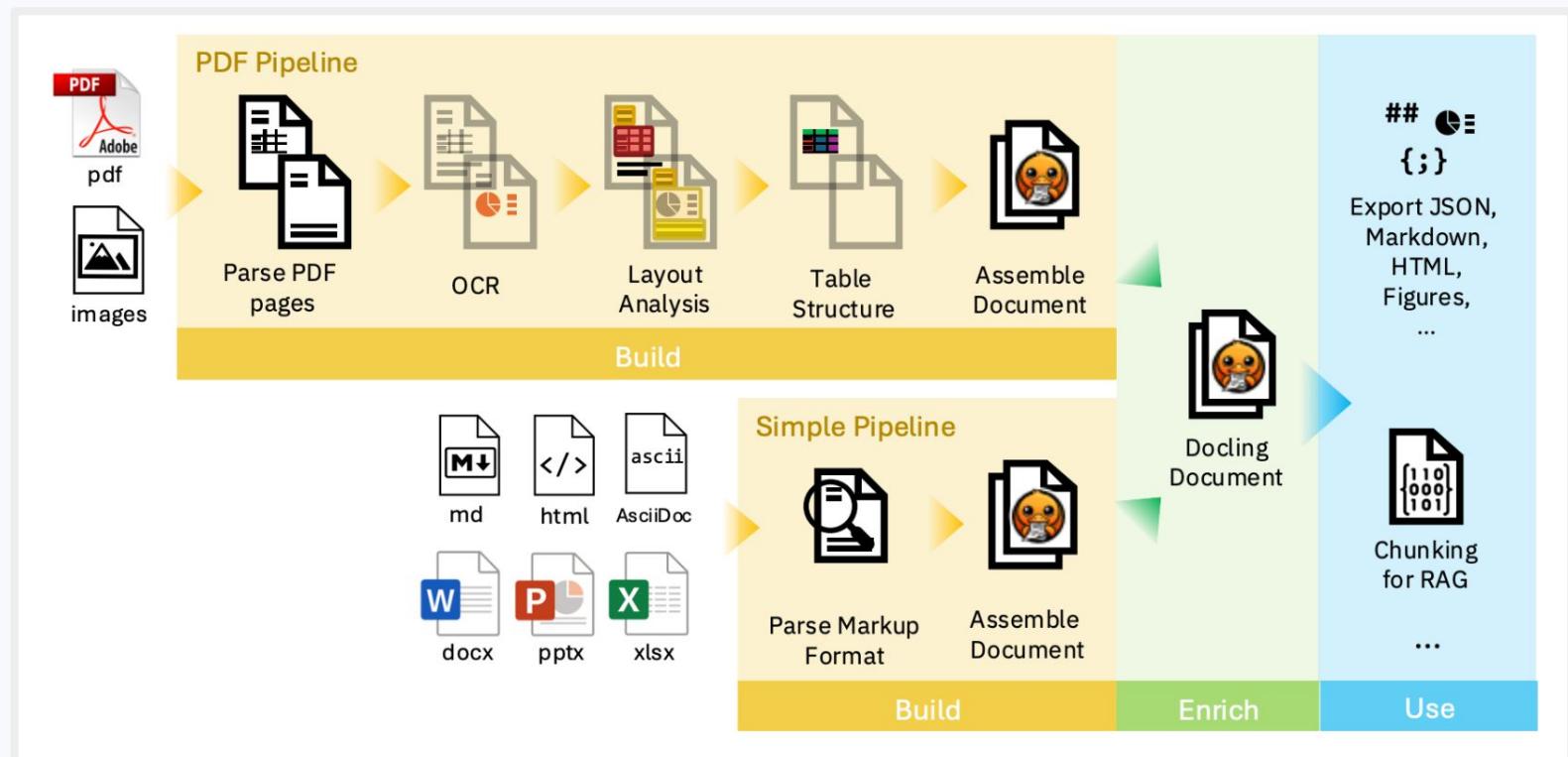
# With existing tooling - example workflow

- Document parsing
- Image & text extraction
- Layout analysis
- Chunking
- Rebuild content (tables / images)
- Text searches (keyword / semantic)





# Example - Docling





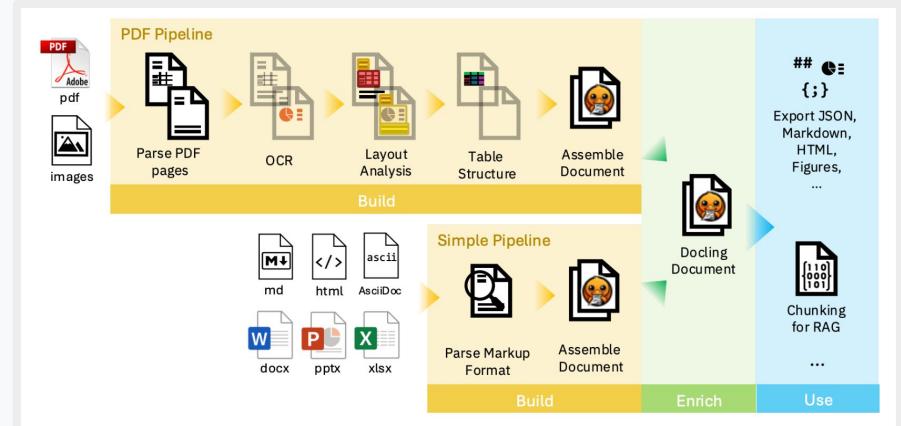
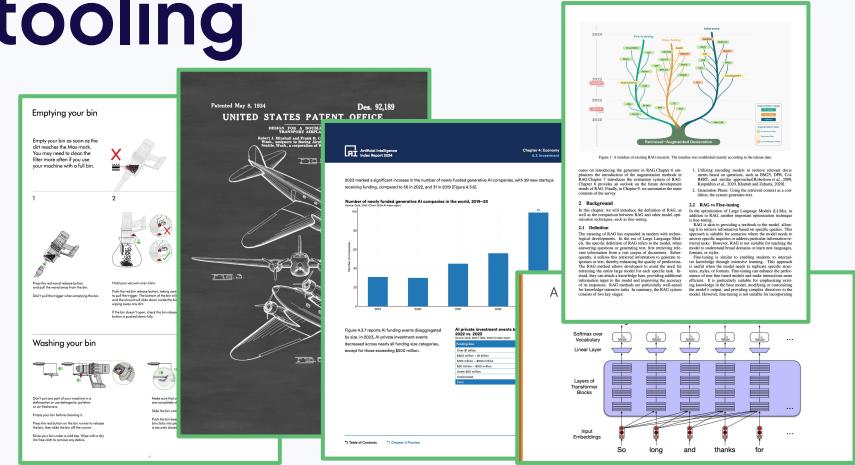
# Demo – Docling

A quick demo 



# Challenges with existing tooling

- Complicated pipeline
  - Image/text borders
  - Relationships between elements
  - Context management
  - Chunking
- Judgement calls required
- One-size-fits-all solution impossible

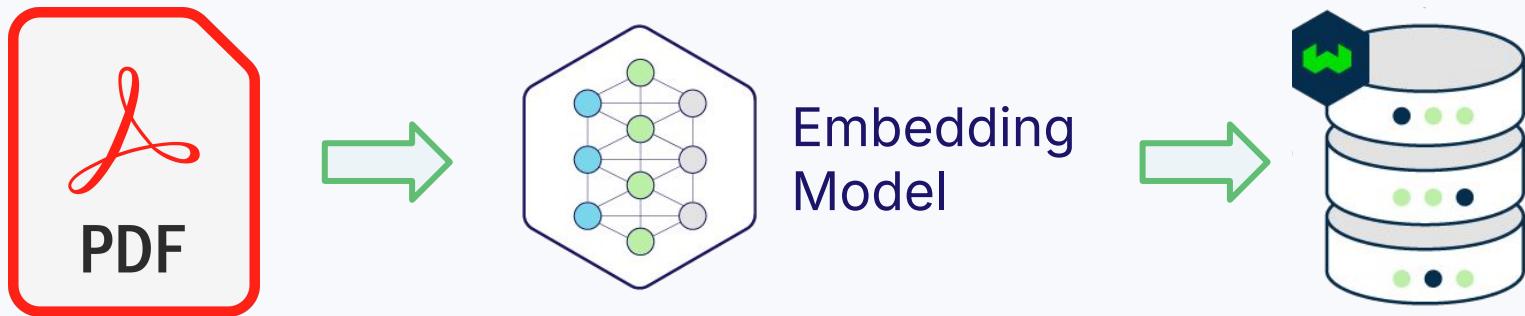




# Another approach

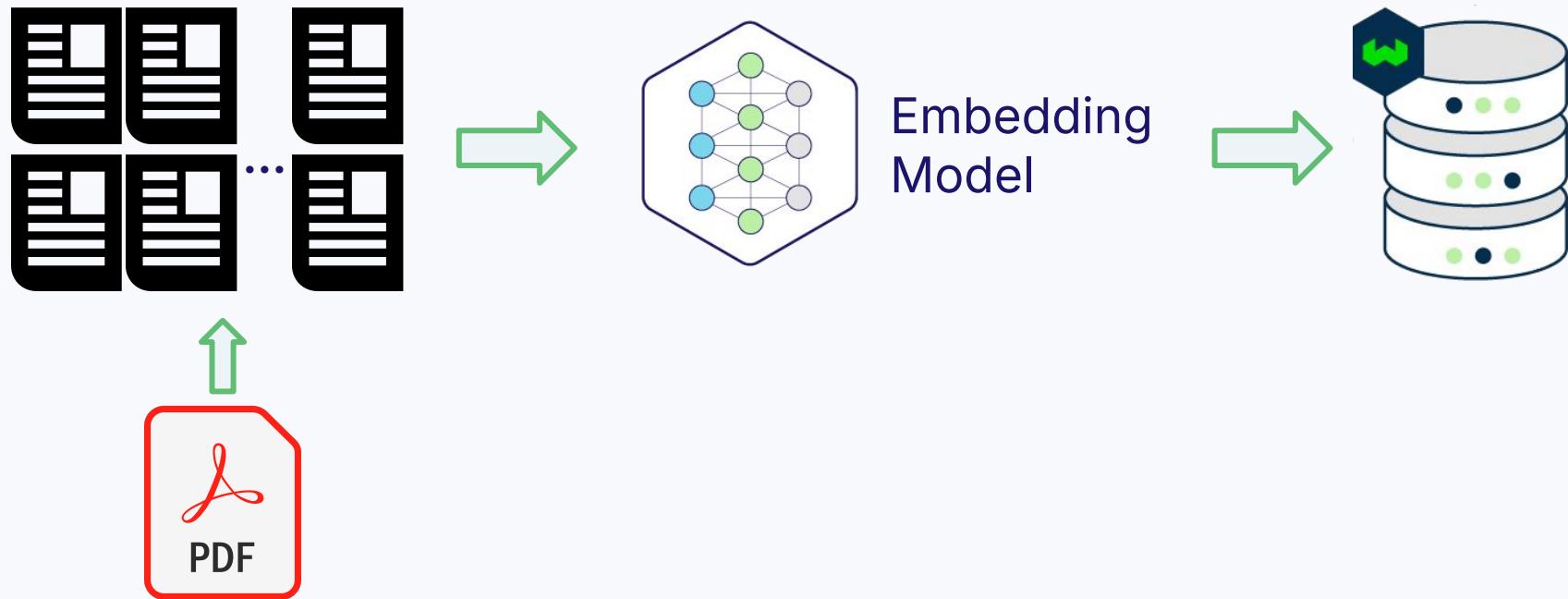


# Just embed the whole page!





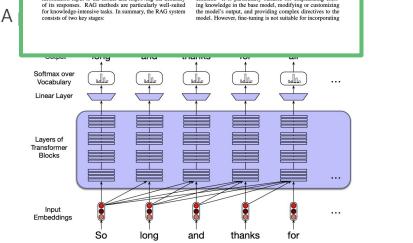
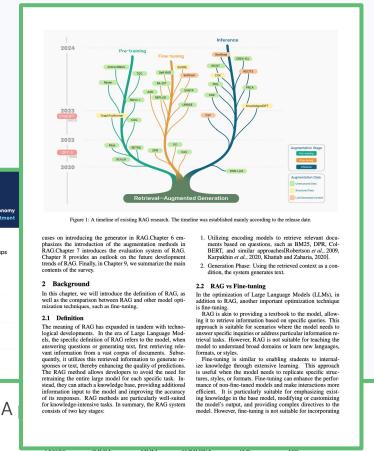
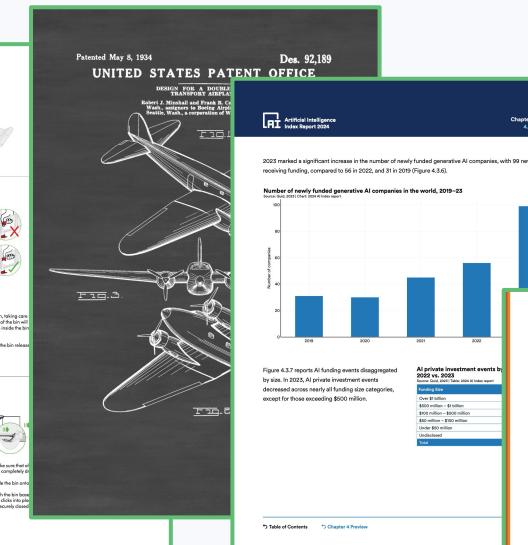
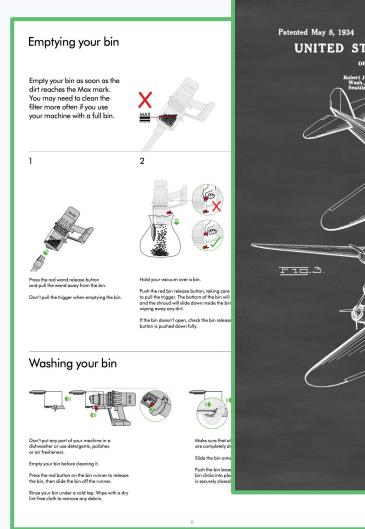
# Just embed the whole page!





# Why is this appealing?

- Respects page layouts
- The data is already chunked (by page)
- Mimics how we read PDFs
- Is robust





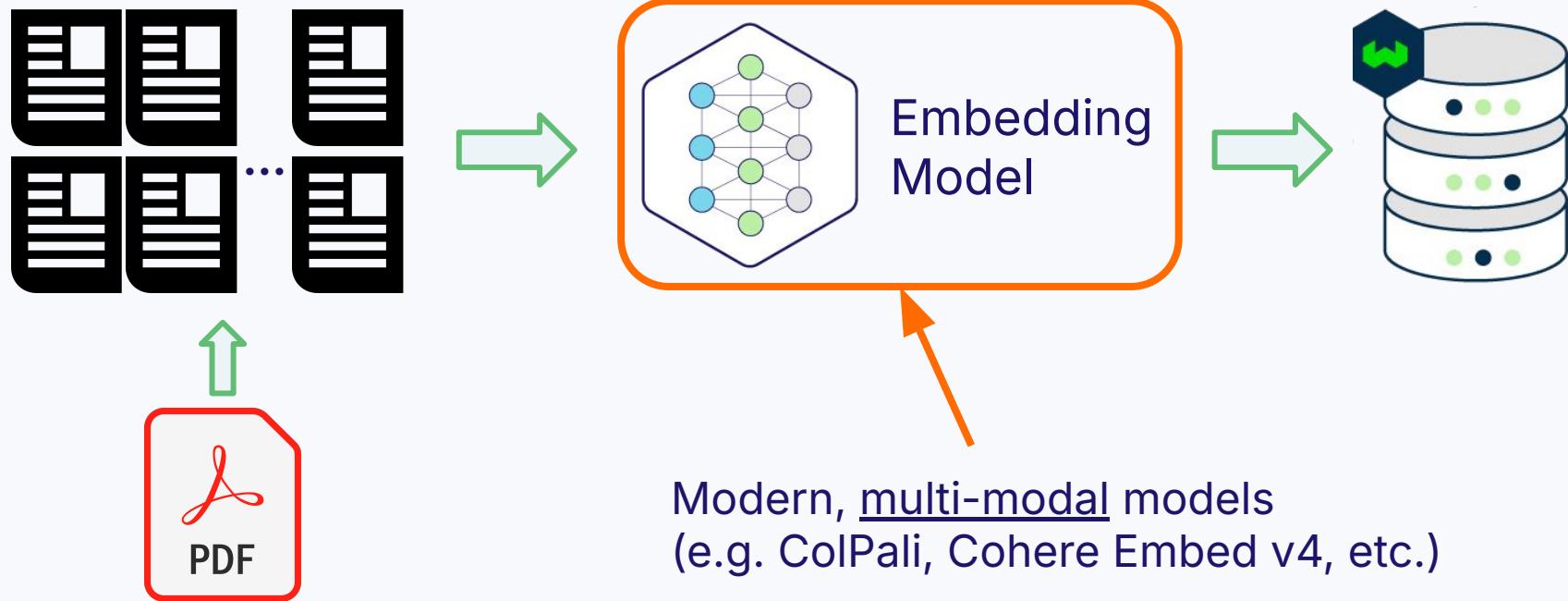
# A Quick Demo

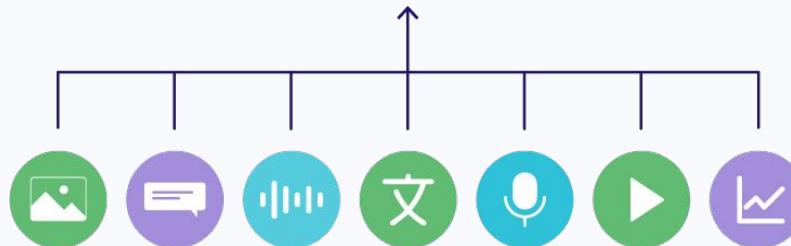
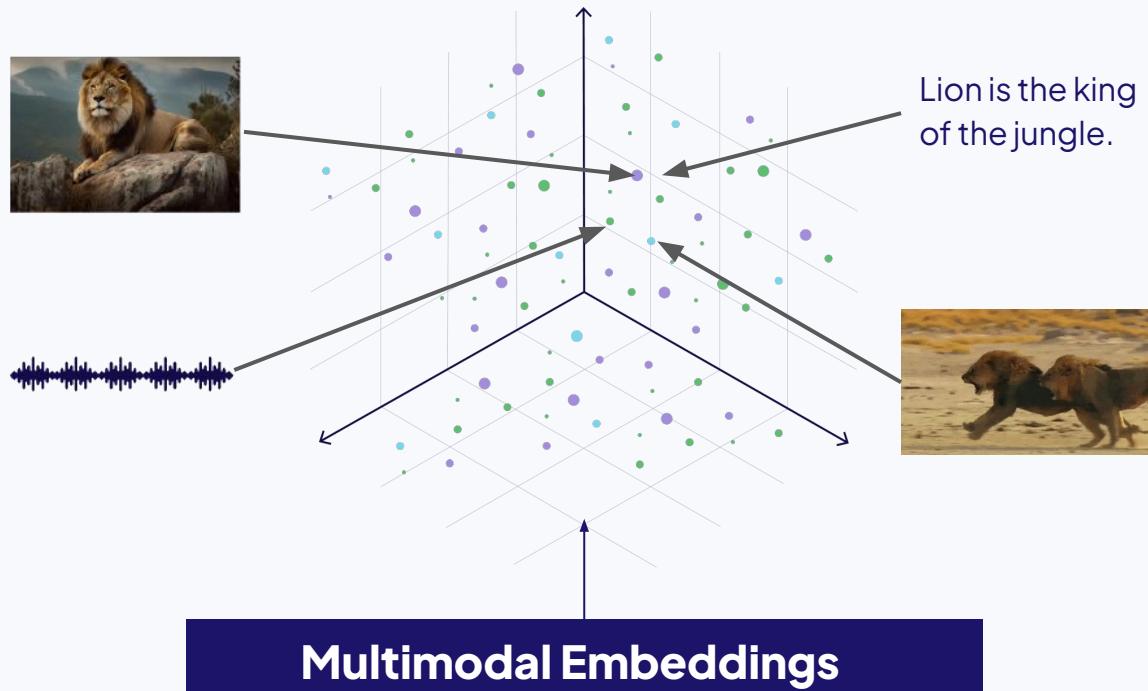


PDF search



# How does this work?







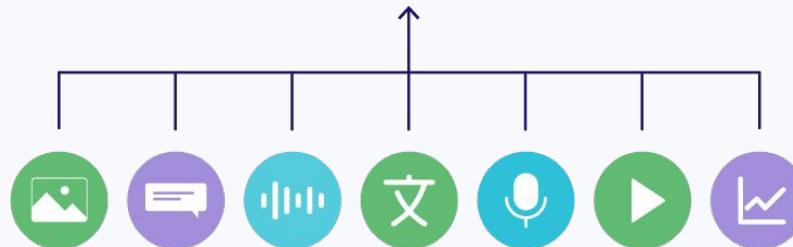
Similar = Close

Dissimilar = Far



Lion is the king  
of the jungle.

## Multimodal Embeddings





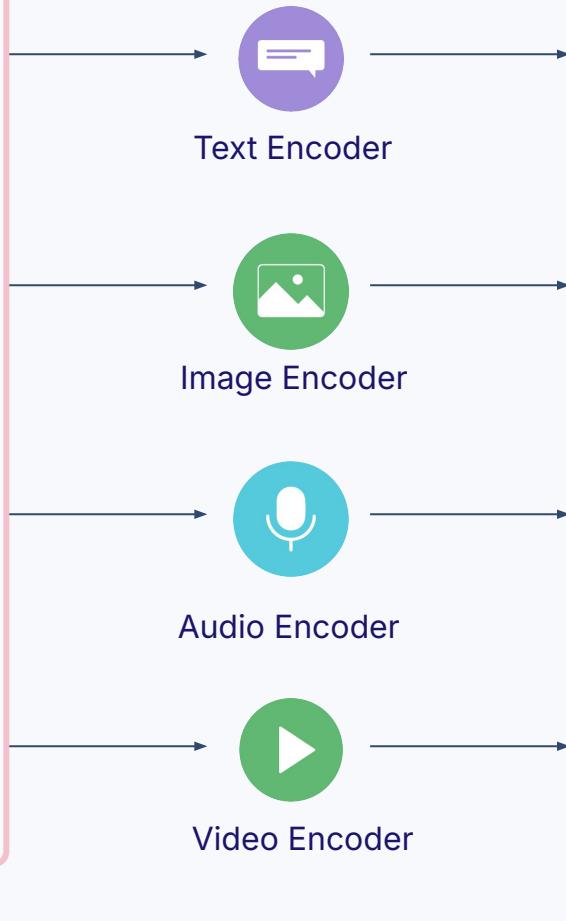
## Start with specialist models

The lion is the king  
of the savannas.





The lion is the king  
of the savannas.



$[-0.10, 0.02, 0.14, \dots, 0.48]$

$[0.61, 0.30, 0.04, \dots, 0.52]$

$[0.06, 0.03, 0.08, \dots, 0.25]$

$[0.02, 0.05, 0.02, \dots, -0.59]$

similar concepts = similar vectors



# Text to Any Search

QUERY

RETRIEVED

Lions roam the savannas



King of the jungle.





# Any to Any Search

## QUERY

Lions roam the savannas.



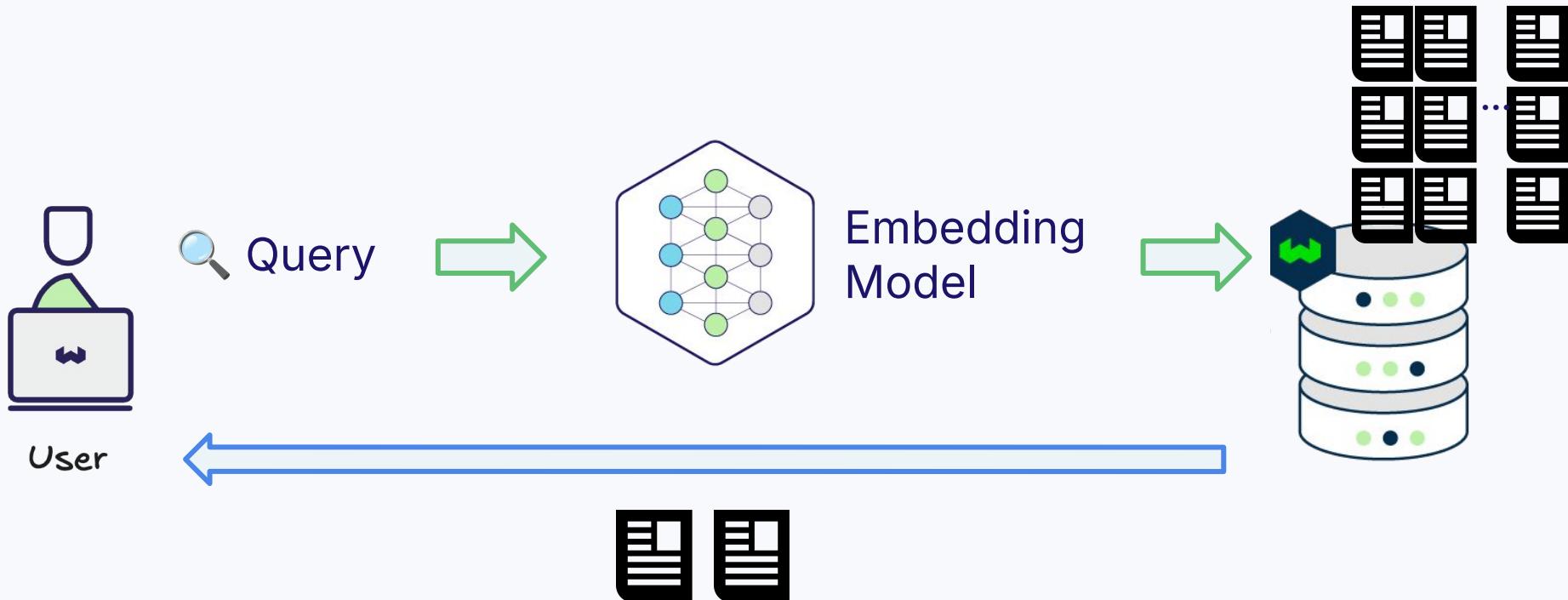
## RETRIEVED

King of the jungle.





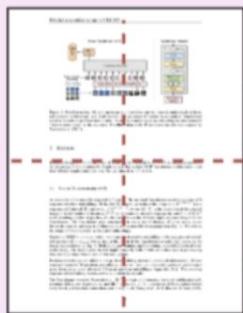
# Multimodal retrieval



# CoIPali model

## CoIPali (ours)

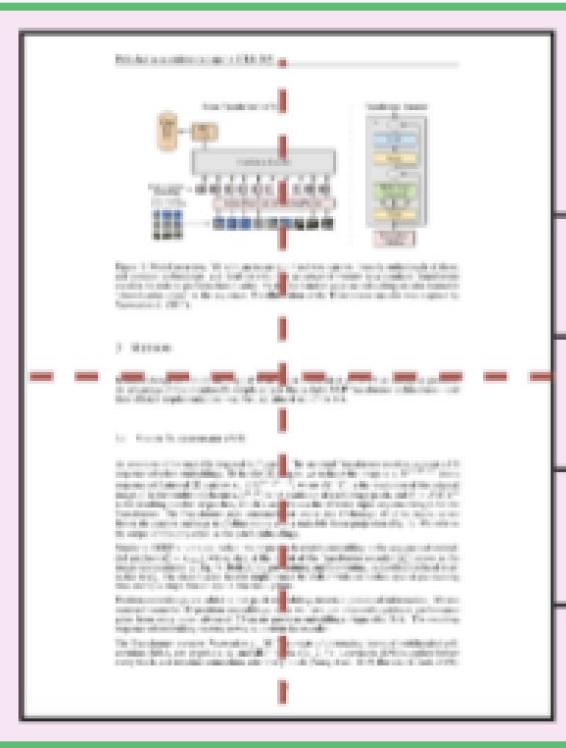
offline



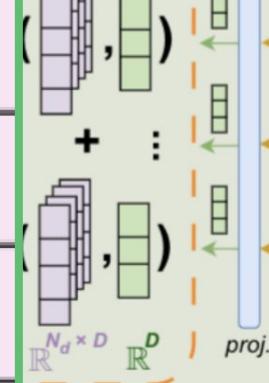
0.81 N

ViT  
encoder

0.39s / page



Visual score



30ms / query

online

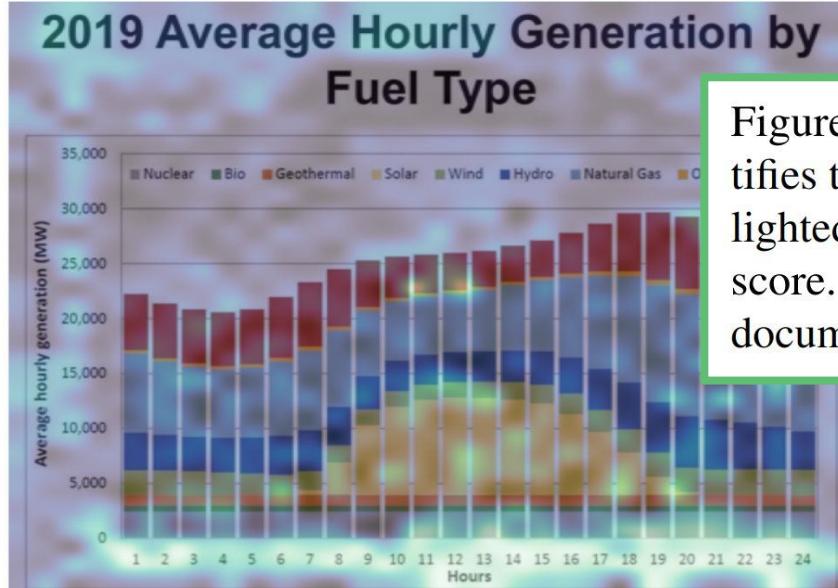
What

are

ViTs?

LLM

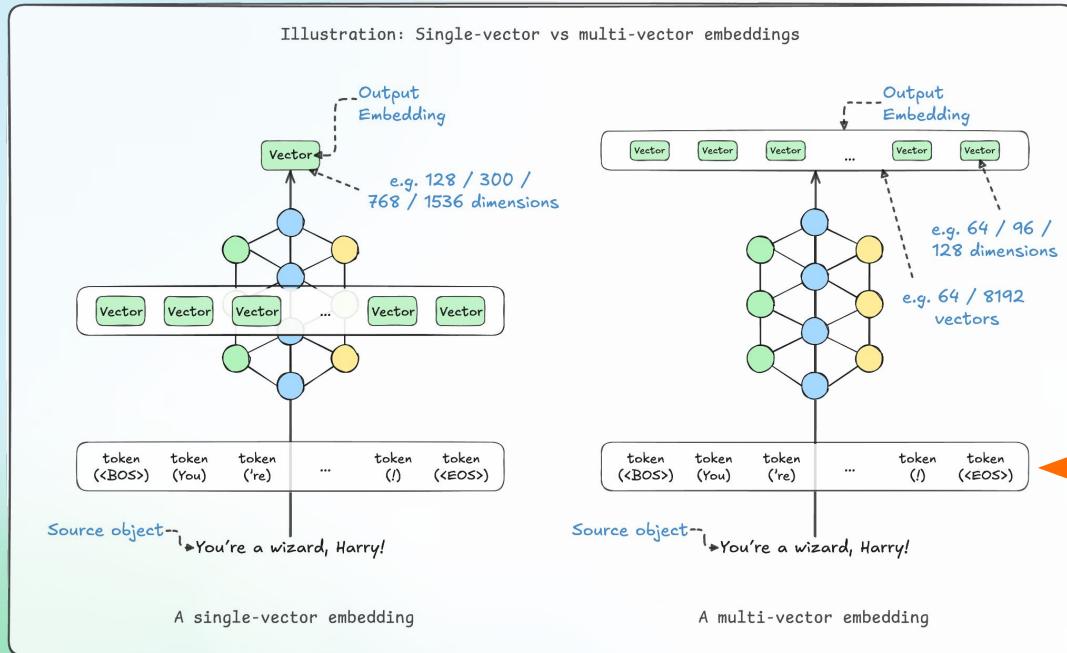
# ColPali model: patches



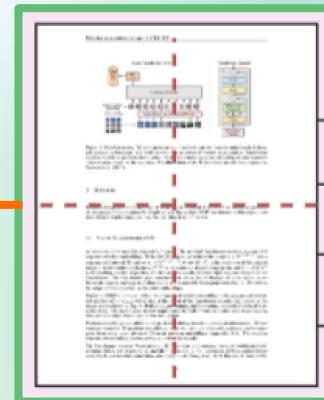
**Query:** "Which hour of the day had the highest overall electricity generation in 2019?"

Figure 1: For each term in a user query, **ColPali** identifies the most relevant document image patches (highlighted zones) and computes a query-to-page matching score. We can then swiftly retrieve the most relevant documents from a large pre-indexed corpus.

# Multi-vector embeddings



Each patch →  
token

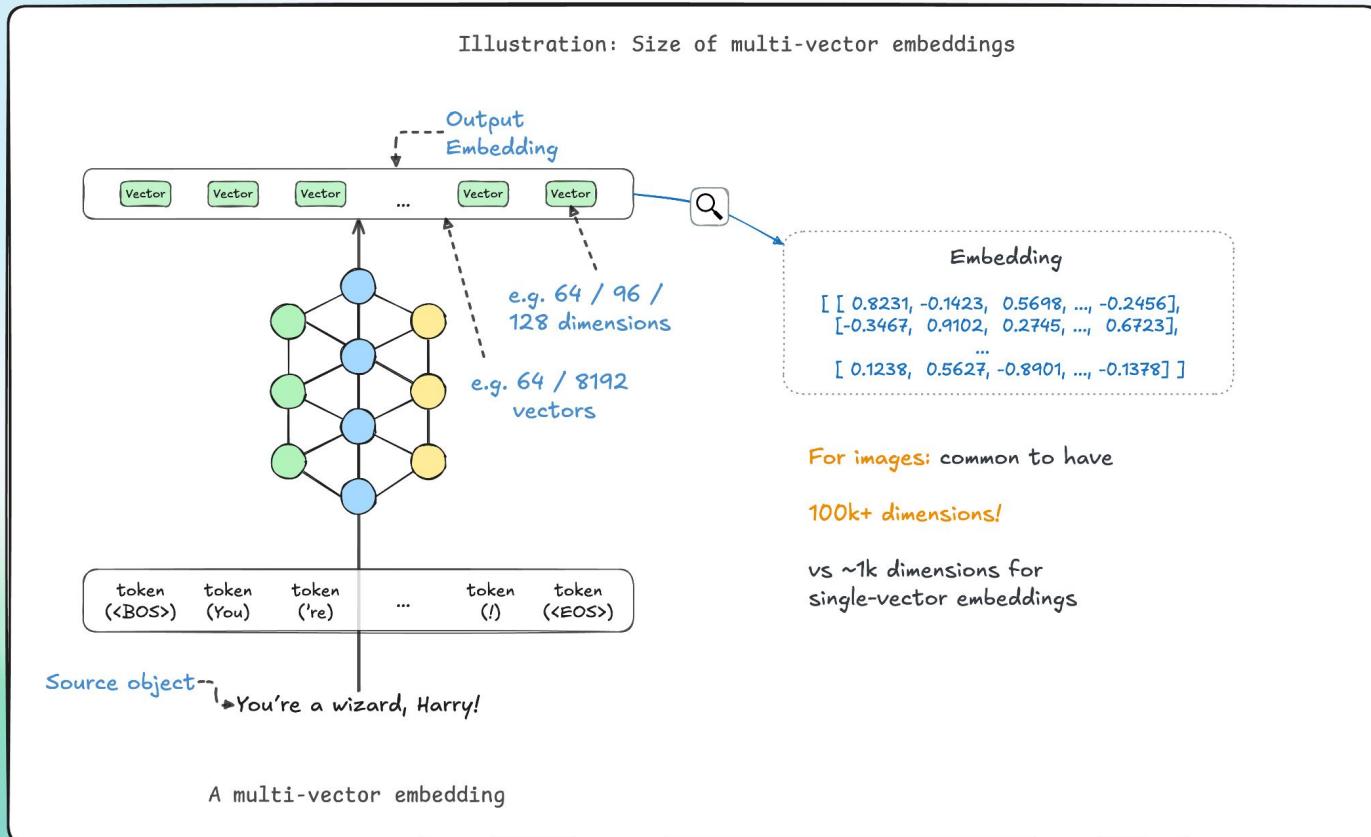




# Scaling multi-vector models

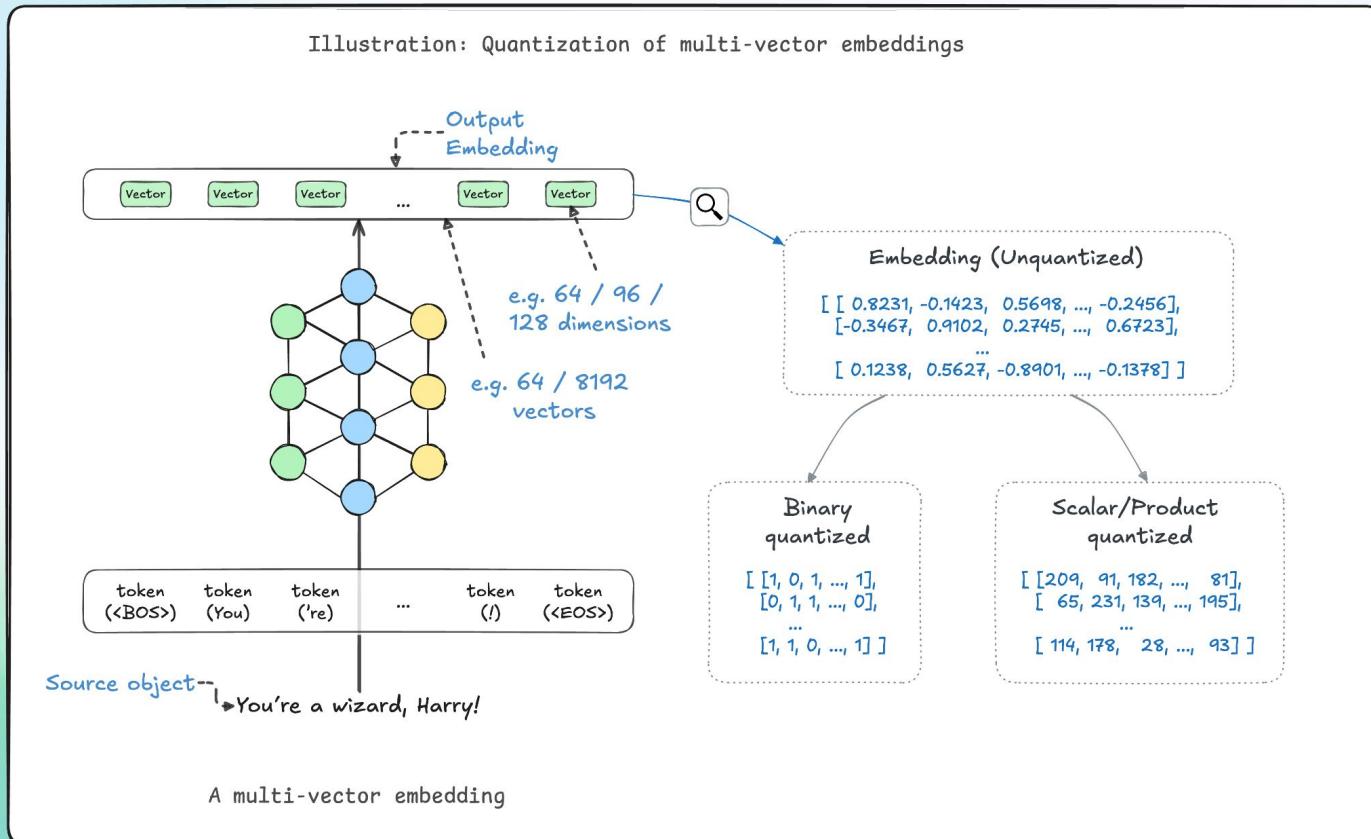


# Multi-vector embeddings



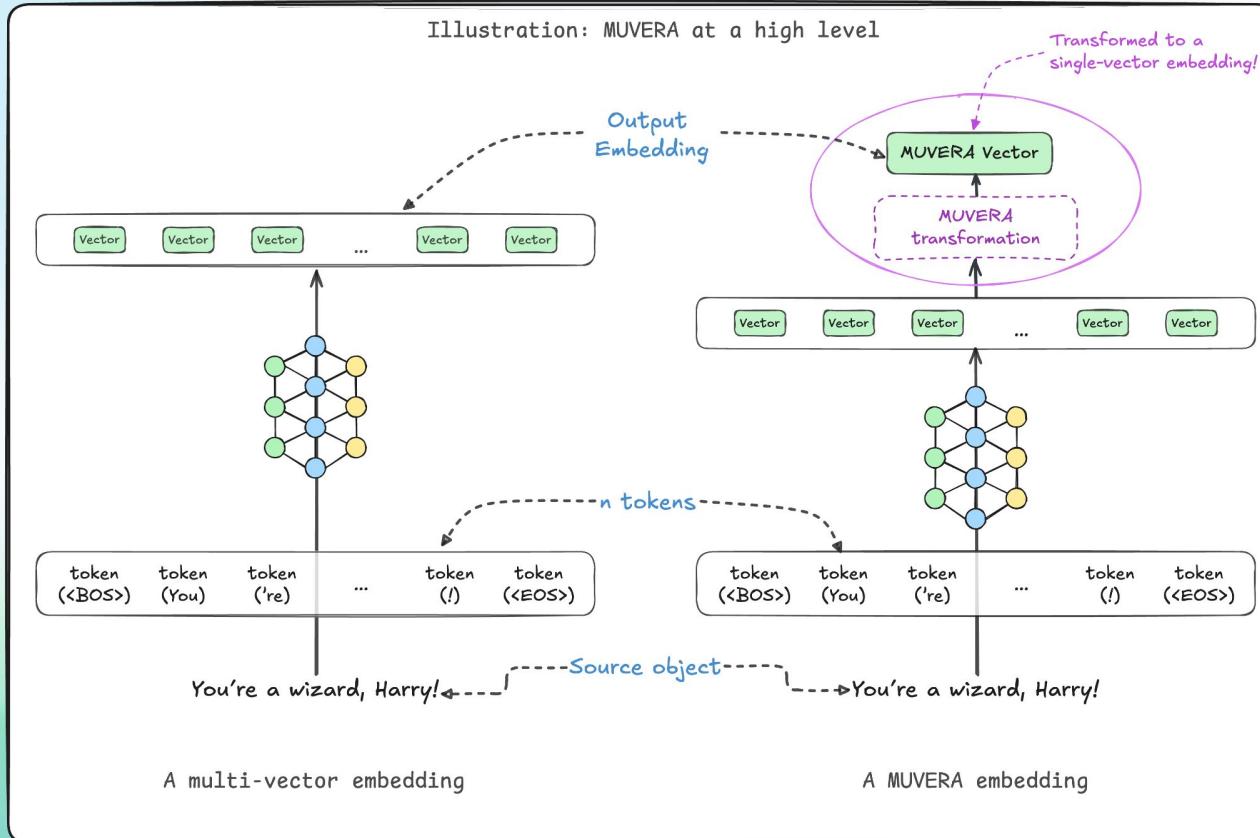


# Multi-vector embeddings: Quantization



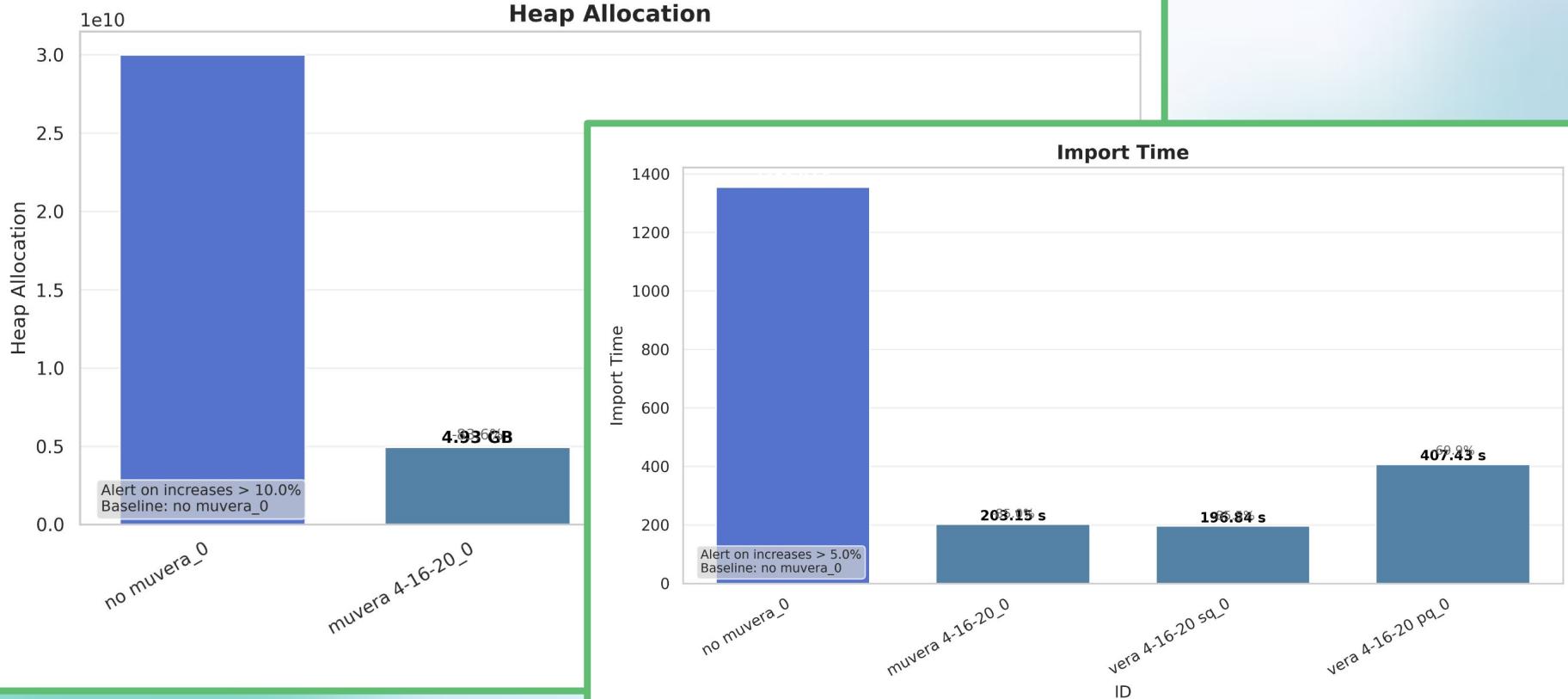


# Multi-vector embeddings: MUVERA

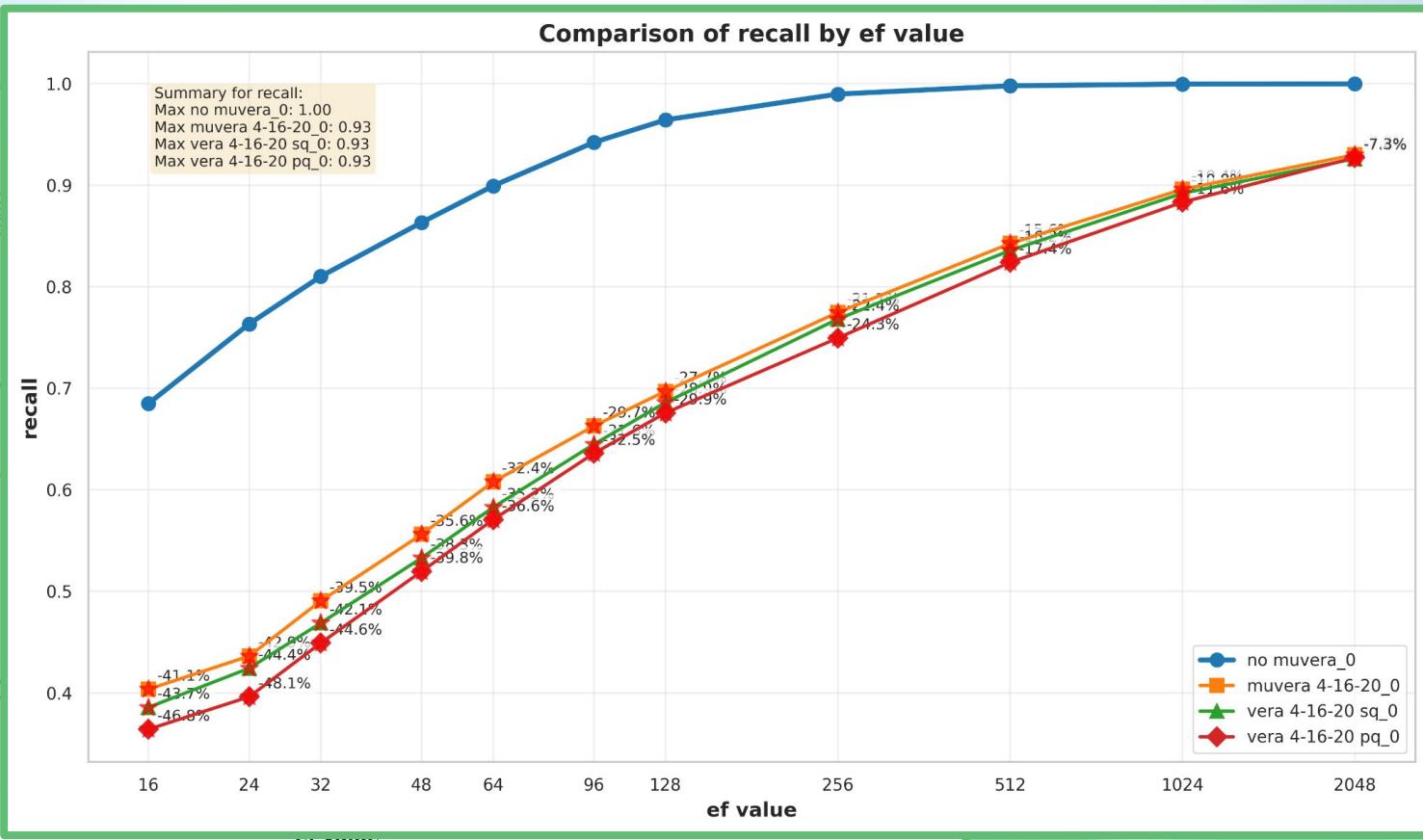
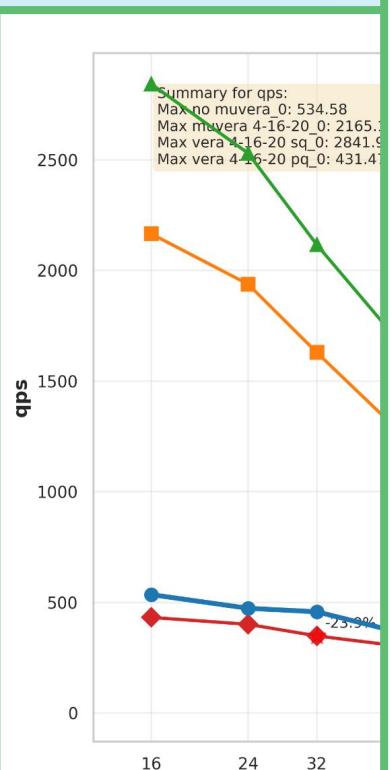




# Multi-vector embeddings: Resources



# Multi-vector embeddings: Quality





# Expanding to RAG

Another quick demo 

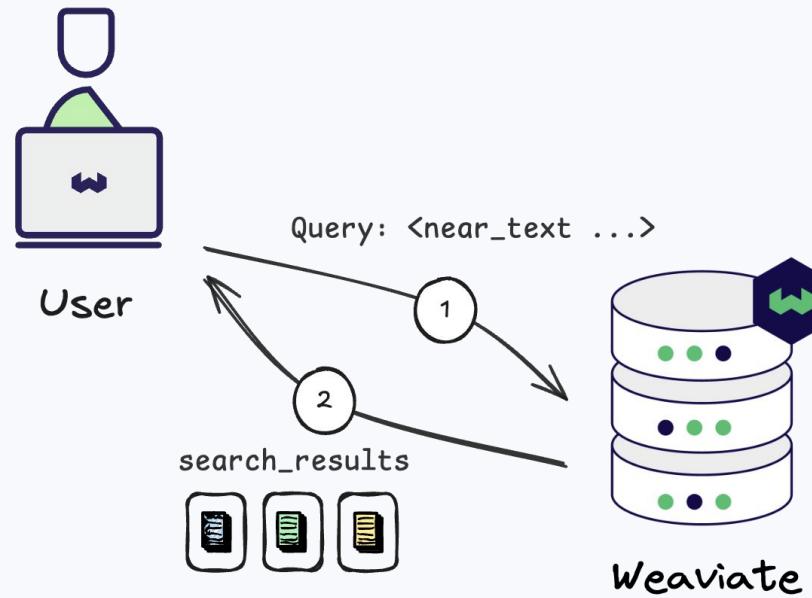


# What's happening here?

Multimodal RAG

# RAG

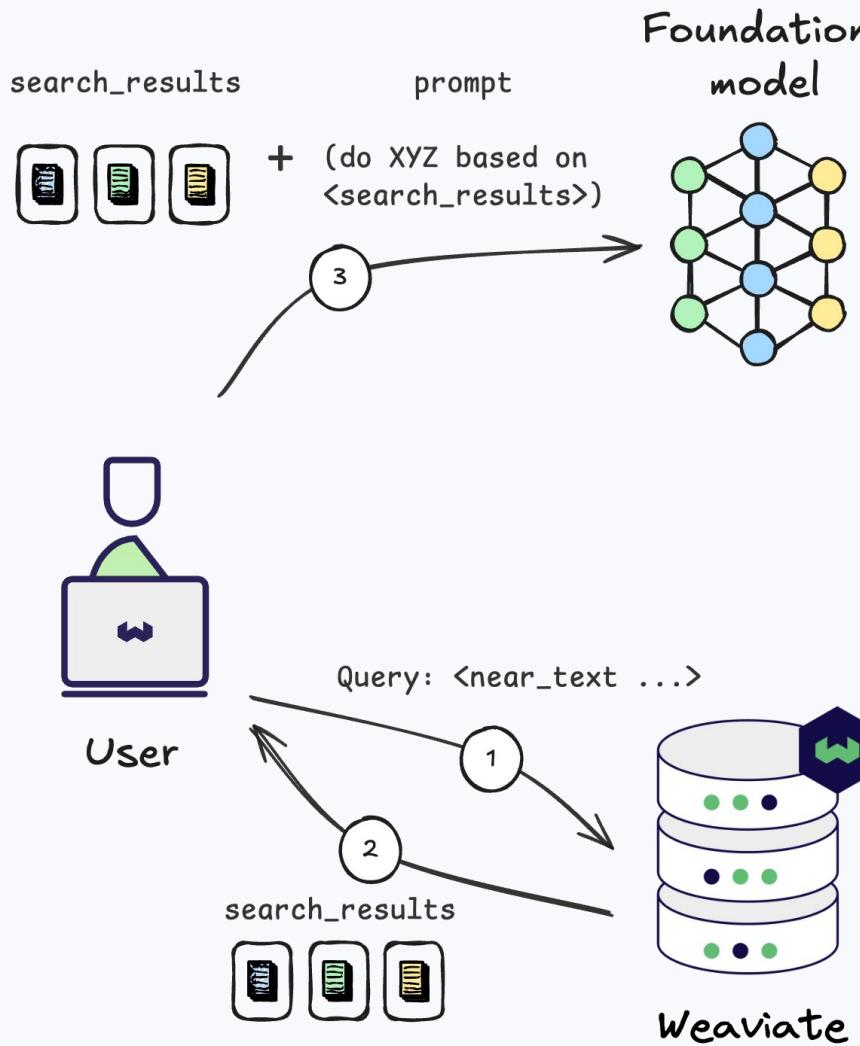
(Retrieval  
augmented  
generation)





# RAG

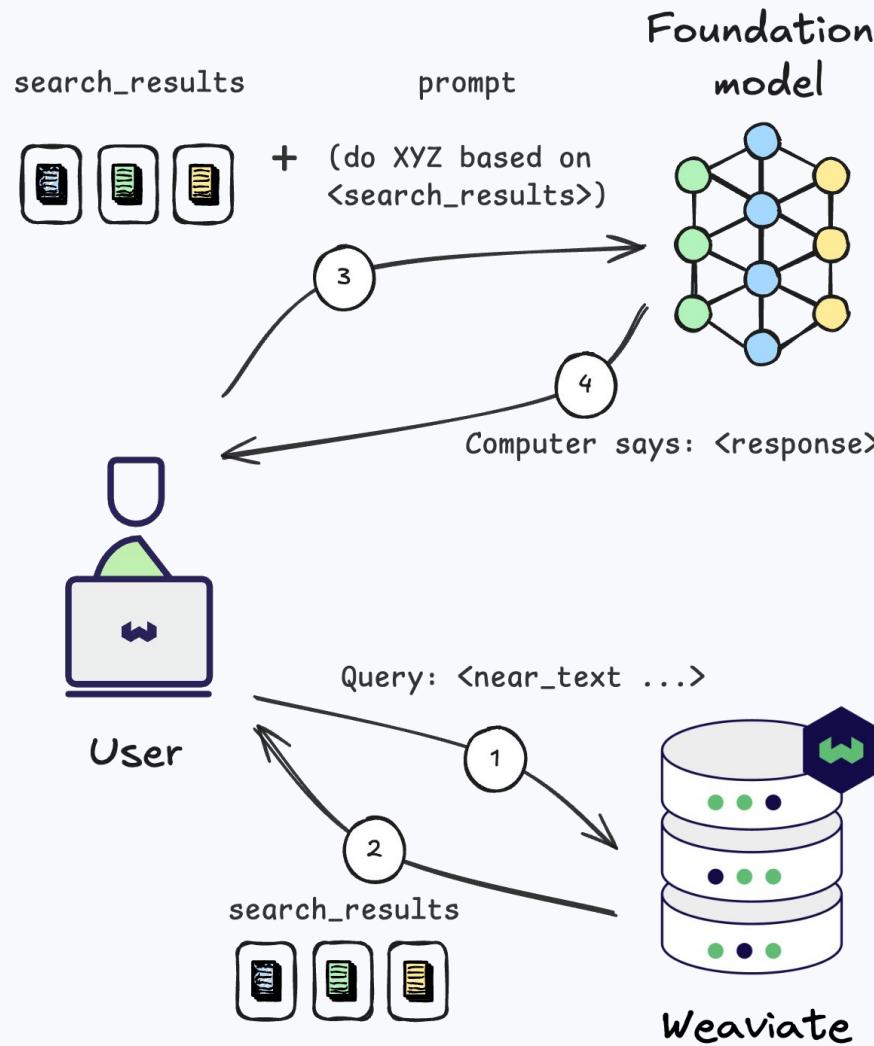
(Retrieval  
augmented  
generation)





# RAG

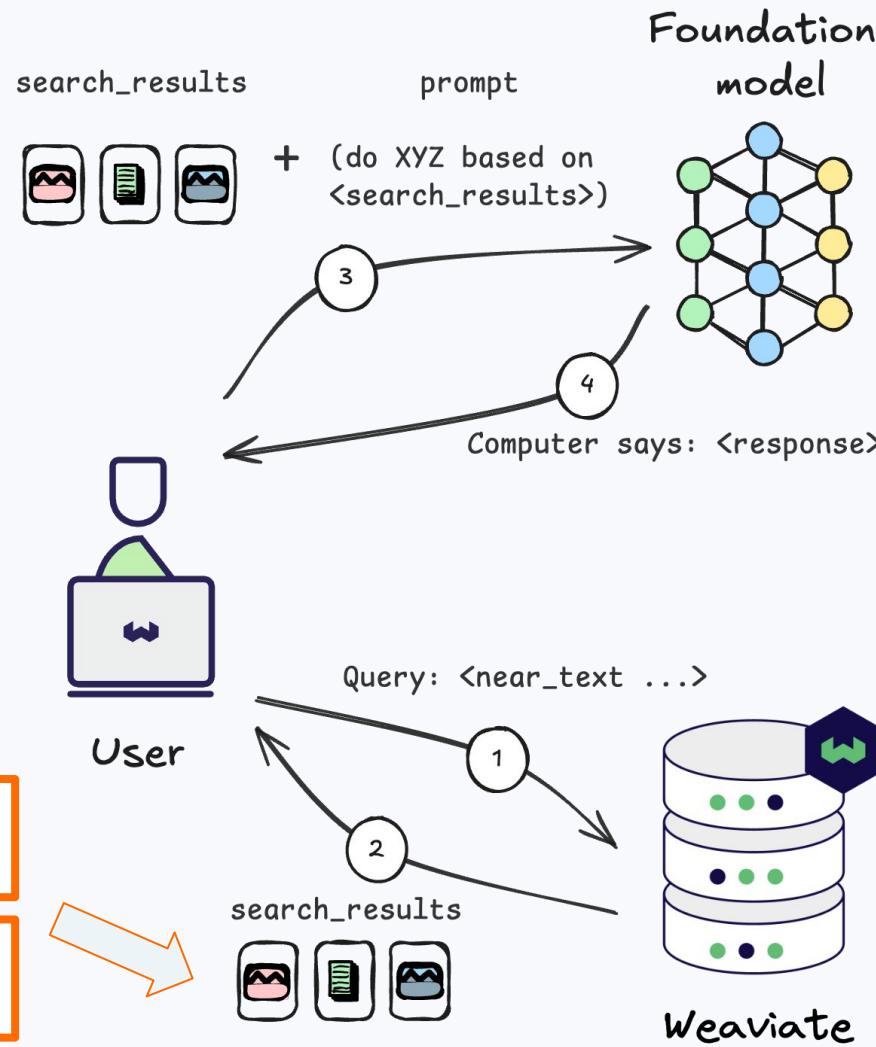
(Retrieval  
augmented  
generation)





# RAG

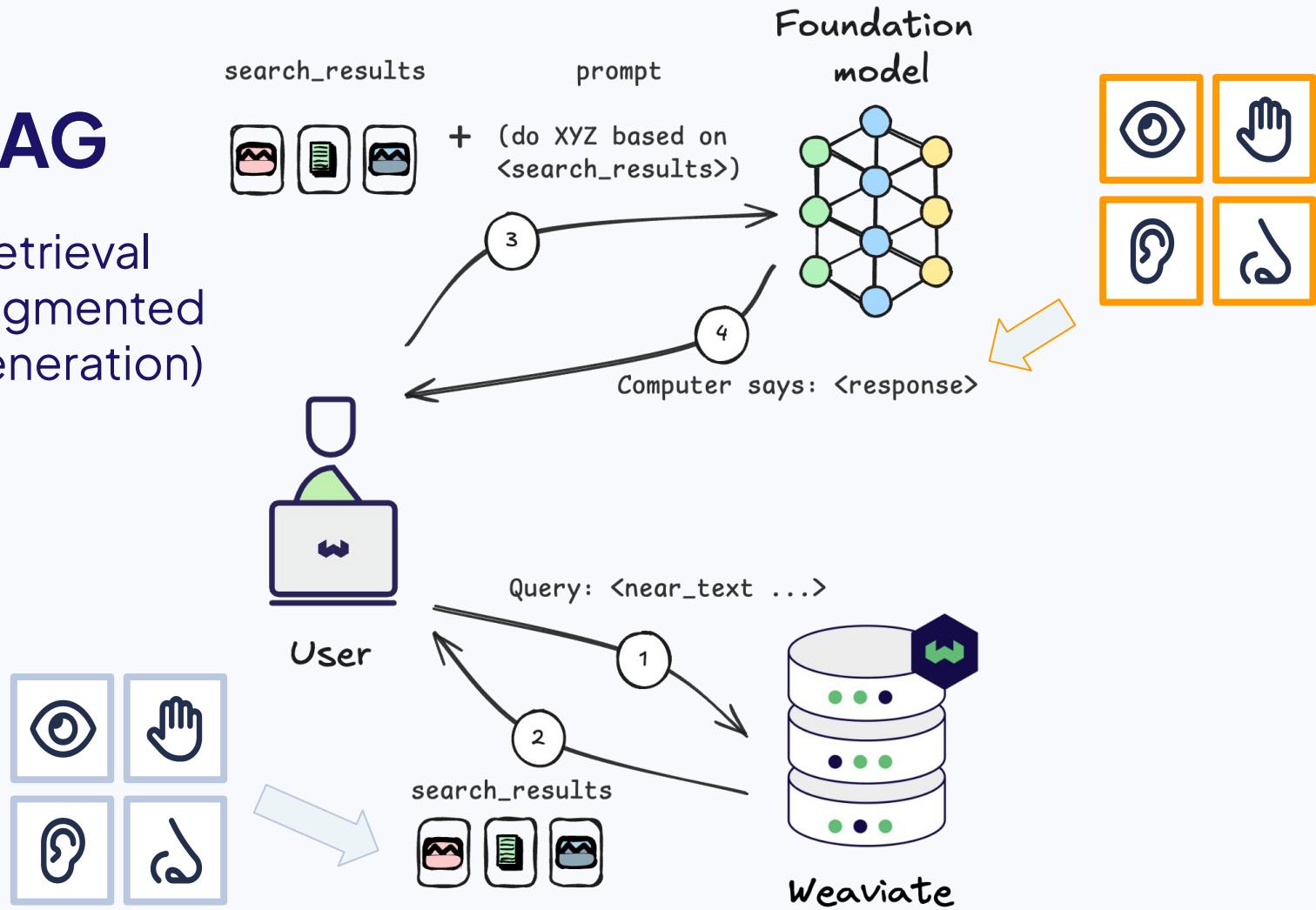
(Retrieval  
augmented  
generation)





# RAG

(Retrieval  
augmented  
generation)





# Suggested solution

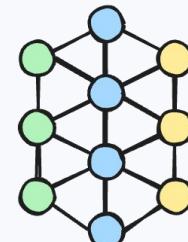


## Foundation model

search\_results

prompt

+ (do XYZ based on  
<search\_results>)



Easier insights from data

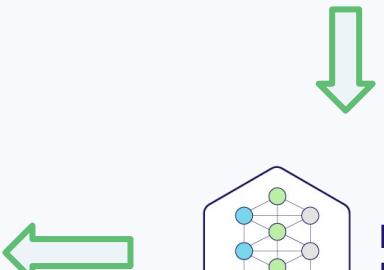
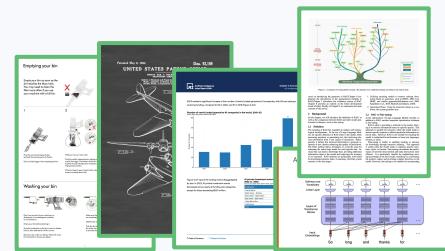
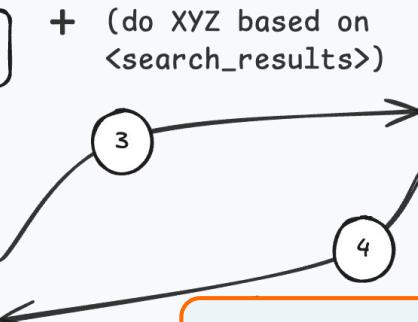
User

Query: <near\_text ...>

search\_results



Weaviate



Embedding  
Model



# Agentic AI Databases



# 1 A little bit of history of how we got to ‘Agentic AI’

## Let’s define “Agents”

- ★ A system that includes LLM capable of making decisions on the best course of action. (an LLM we’ve offloaded the burden of planning to)
- ★ A system that can execute and interact with its environment based on generated course of action.



## 1 A little bit of history of how we got to ‘Agentic AI’

### Let’s remember

- ★ Language models do not know ✨ everything ✨..BUT..
- ★ They work remarkably well with good instructions and the right context
- ★ The instruction defines the task: QA, summarization, planning...?



1

# A little bit of history of how we got to ‘Agentic AI’

## ReAct: Reason & Act

★ Language models do not know  
★ how★ to do everything



### REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS

Shunyu Yao<sup>\*,1</sup>, Jeffrey Zhao<sup>2</sup>, Dian Yu<sup>2</sup>, Nan Du<sup>2</sup>, Izhak Shafran<sup>2</sup>, Karthik Narasimhan<sup>1</sup>, Yuan Cao<sup>2</sup>

<sup>1</sup>Department of Computer Science, Princeton University

<sup>2</sup>Google Research, Brain team

<sup>1</sup>{shunyuy, karthikn}@princeton.edu

<sup>2</sup>{jeffreyzhao, dianyu, dunan, izhak, yuancao}@google.com

#### ABSTRACT

While large language models (LLMs) have demonstrated impressive performance across tasks in language understanding and interactive decision making, their abilities for reasoning (e.g. chain-of-thought prompting) and acting (e.g. action plan generation) have primarily been studied as separate topics. In this paper, we explore the use of LLMs to generate both reasoning traces and task-specific actions in an interleaved manner, allowing for greater synergy between the two: reasoning traces help the model induce, track, and update action plans as well as handle exceptions, while actions allow it to interface with and gather additional information from external sources such as knowledge bases or environments. We apply our approach, named ReAct, to a diverse set of language and decision making tasks and demonstrate its effectiveness over state-of-the-art baselines in addition to improved human interpretability and trustworthiness. Concretely, on question answering (HotpotQA) and fact verification (Fever), ReAct overcomes prevalent issues of hallucination and error propagation in chain-of-thought reasoning by interacting with a simple Wikipedia API, and generating human-like task-solving trajectories that are more interpretable than baselines without reasoning traces. Furthermore, on two interactive decision making benchmarks (ALFWORLD and WebShop), ReAct outperforms imitation and reinforcement learning methods by an absolute success rate of 34% and 10% respectively, while being prompted with only one or two in-context examples.

#### 1 INTRODUCTION

A unique feature of human intelligence is the ability to seamlessly combine task-oriented actions with verbal reasoning (or inner speech, Alderson-Day & Fernyhough, 2015), which has been theorized to play an important role in human cognition for enabling self-regulation or strategization (Vygotsky, 1987; Luria, 1965; Fernyhough, 2010) and maintaining a working memory (Baddeley, 1992). Con-

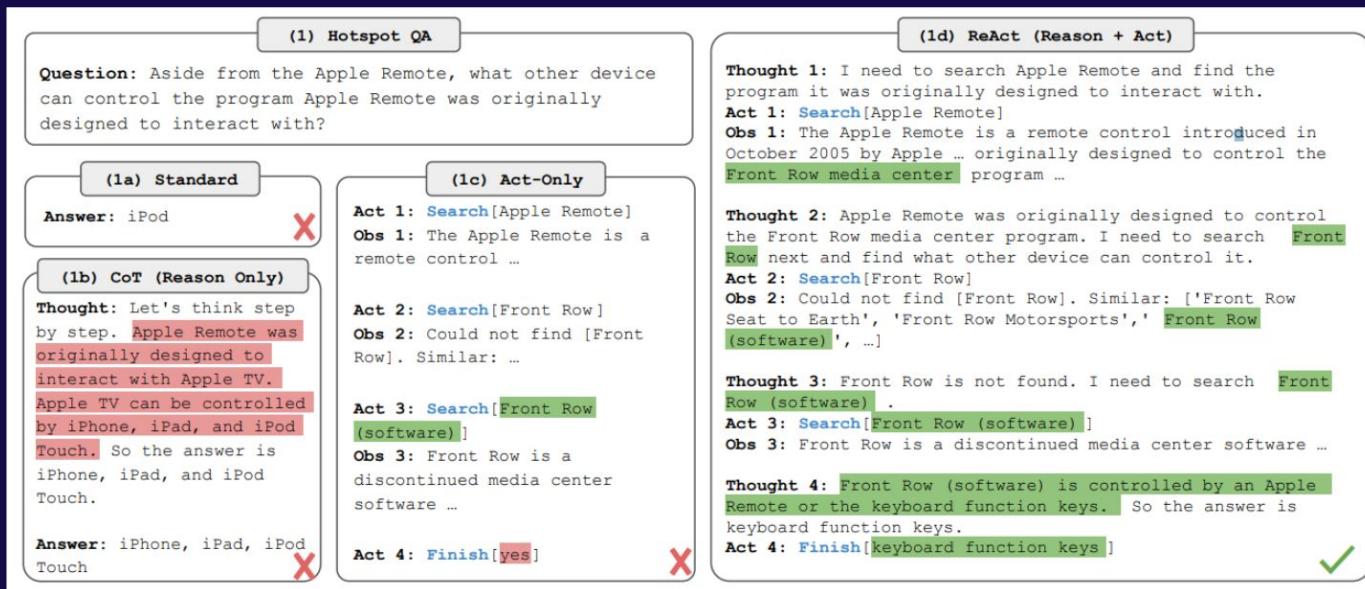


## 1

# A little bit of history of how we got to ‘Agentic AI’

## ReAct: Reason & Act

★ Given the right instruction, and access to the right ⚡tools⚡ – they may reason about what action to take to solve a task





## 2 Functions/Tool Calling

- ★ Models can “reason” about which tool to run, and..
- ★ They can produce an output about *how* to run the tool

## 2

## Functions/Tool Calling



Name: weaviate\_docs\_rag

Description: Useful for finding answers in Weaviate documentation  
Inputs: {query: str}

Outputs: {answer:str}

Name: weather\_tool

Description: Useful for weather information  
Inputs: {city: str, country: str, date: datetime}  
Outputs: {temp: int, unit: str, location: str}



## 2 Functions/Tool Calling



**Name:** weavite\_docs\_rag  
**Description:** Useful for finding answers in Weavite documentation  
**Inputs:** {query: str}  
**Outputs:** {answer:str}

**Name:** weather\_tool  
**Description:** Useful for weather information  
**Inputs:** {city: str, country: str, date: datetime}  
**Outputs:** {temp: int, unit: str, location: str}





## 2 Functions/Tool Calling



Run: weather\_tool  
Inputs: {city: Berlin, country: DE, date: 12/12/2024}

**Name:** weavite\_docs\_rag  
**Description:** Useful for finding answers in Weavite documentation  
**Inputs:** {query: str}  
**Outputs:** {answer:str}

**Name:** weather\_tool  
**Description:** Useful for weather information  
**Inputs:** {city: str, country: str, date: datetime}  
**Outputs:** {temp: int, unit: str, location: str}

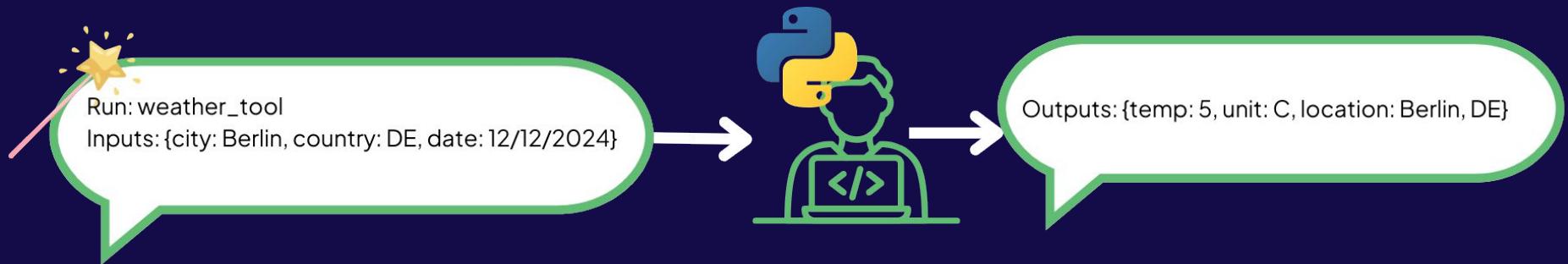
The execution is up to us





## 2 Functions/Tool Calling

### Tool Calling Step





## 2 Functions/Tool Calling

### Observation Step



Outputs: {temp: 5, unit: C, location: Berlin, DE}





## 2 Functions/Tool Calling

### Thought & Answer Step



The question was about what the weather is like in Berlin and it seems and I have the answer: 5C



The weather in Berlin is 5 degrees Celcius



## 2

## Functions/Tool Calling

```
from openai import OpenAI

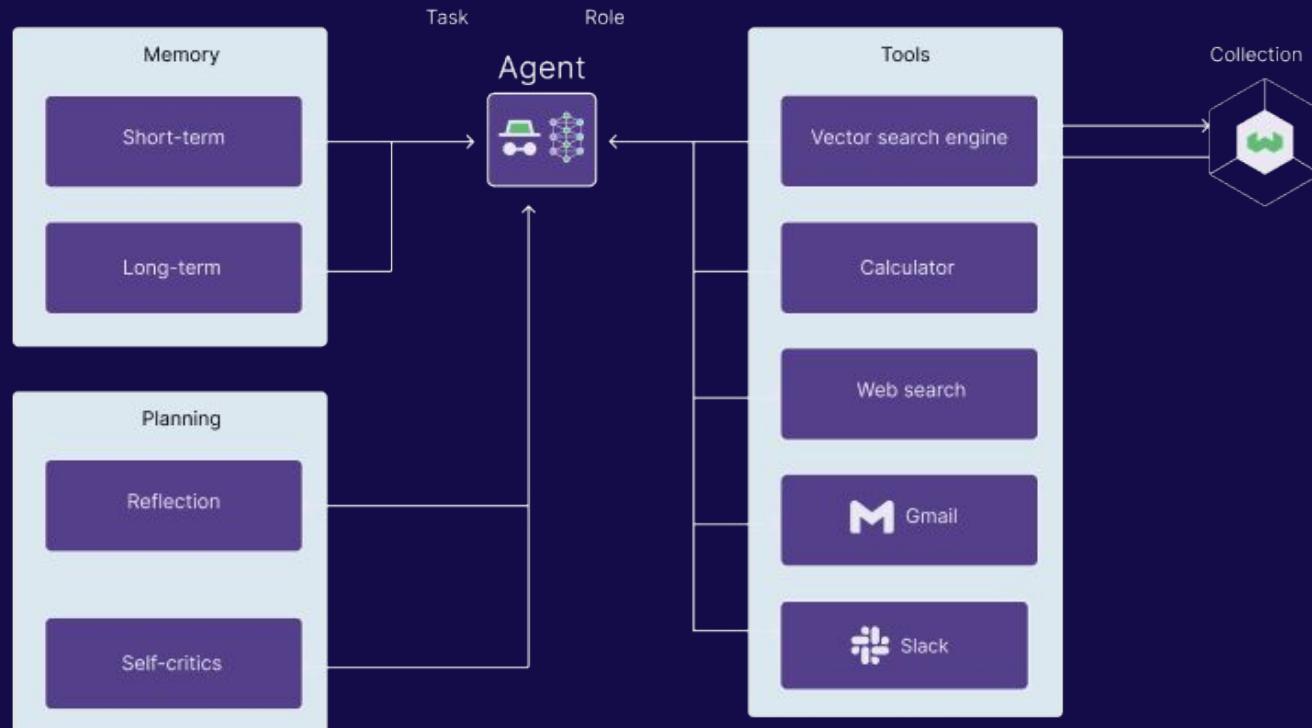
client = OpenAI()

tools = [
    {
        "type": "function",
        "function": {
            "name": "get_delivery_date",
            "description": "Get the delivery date for a customer's order. Call this whenever you\n                need to know the delivery date, for example when a customer asks\n                    'Where is my package?'",
            "parameters": {
                "type": "object",
                "properties": {
                    "order_id": {
                        "type": "string",
                        "description": "The customer's order ID.",
                    },
                },
                "required": ["order_id"],
                "additionalProperties": False,
            },
        },
    }
]
```



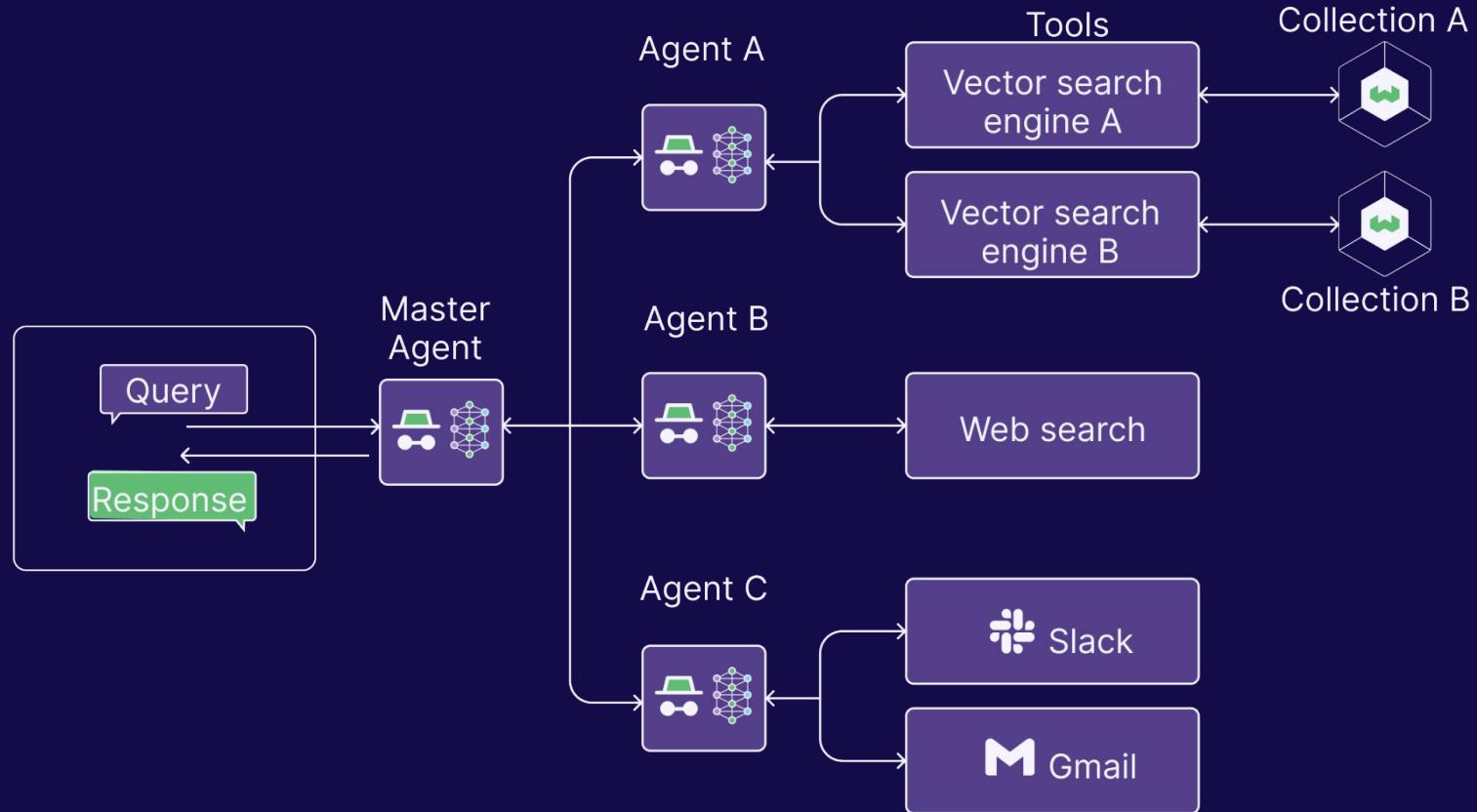


# ✨ Agents ✨



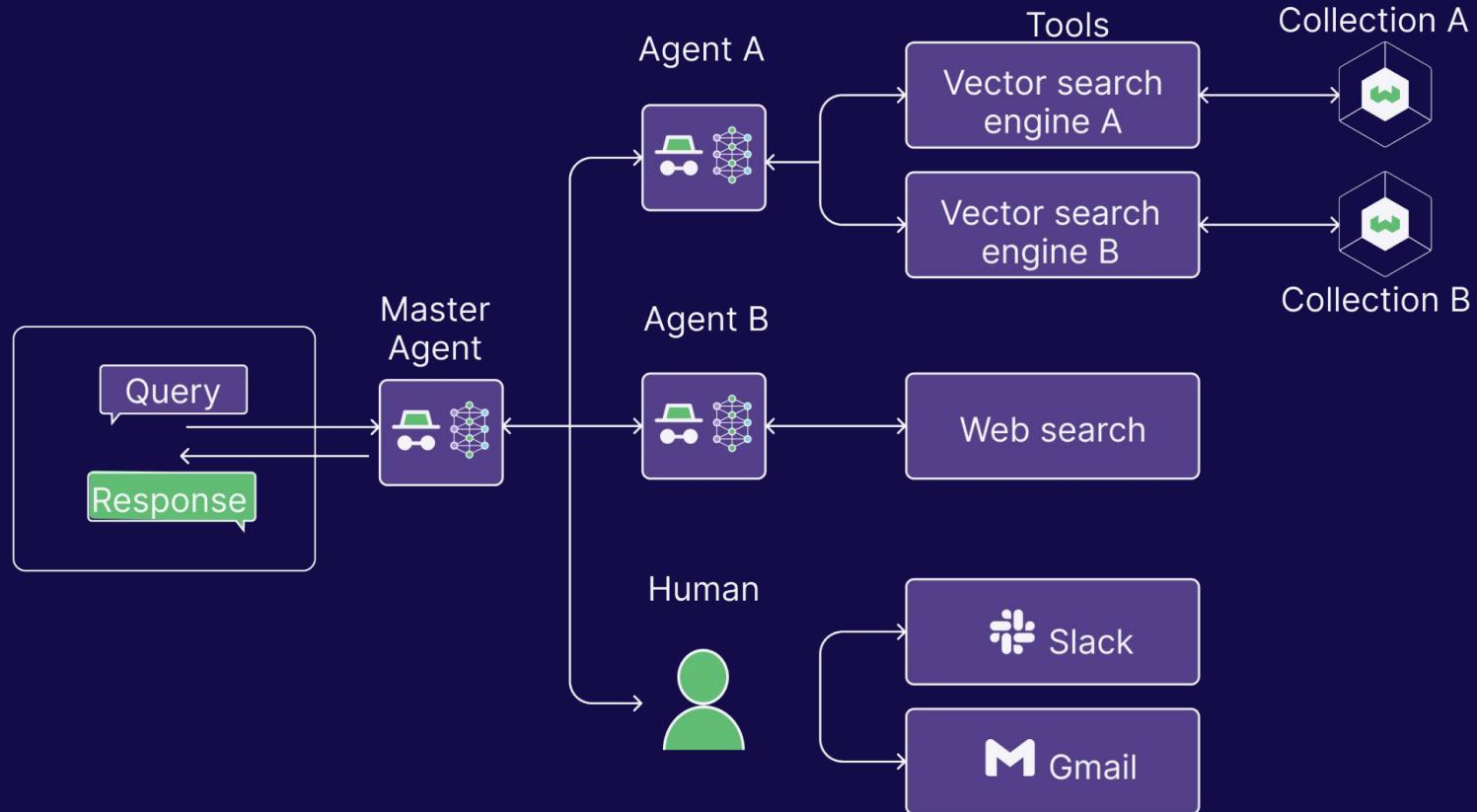


# ✨ Agents ✨





# ✨ Agents ✨





# Quick notes

Weaviate at a high level +  
Features we didn't cover



# (Multi-) Tenancy

Single vs multi-tenancy architectures. Use cases for each, and how to manage resources



## Single tenant

When your use case requires all data in a single collection

- Store and search through **Multi-billion scale** collections
- Process millions of updates on a daily basis
- Smart (re) vectorization techniques
- Combine keyword, and vector search for powerful **hybrid search**
- Sharding and replication allows for tuning throughput
- Fast filtering capabilities to narrow down results



## Multi tenant

When your use case needs to isolate data per user

- Scales to **millions** of tenants
- Full isolation
- Tenants can be **Cold, Warm, or Hot**, depending on activity
- Control tenant status through APIs or automatically
- Store Frozen tenants on cloud storage to reduce cost
- Allows for on-disk tenants and in-memory depending on use case



# Split Data per (Domain)

## User

Each user with their own data

## Customer

White labelled-solutions

## Team / Department

Separate HR docs from legal docs



# Split Data per (Resource)

## PDF

Tenant for PDF chunks

## Video

Tenant for Video chunks

## Group of Documents

Tenant per topic

## Videos

Tenant per set of videos  
(i.e. course lessons)



# Split Data per (Time)

**Daily**

Tenant per day

**Monthly**

Tenant per month



# Built-in Multi-tenancy

## Smaller Indexes

Index per tenant

## Faster queries

Less to search

## GDPR compliant

For User tenants

## Tenant memory management

Hot and Cold

## Flat Index – friendly

Save RAM

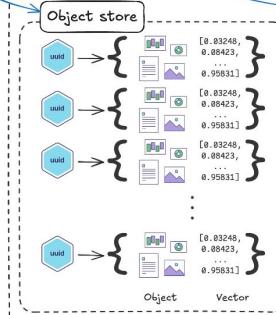


Each collection stores  
source objects + vectors

Each collection has  
sets of indexes

Each collection has its own  
metadata / configuration

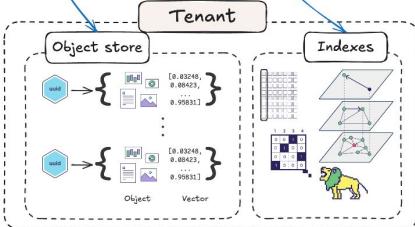
### Collection (without multi-tenancy)



Each tenant has its own  
object & vector store

Each tenant's indexes are  
separate from the others'

### Collection (multi-tenant)

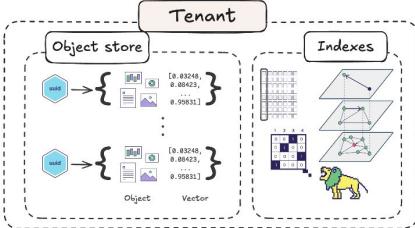


#### Metadata

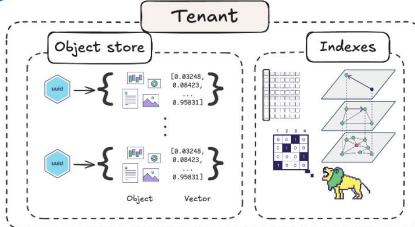
- Data schema
  - Integration settings
    - Vectorizer
    - Reranker

- Indexing settings
  - Vector index
  - Inverted index

- Replication settings
- And more...



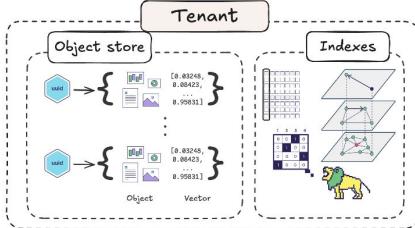
Tenants share the metadata /  
Collection configuration



#### Metadata

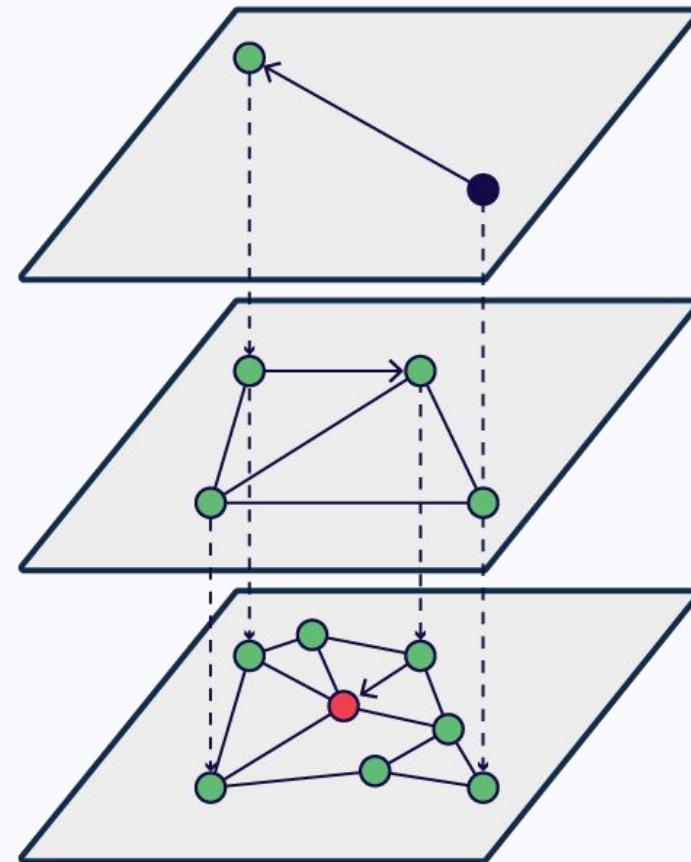
- Data schema
  - Integration settings
    - Vectorizer
    - Reranker

- Indexing settings
  - Vector index
  - Inverted index



# Indexing techniques

- **HNSW:** high-performance, in-memory. Scales well, even for very large data sets.
- **Flat:** disk-based indexes that perform brute-force vector searches. Best for small data sets.
- **Dynamic:** Automatically switches from flat to HNSW when a collection (or tenant) reaches a threshold size.



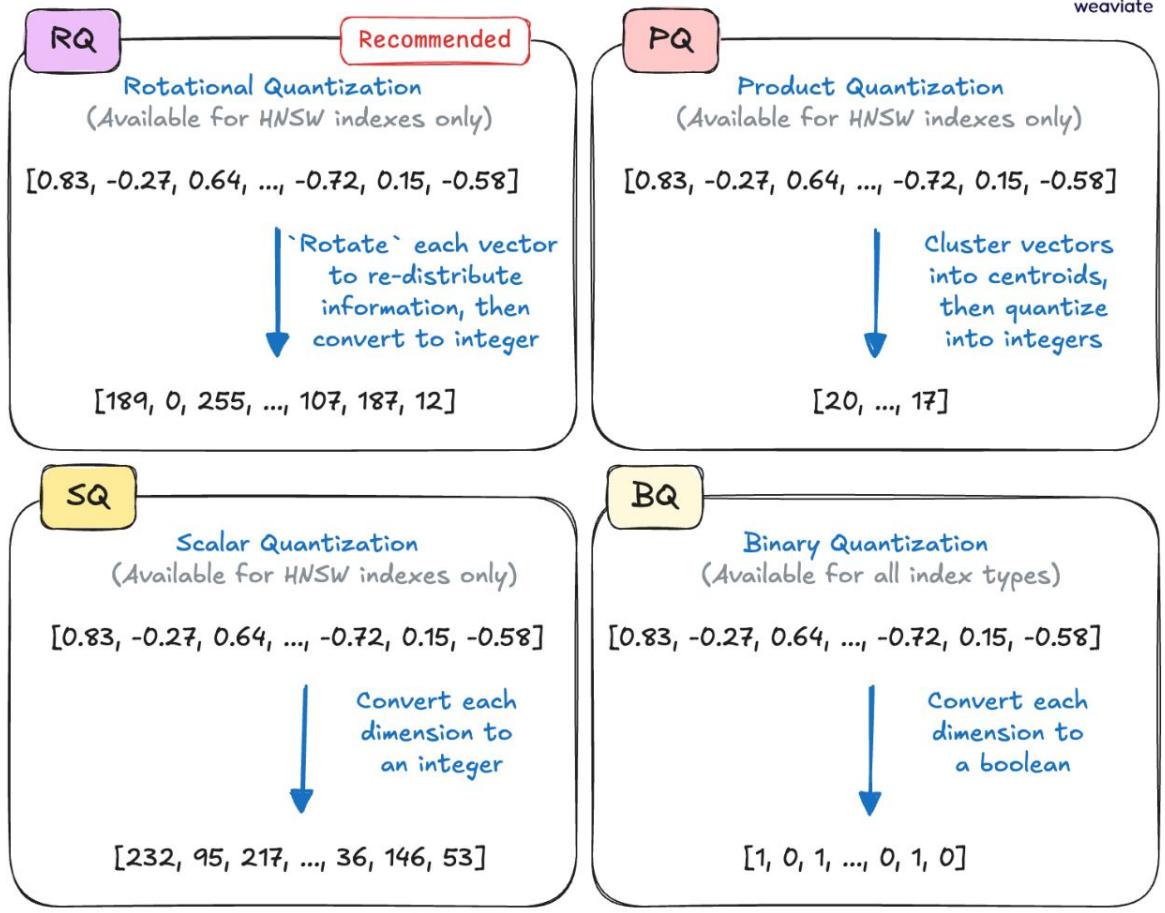
# Quantization at a glance

Vectors can be quantized to reduce their size and improve search speeds.

Weaviate offers a choice of four quantization algorithms:

- Rotational
- Product
- Scalar
- Binary

Generally, RQ (or PQ) is the best choices for most users.



 1024 gib

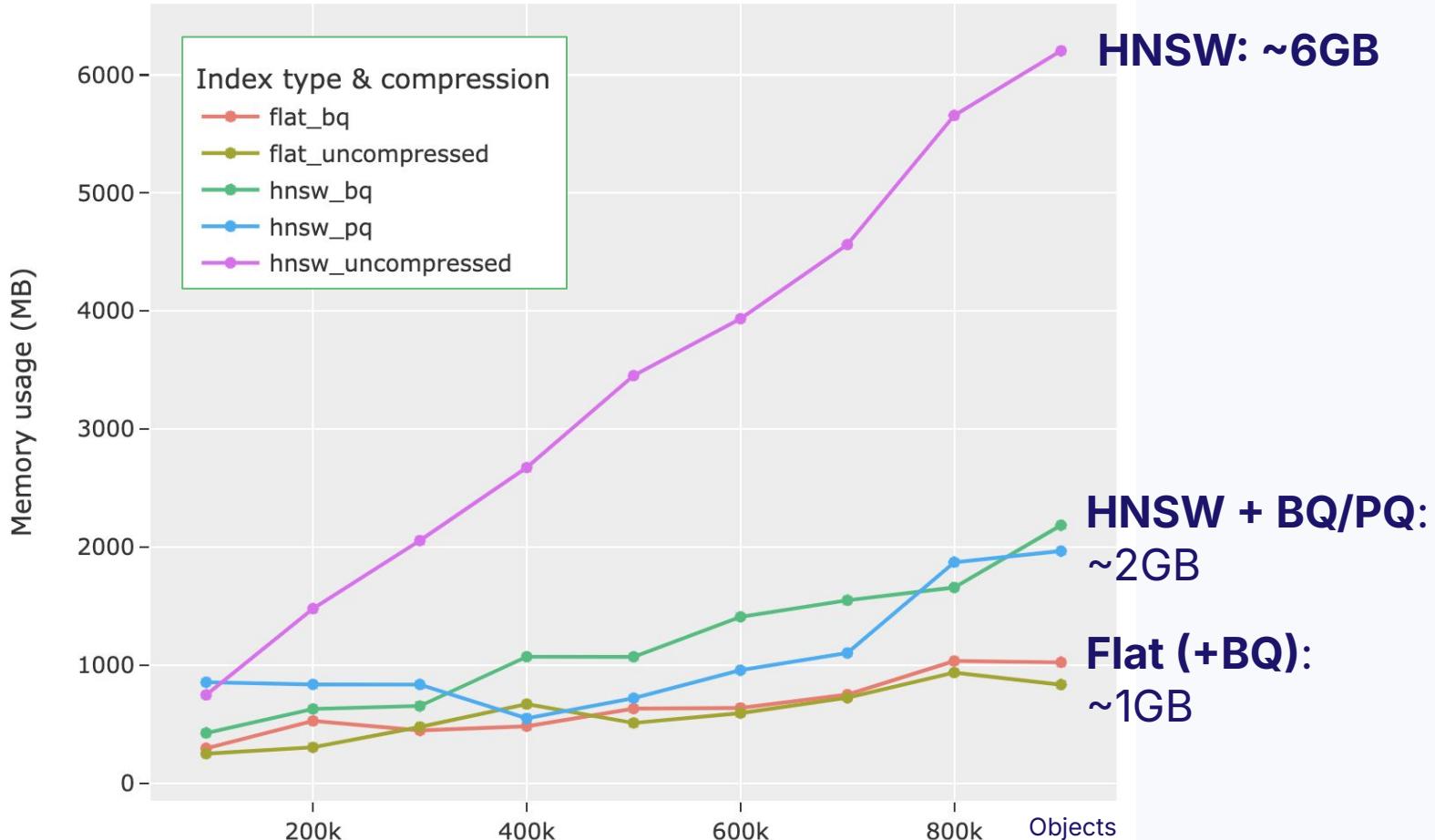
19 matches

Instance name ▽	On-Demand hourly rate ▲	vCPU ▽	Memory ▽	Storage ▽
x2gd.metal	\$5.344	64	1024 GiB	2 x 1900 SSD
x2gd.16xlarge	\$5.344	64	1024 GiB	2 x 1900 SSD
hpc6id.32xlarge	\$5.70	64	1024 GiB	4 x 3800 NVMe SSD
x2iedn.8xlarge	\$6.669	32	1024 GiB	1 x 950 NVMe SSD
x2idn.16xlarge	\$6.669	64	1024 GiB	1 x 1900 NVMe SSD

AWS Spot pricing, Nov 2024



# Example: Index type & quantization vs memory footprint





Instance type

Memory

**\$58k / year**

View



**VS.**

**\$23k / year**

vCPU

32

< 1 2 >

Instance name	Demand hourly rate	vCPU	Memory	Storage	Network performance
x2iedn.8xlarge	\$6.669	32	1024 GiB	1 x 950 NVMe SSD	25 Gigabit
x1e.8xlarge	\$6.672	32	976 GiB	1 x 960 SSD	Up to 10 Gigabit
x2gd.8xlarge	\$2.672	32	512 GiB	1 x 1900 SSD	12 Gigabit

AWS Spot pricing, Nov 2024



# Q & A

Ask us anything!



# Thank you



[weaviate.io](https://weaviate.io)



[weaviate/weaviate](https://github.com/weaviate/weaviate)



[linkedin](#)



[@weaviate\\_io](#)