

# Understanding Vector Databases

PyLadies Meetup NYC

# Jessica Garson

@JessicaGarson

@JessicaGarson@macaw.social

Senior Developer Advocate



**This talk is aimed for developers who may have  
some experience with machine learning but looking  
to skill up.**

This talk will explore the traditional  
"Classic" search, highlight its constraints,  
and provide an overview of vector  
databases.

# Elasticsearch: You Know, for Search



elasticsearch

lucene

What's the well-known sentence of this scene  
in Star Wars?



66

These are not the droids  
you are looking for.

```
GET /_analyze  
{  
  "char_filter": [ "html_strip" ],  
  "tokenizer": "standard",  
  "filter": [ "lowercase", "stop", "snowball" ],  
  "text": "These are <em>not</em> the droids  
you are looking for."  
}
```

```
"char_filter": "html_strip"
```

These are **<em>not</em>** the droids you are looking for.



These are not the droids you are looking for.

"tokenizer": "**standard**"

These are not the droids you are looking for.

These  
are  
not  
the  
droids  
you  
are  
looking  
for



"filter": "lowercase"

These  
are  
not  
the  
droids  
you  
are  
looking  
for



these  
are  
not  
the  
droids  
you  
are  
looking  
for

```
"filter": "stop"
```

These  
are  
not  
the  
droids  
you  
are  
looking  
for

→

droids  
you  
looking

"filter": "snowball"

droids  
you  
looking



droid  
you  
look

These are <em>not</em> the **droids you** are **looking** for.

```
{ "tokens": [ {  
    "token": "droid",  
    "start_offset": 27, "end_offset": 33,  
    "type": "<ALPHANUM>", "position": 4  
} , {  
    "token": "you",  
    "start_offset": 34, "end_offset": 37,  
    "type": "<ALPHANUM>", "position": 5  
} , {  
    "token": "look",  
    "start_offset": 42, "end_offset": 49,  
    "type": "<ALPHANUM>", "position": 7  
} ] }
```

## PROS

- Fast and cheap
- Scalable
- Easy to understand

## CONS

- Synonyms and homonyms
- Images, videos, audio
- User intent

# Semantic Search: Meaning, not literal matches

# X-wing starfighter squadron



**What ships and crews  
do I need to destroy an  
almost finished death  
star?**

**Or is there a secret  
weakness?**



# Elasticsearch: You Know, for **Vector** Search



# What do you need to get started with vector search?

- Elasticsearch instance optimized for machine learning (8.13 or later)
- Python 3.8 (or later)
- Elasticsearch Python Client

A close-up photograph of a black cat's face, partially hidden in dense, dark green grass. The cat has intense, glowing green eyes that reflect light. Its fur is dark and shiny, and its whiskers are clearly visible. The background is out of focus, creating a sense of depth.

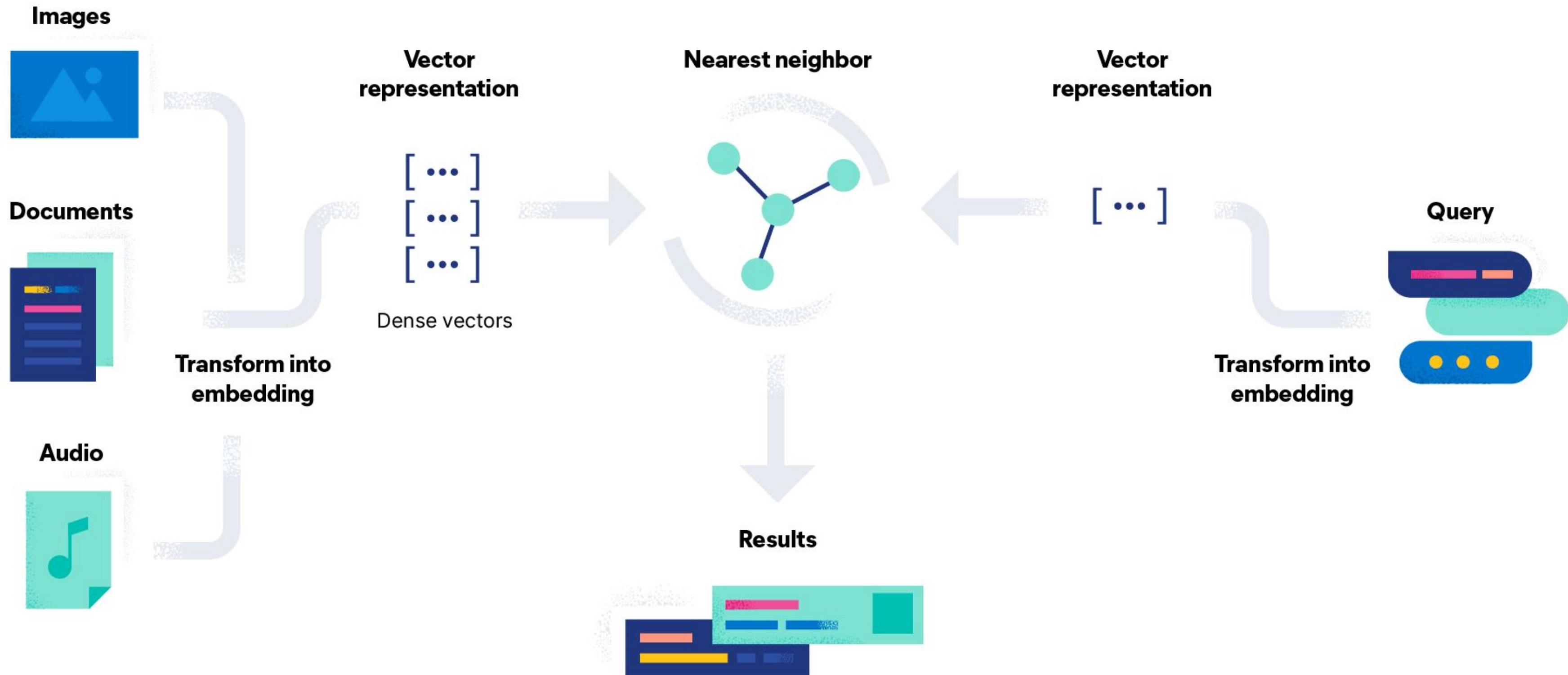
What is a vector  
database?

# Why use vector search?

Vector search provides the foundation for implementing semantic search for text or similarity search for images, videos, or audio.

A vector database allows you to find similar data quickly.

# Searching in a vector database



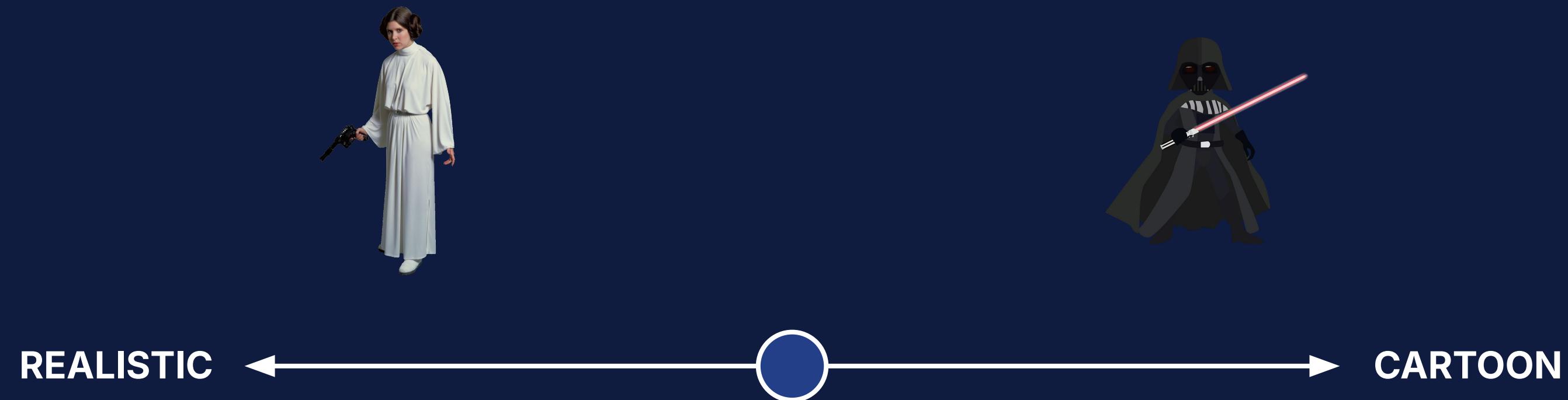
# What's a vector?

# What are embeddings?

Embeddings leverage a machine learning model to translate text into numbers, allowing you to perform vector searches.

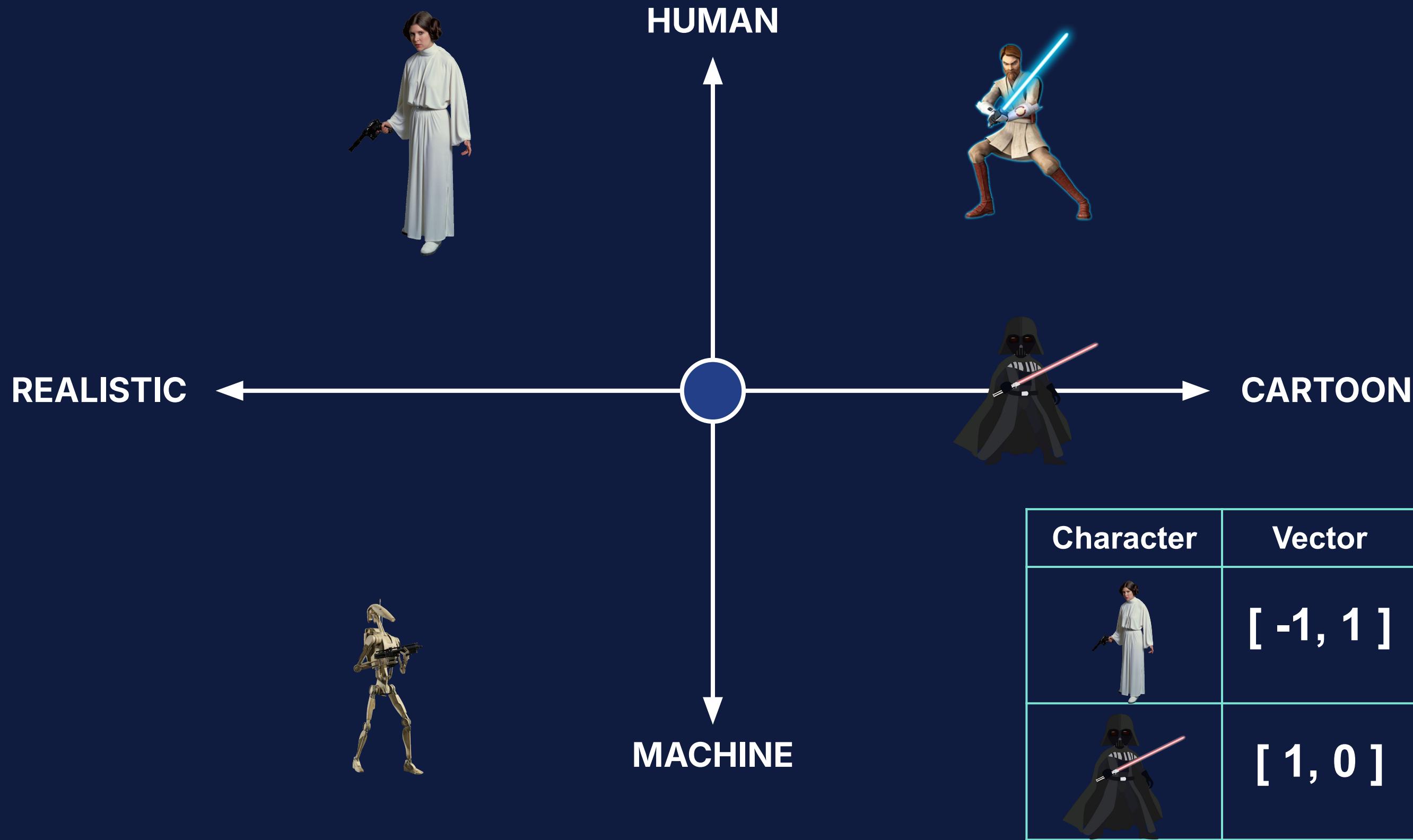
# Embeddings represent your data

Example: 1-dimensional vector

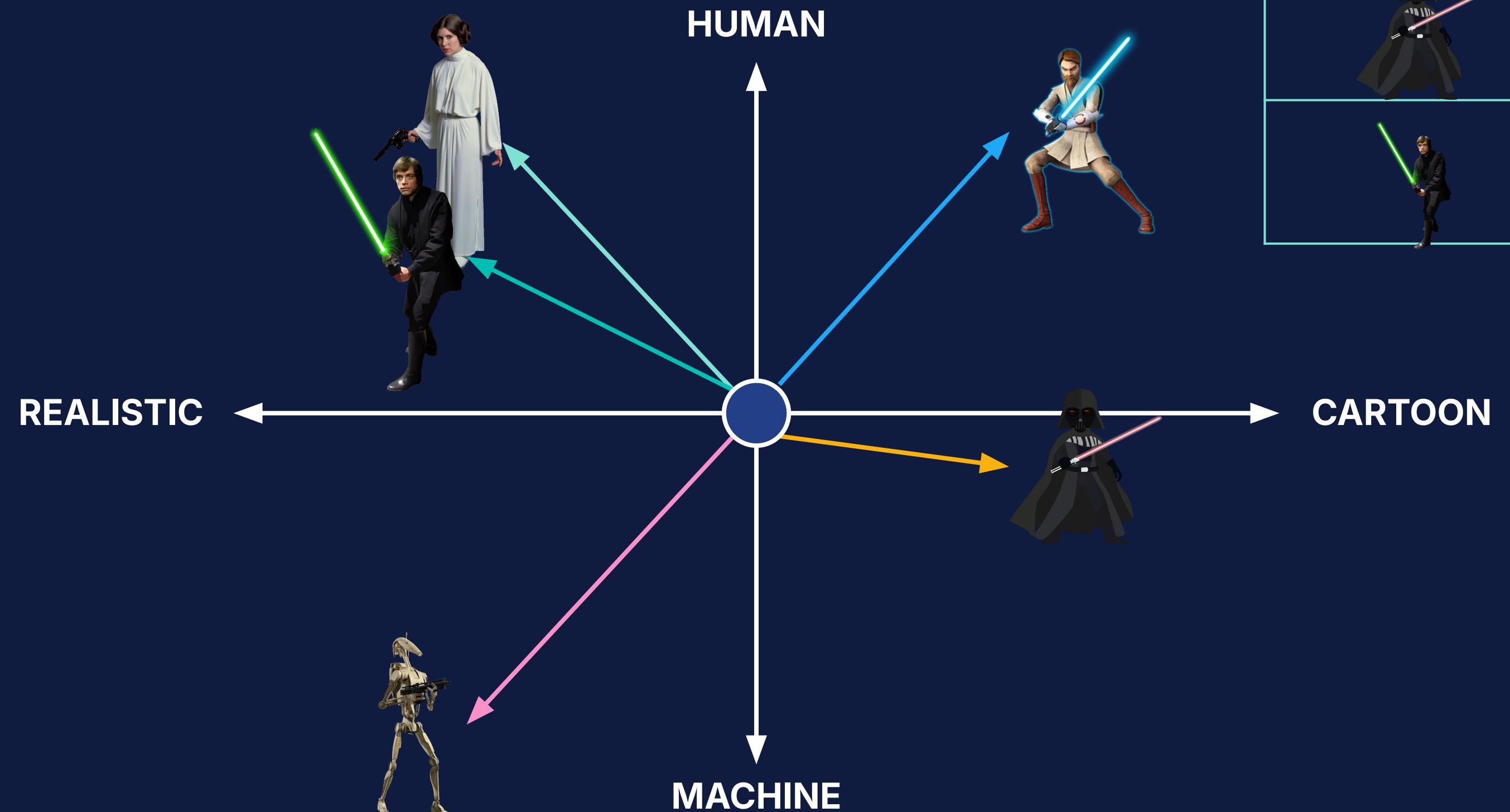


Character	Vector
	[ -1 ]
	[ 1 ]

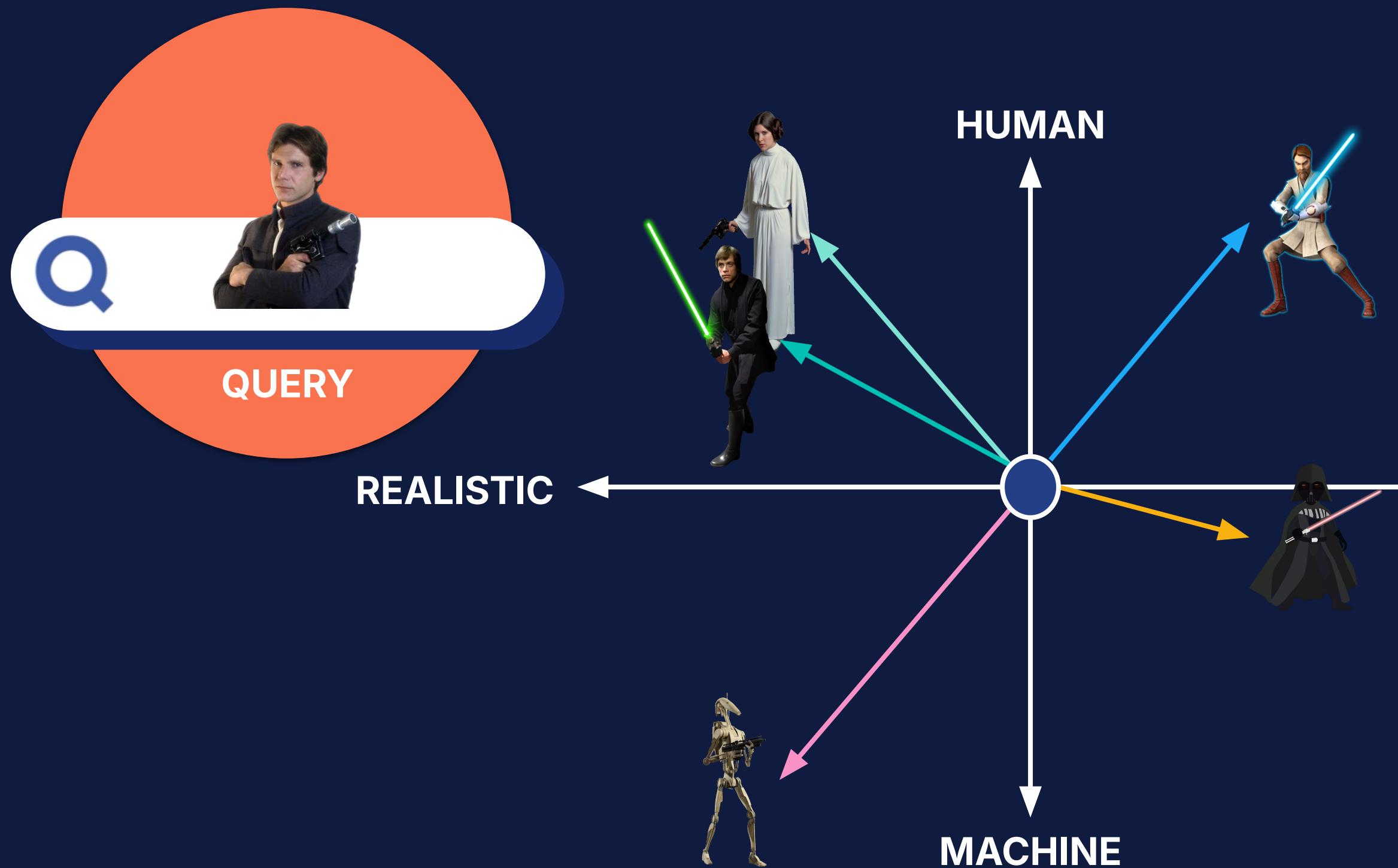
# Multiple dimensions represent different data aspects



# Similar data is grouped together



# Vector search ranks objects by similarity (relevance) to the query



Relevance	Result
Query	
1	
2	
3	
4	
5	

# How to select a model?

# Choice of Embedding Model

## Start with Off-the Shelf Models

- Text data: Hugging Face (like Microsoft's E5)
- Images: OpenAI's CLIP

## Extend to Higher Relevance

- Apply hybrid scoring
- Bring Your Own Model: requires expertise + labeled data

# What is chunking?

Chunking refers to the process of breaking content up.

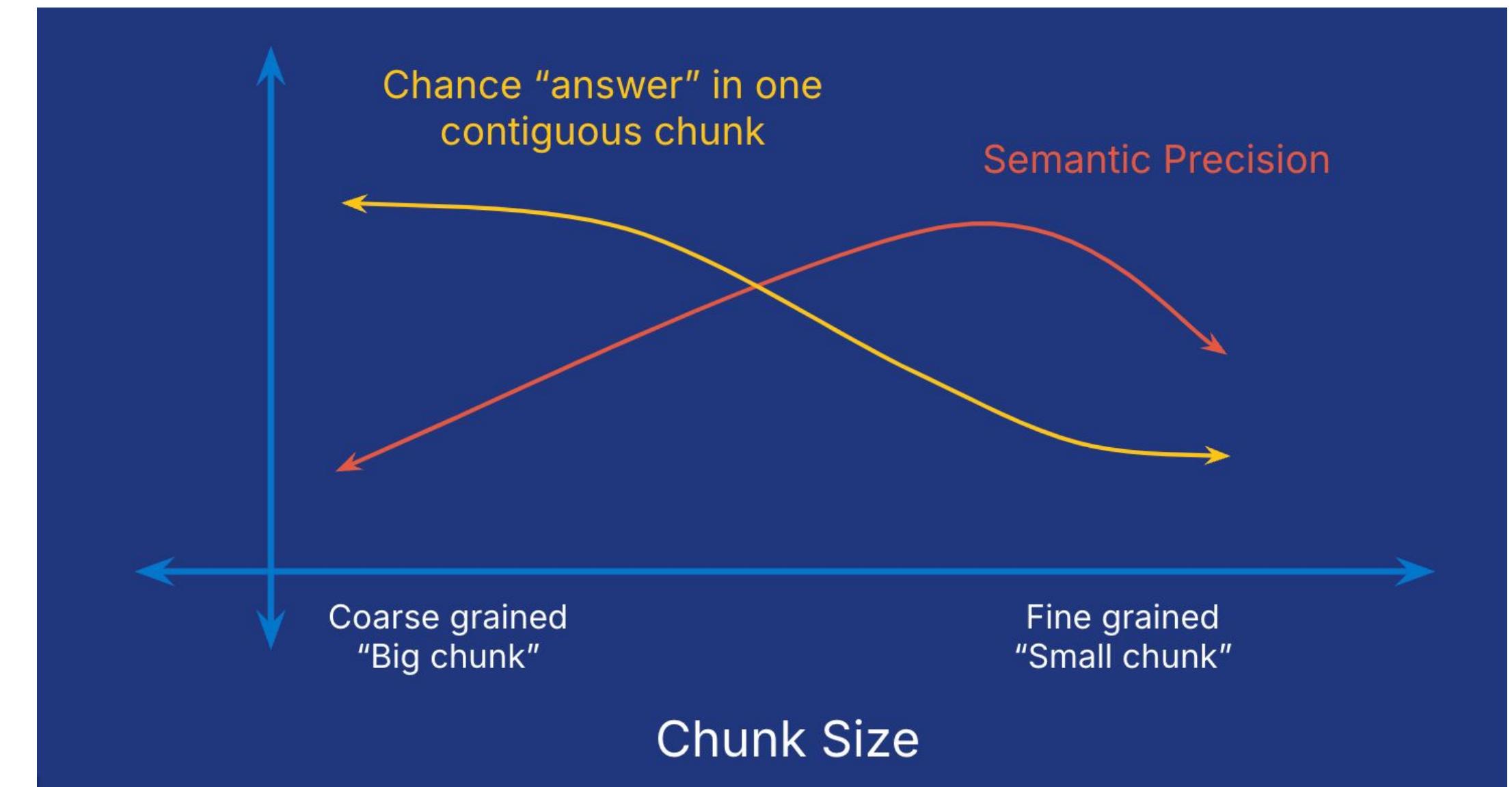
It is necessary for representing complex texts in vector embeddings.



```
helpers.bulk(es, actions, pipeline="text-embedding", chunk_size=1000)
```



# Impact of chunk size



# How and when to run a hybrid search?

# Demo



# Demo

[https://www.elastic.co/search-labs/blog/elastic-vector-data  
base-practical-example](https://www.elastic.co/search-labs/blog/elastic-vector-data-base-practical-example)



# Let's take a look at the book data

# Pipeline



```
import json
from getpass import getpass

from elasticsearch import Elasticsearch, helpers
from sentence_transformers import SentenceTransformer
```



```
cloud_endpoint = getpass("Elastic deployment Cloud Endpoint: ")
cloud_api_key = getpass("Elastic deployment API Key: ")
INDEX_NAME = "books-pipeline"

es = Elasticsearch(
    hosts=cloud_endpoint,
    api_key=cloud_api_key,
)
```

```
● ● ●  
  
resp = es.ingest.put_pipeline(  
    id="text-embedding",  
    description="converts book description text to a vector",  
    processors=[  
        {  
            "inference": {  
                "model_id": "sentence-transformers_msмарко-minilm-l-12-v3",  
                "input_output": [  
                    {  
                        "input_field": "book_description",  
                        "output_field": "description_embedding",  
                    }  
                ],  
            }  
        },  
        {  
            "set": {  
                "description": "Index document to 'failed-<index>'",  
                "field": "_index",  
                "value": "failed-{{{_index}}}",  
            }  
        },  
        {  
            "set": {  
                "description": "Set error message",  
                "field": "ingest.failure",  
                "value": "{{_ingest.on_failure_message}}",  
            }  
        },  
    ],  
)  
  
print(resp)
```



```
● ● ●

mappings = {
    "mappings": {
        "properties": {
            "book_title": {"type": "text"},
            "author_name": {"type": "text"},
            "rating_score": {"type": "float"},
            "rating_votes": {"type": "integer"},
            "review_number": {"type": "integer"},
            "book_description": {"type": "text"},
            "genres": {"type": "keyword"},
            "year_published": {"type": "integer"},
            "url": {"type": "text"},
        }
    }
}

# Delete any previous index
es.indices.delete(index=INDEX_NAME)
es.indices.create(index=INDEX_NAME, body=mappings)
print(f"Index '{INDEX_NAME}' created.")
```





```
file_path = "../data/books.json"
with open(file_path, "r") as file:
    books = json.load(file)

# create an array of index actions, with each element holding one document
actions = [
    {"_index": INDEX_NAME, "_id": book.get("id", None), "_source": book}
    for book in books
]

try:
    helpers.bulk(es, actions, pipeline="text-embedding", chunk_size=1000)
    print(f"Successfully added {len(actions)} books into the '{INDEX_NAME}' index.")

except helpers.BulkIndexError as e:
    print(f"Error occurred while ingesting books: {e}")
    print(e.errors)
```

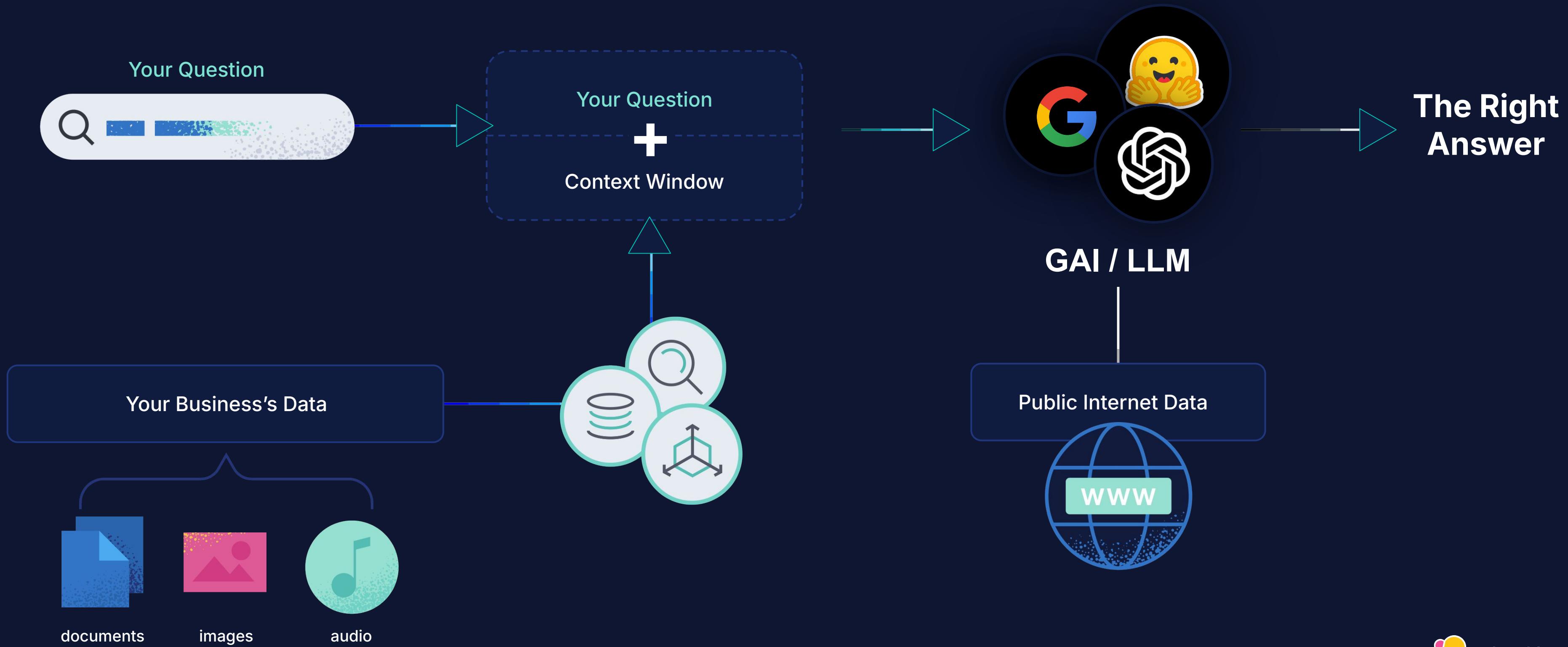


# What is RAG?

Retrieval-augmented generation (RAG) integrates external information retrieval into generating responses by Large Language Models (LLMs).

# Vector Databases vs RAG

# Retrieval Augmented Generation



# Elastic's playground

A low-code interface for you to explore LLMs



```
es_client = Elasticsearch(  
    getpass.getpass("Host: " ),  
    api_key=getpass.getpass("API Key: " ),  
)
```



```
embedding = OpenAIEmbeddings(model="text-embedding-3-large")
elastic_vector_search = ElasticsearchStore(
    index_name="vr_tour_data",
    es_connection=es_client,
    embedding=embedding,
)
```





```
translated_texts = df['translated_text'].tolist()
combined_text = "\n".join(translated_texts)
text_splitter = CharacterTextSplitter(chunk_size=1100, chunk_overlap=0)
docs = text_splitter.split_documents([Document(page_content=combined_text)])
```



# Next steps



RESOURCES FOR DEVELOPERS | BY DEVELOPERS LIKE YOU!

# Search Labs



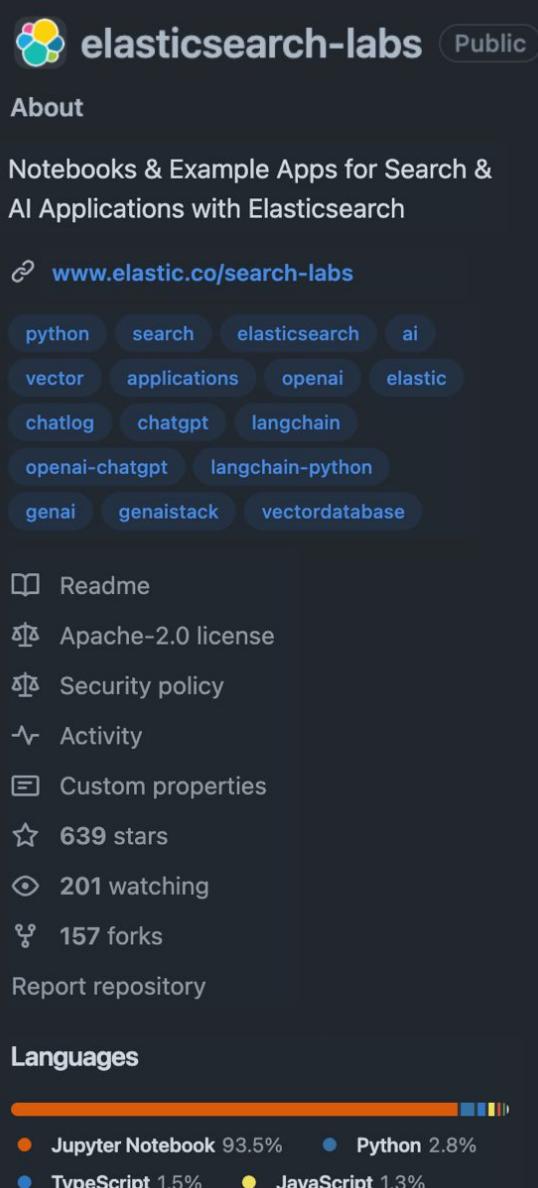
ML Research December 1, 2023

## RAG evaluation metrics: A journey through metrics

Explore RAG evaluation metrics like BLEU score, ROUGE score, PPL, BARTScore, and more. Discover...

By: Quentin Herreros, Thomas Veasey and Thanos Papaoikonomou

[elastic.co/search-labs](https://elastic.co/search-labs)



elasticsearch-labs Public

About

Notebooks & Example Apps for Search & AI Applications with Elasticsearch

[www.elastic.co/search-labs](https://www.elastic.co/search-labs)

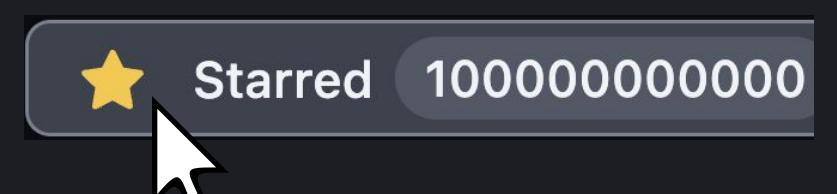
python search elasticsearch ai  
vector applications openai elastic  
chatlog chatgpt langchain  
openai-chatgpt langchain-python  
genai gennai vectordatabase

Readme  
Apache-2.0 license  
Security policy  
Activity  
Custom properties  
639 stars  
201 watching  
157 forks  
Report repository

Languages

Jupyter Notebook 93.5% Python 2.8%  
TypeScript 1.5% JavaScript 1.3%

[github.com/elastic/elasticsearch-labs](https://github.com/elastic/elasticsearch-labs)



Tutorials  
GitHub Example Code  
Integrations  
Blogs





ElasticON is back!



# ElasticON New York

Nov. 13, 2024  
Manhattan Center



<https://www.elastic.co/events/elasticon/new-york-city>



# Elastic Contributors



This community program is designed to recognize and reward the hard work of our awesome contributors. Join this friendly competition and earn points for:

Code contributions | Presentations | Video tutorials | Event organization  
Translations | Technical Q&A | Written content | Content validation

The top contributors will win cool prizes such as **Elastic trainings, certificate exams, cloud credits** and more. Check out the rules on our [website](#) and start submitting today!



[elastic.co/community/contributor](https://elastic.co/community/contributor)

The current cycle runs from February 1st 2024 to January 31st 2025

# Get involved

**The community would love to hear from you!**

If you have a cool use case to share at a meetup, please let us know! We would love to make that happen.

And if your company would like to host a meetup - we have user groups in other cities as well, please let us know.

Send an email to [meetups@elastic.co](mailto:meetups@elastic.co) or [ully@elastic.co](mailto:ully@elastic.co)

**Let me know if this talk inspires you to build  
anything. I'm @JessicaGarson on most  
platforms.**

If you have any questions on our **Discuss forums** and the community **Slack channel**.

# Thank you!