



DEV DAY

# Vector Search: Beginner to Pro

WiFi Password: loveyourdevelopers



# Hands-on session

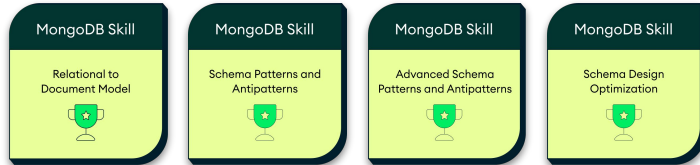
- There are no stupid questions
- Actively discuss
- Be respectful



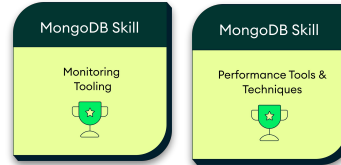
# MongoDB Skill Badges



## Data Modeling



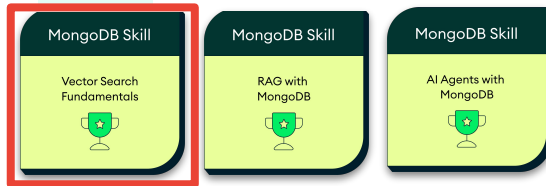
## Monitoring/Tuning/Automation



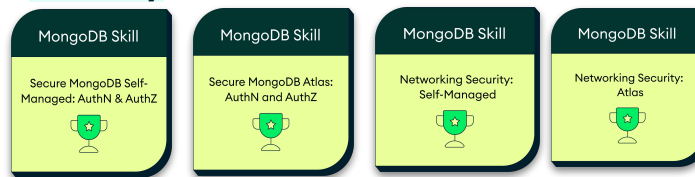
## Performance at Scale



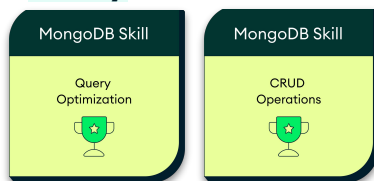
## Gen AI



## Security



## Query



## Aggregation



## Sharding



## Indexes



## Architecture



## Search



# Skill Badge Benefits

## **Formal Validation**

Provide a way to quickly learn and validate specific MongoDB skills

Elevate in current roles and increase appeal for future positions.

## **Career Advancement**

## **Digital Badges**

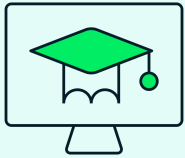
Showcase achievements with a Credly badge and inclusion in the Talent Directory.

MongoDB Skill



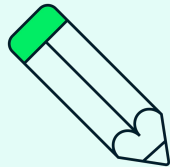


# Free, Focused Credentials



## Learn

Learn and practice during this workshop.



## Validate

Showcase your knowledge in a short 10-question skill check.



## Earn

Claim your Credly badge and share it on social media.

**Earn a  
badge in  
60 – 90  
minutes.**






# Agenda

- What is vector search?
- What are embeddings?
- Vector search in MongoDB
- Tuning vector search queries
- Q&A



# Lab instructions

 , 	Hands-on sections
<CODE_BLOCK_N>	Code placeholders
	Reference documentation for the code blocks in the cell that follows

# Lab setup

**PASSKEY** = “replace-me”

## TODO:

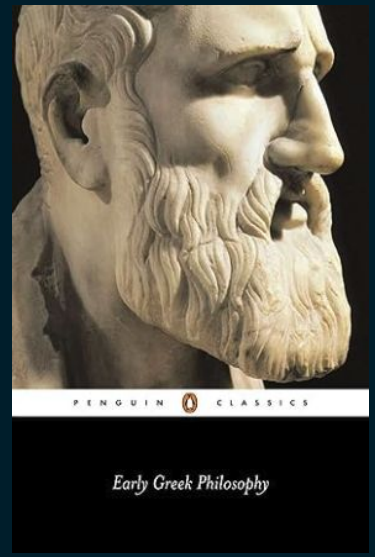
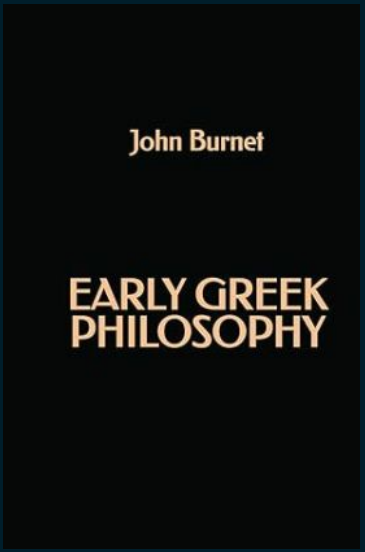
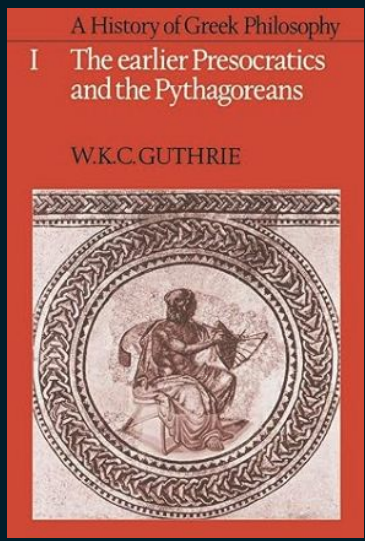
- Navigate to the lab:  
[mdb.link/devrel-vector-search](https://mdb.link/devrel-vector-search)
- Navigate to the lab setup:  
[mdb.link/lab-setup](https://mdb.link/lab-setup)
- Complete the 🙌 modules

A decorative graphic on the left side of the slide, consisting of a light purple rounded rectangular shape and a thin green line that curves from the bottom left towards the top right.

# Lexical vs. Vector Search

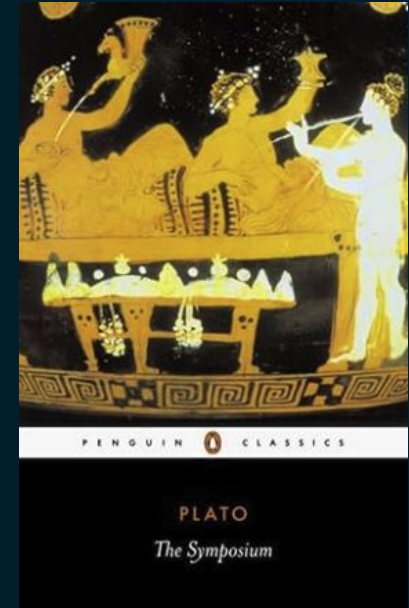
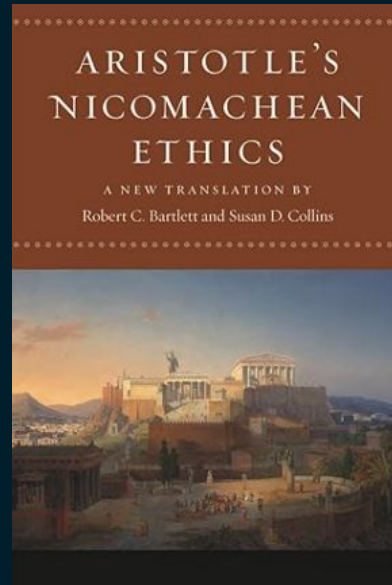
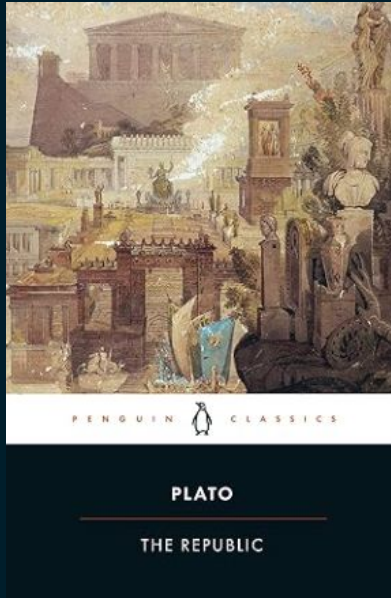


🔍 Greek philosophy



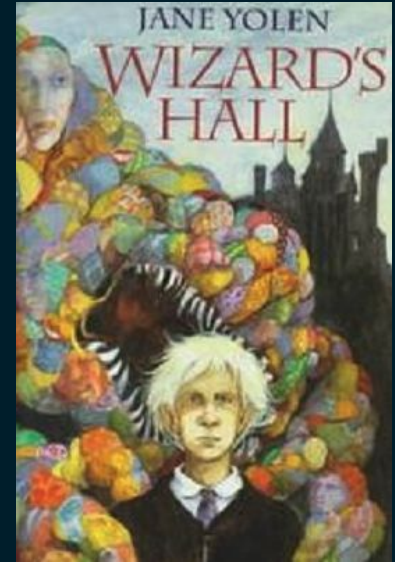
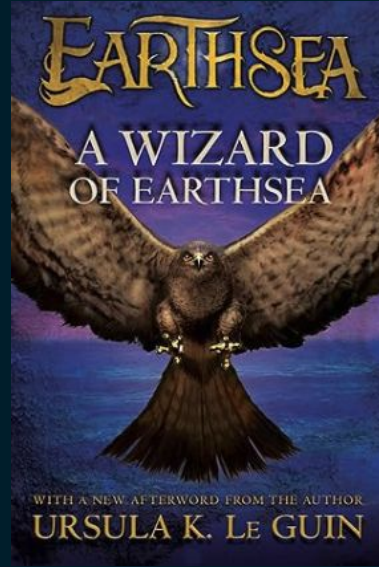
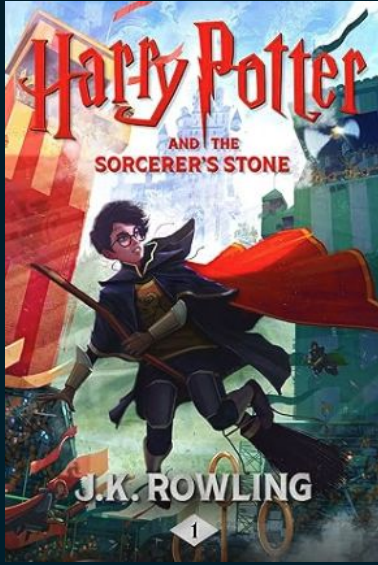


🔍 Greek philosophy

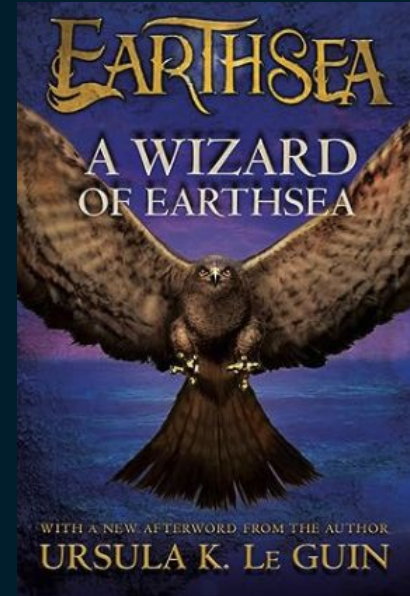
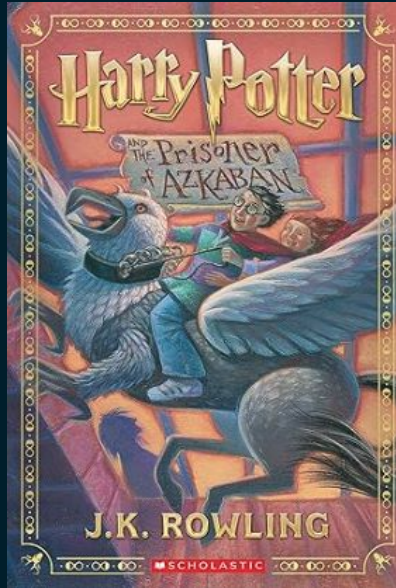
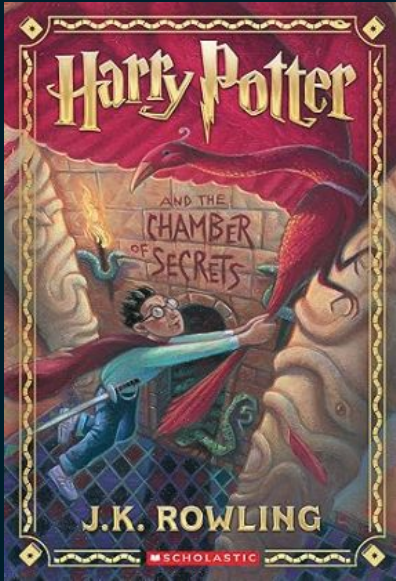




🔍 Fantasy novel about a boy in a school of magic



🔍 harry\_potter\_1\_book\_cover.png



# Lexical search

What?

- Keyword search

When?

- Your text corpus closely matches how users search
- First pass at text-based relevancy

# Vector search



What?

- Semantic similarities

When?

- “Vocabulary gap” between corpus and how users search
- Text, image, audio, video search

# Glossary

## Embedding

Representation into vectors of the semantic meaning produced by LLM

## Embedding Model

AI model that understands the human language meaning of large input

## Vector and Dimensions

A numerical array representing a point in a multidimensional space, often used as an angle to the axis, and a length

## Vector Indexing

An inverted index on vectors to find and score embedding with similarity





# Vectors



**Vectors** = numerical representations of unstructured data.

- encoded as an array
- each element represents a dimension
- generated using embedding models



# Vectors example: store



SKILL



[aisle, bay]





[aisle, bay]



[aisle, bay, brand, color]



Thousands of dimensions 🧠



A decorative graphic on the left side of the slide consists of a light purple rounded rectangular shape at the bottom, with a thin green line that curves upwards from its top edge and then extends vertically to the top of the slide.

# Embeddings



# Converting text to vectors

If you're new to machine learning, you might wonder how plain text gets transformed into vectors — mathematical objects that power semantic search.

Let us build some intuition with a simple example using character bigrams.

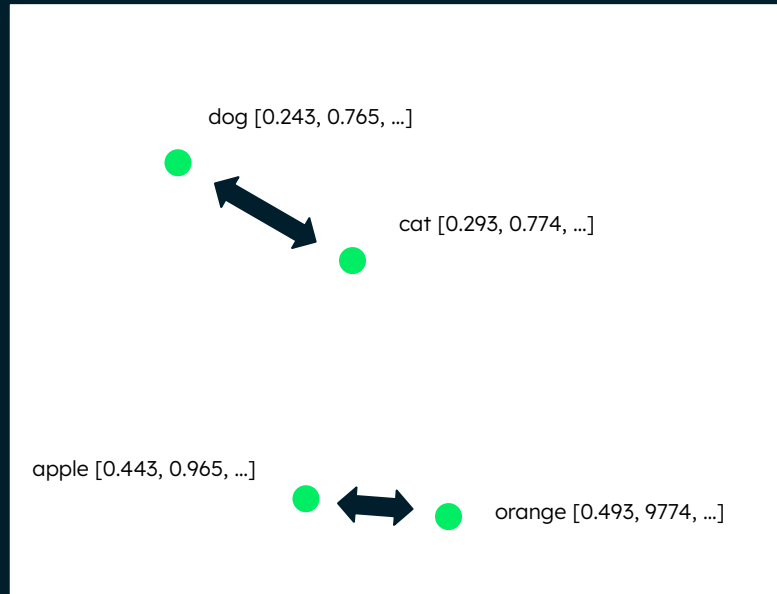
“Aircraft” -> "ai", "ir", "rc", "cr", "ra", "af", "ft"

How many possible bigrams of alphabets (Dimensions) -  $26*26= 676$

Vector representation of “Aircraft” : { "ai":1, "ir":1, "rc":1, "cr":1, "ra":1, "af":1, "ft":1, "aa": 0, "ab":0, ....}



# Embeddings

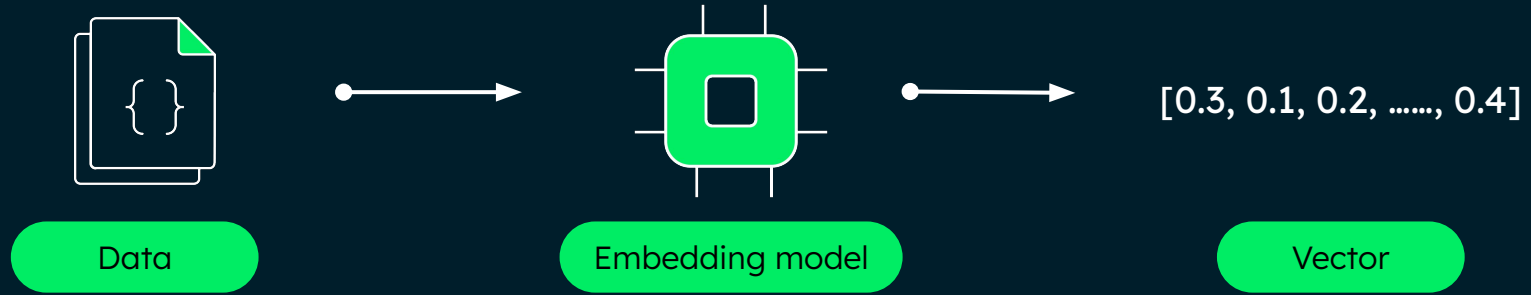


Numeric, multi-dimensional representation of a piece of information

- Capture semantic qualities of data
- Semantically similar data ends up close together in vector space



# How to embed data





Text

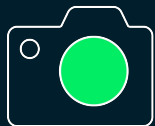


VOYAGE AI

Multimodal  
Embedding model

[0.3, 0.1, 0.2, ....., 0.4]

Vector



Image



VOYAGE AI

Multimodal  
Embedding model

[0.3, 0.1, 0.2, ....., 0.4]

Vector



Video



VOYAGE AI

Multimodal  
Embedding model

[0.3, 0.1, 0.2, ....., 0.4]

Vector

Project Overview ⚙️

## DATABASE ▲

Clusters

Search &amp; Vector Search

Data Explorer

Backup

## STREAMING DATA ▲

Stream Processing

Triggers

## SERVICES ▲

AI Models

Migration

Data Federation

Visualization

## Supercharging Search and Retrieval for Unstructured Data

## Voyage AI Embedding &amp; Reranking Models

Best-in-class models built to deliver faster, more accurate, and scalable retrieval for semantic search and RAG

[+ Create model API key](#)

## Getting Started with Embedding and Reranking API Service

1

## Create API key

Start by creating your API key. We recommend setting it as an environment variable.

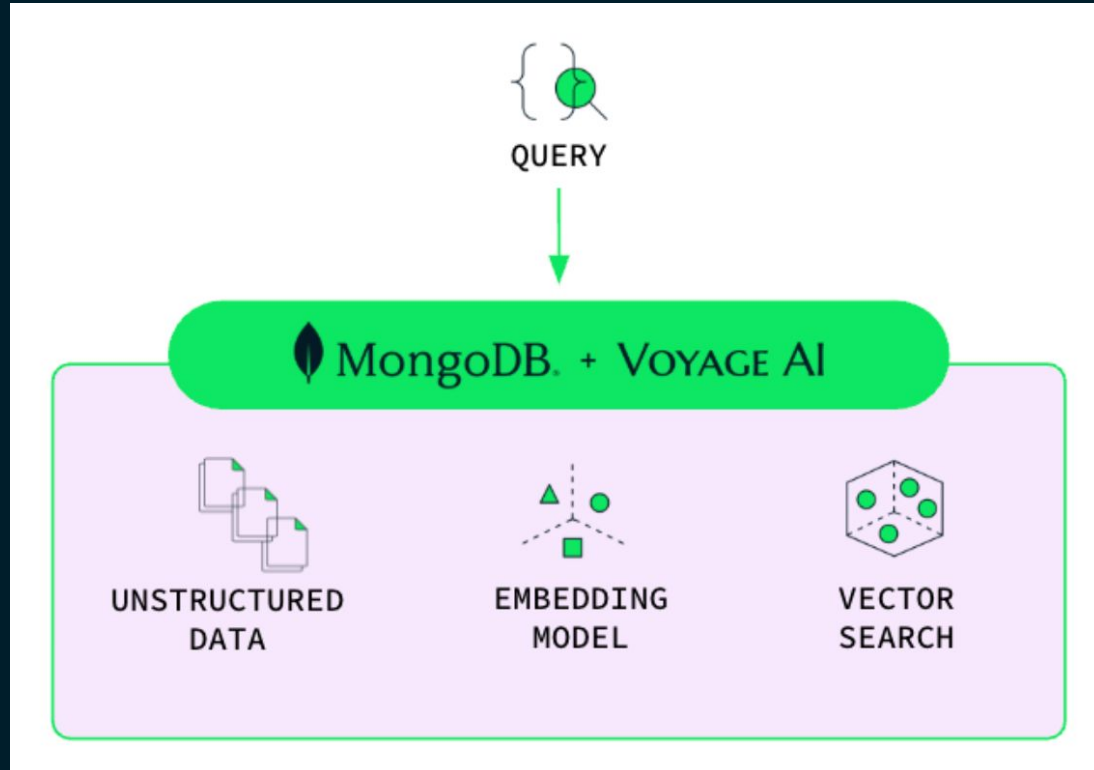
[Create API Key](#)

2

## Create your first embedding with Python

Install the Python Package and run a quick test to generate embeddings using our state-of-the-art models.

# Auto-embedding in MongoDB Atlas Vector Search





# How to choose an embedding model

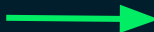


[mdb.link/embedding-models](https://mdb.link/embedding-models)



# Adding embeddings to existing data

```
_id: "0028608488"  
title: "David Copperfield's Tales of the  
Impossible"  
cover:  
"https://images.isbndb.com/covers/22/86/97  
80061052286.jpg"  
year: 1995  
pages: 385  
synopsis: "David Copperfield, Arguably The  
Greatest Illusionist-magician..."
```



```
_id: "0028608488"  
title: "David Copperfield's Tales of the  
Impossible"  
cover:  
"https://images.isbndb.com/covers/22/86/97800  
61052286.jpg"  
year: 1995  
pages: 385  
synopsis: "David Copperfield, Arguably The  
Greatest Illusionist-magician..."  
embedding: Array  
  0: 0.03898080065846443  
  1: -0.05879044905304909  
  2: 0.04323239979442215  
  .  
  .  
  512: 0.034234036451233547
```

# Recap



SKILL



- Embeddings are an array of numbers that capture semantic attributes of data.
- Embeddings are generated by specialized ML models.
- Different embedding models produce embeddings of different sizes.
- Embeddings can be added in-place into existing MongoDB documents.

# Lab setup

**PASSKEY** = “+8wWcA”

## TODO:

- Navigate to the lab:  
[mdb.link/devrel-vector-search](https://mdb.link/devrel-vector-search)
- Navigate to the lab setup:  
[mdb.link/lab-setup](https://mdb.link/lab-setup)
- Complete the 🙌 modules

# In this exercise, you will:

- Import a dataset into MongoDB.
- Add embeddings to your MongoDB data.

## **TODO:**

- Navigate to [mdb.link/devrel-vector-search](https://mdb.link/devrel-vector-search)
- Complete **Step 2-4** in the notebook

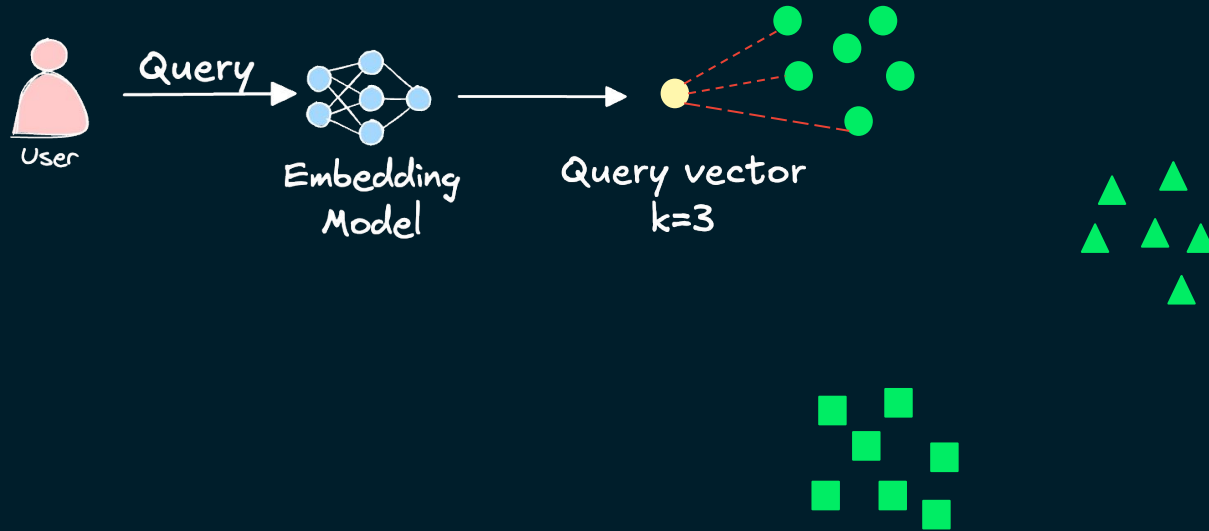
A decorative graphic on the left side of the slide consists of a light purple rounded rectangular shape at the bottom, with a thin green line that curves upwards from its top edge, then turns vertically to the top of the slide.

# Vector Search

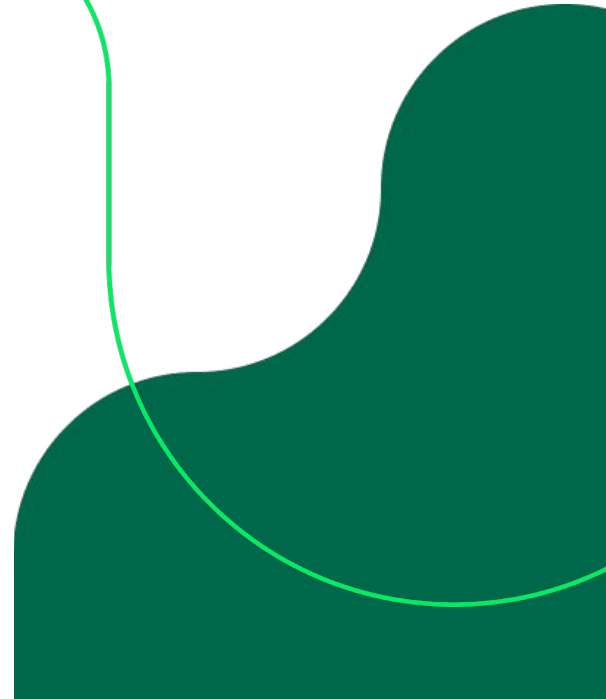


# Vector search

Search based on intent/meaning using embeddings



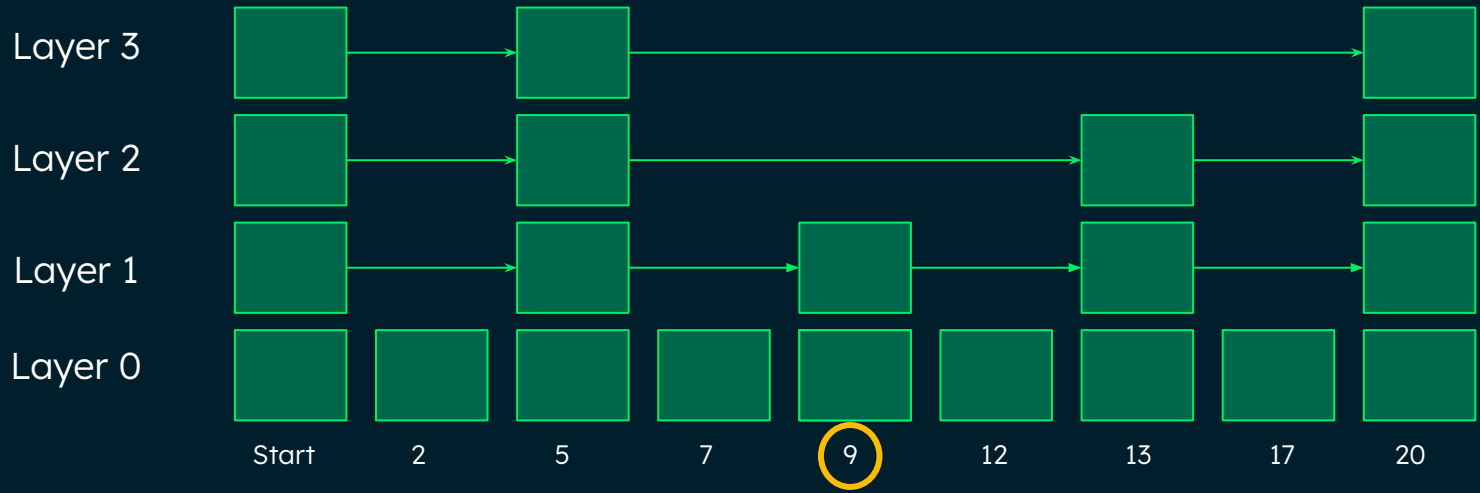
# Indexing Algorithms



# Skip List



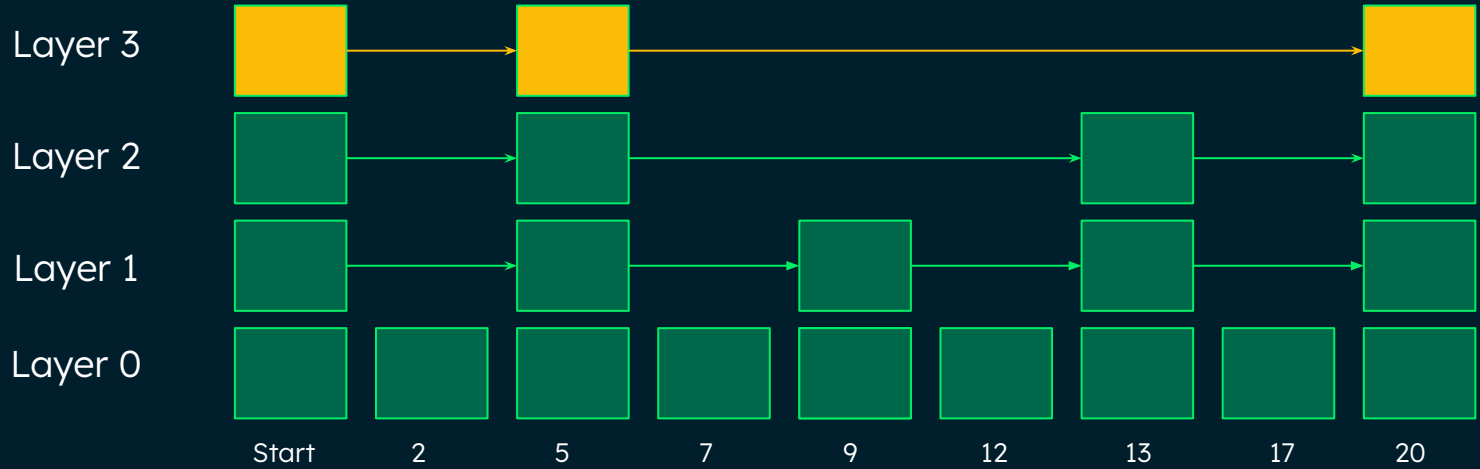
Find 9



# Skip List



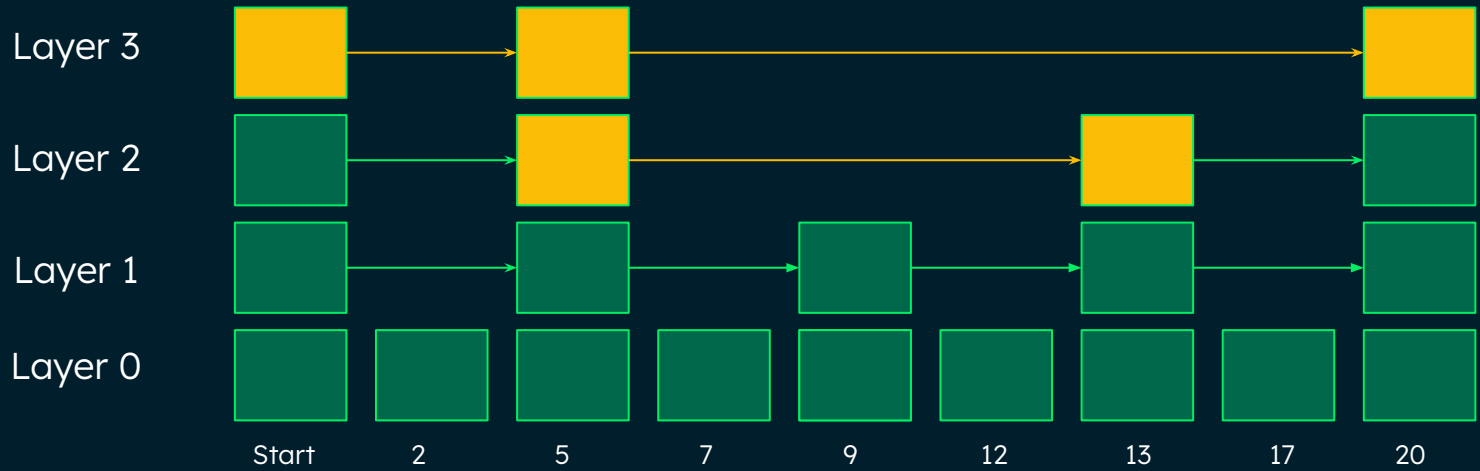
## Find 9



# Skip List



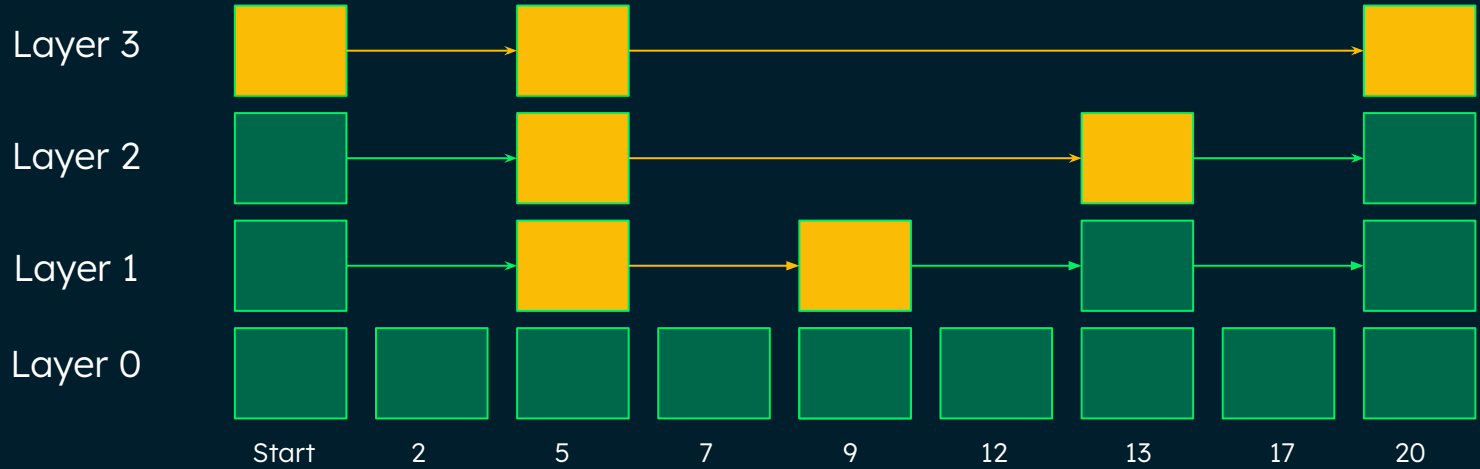
Find 9



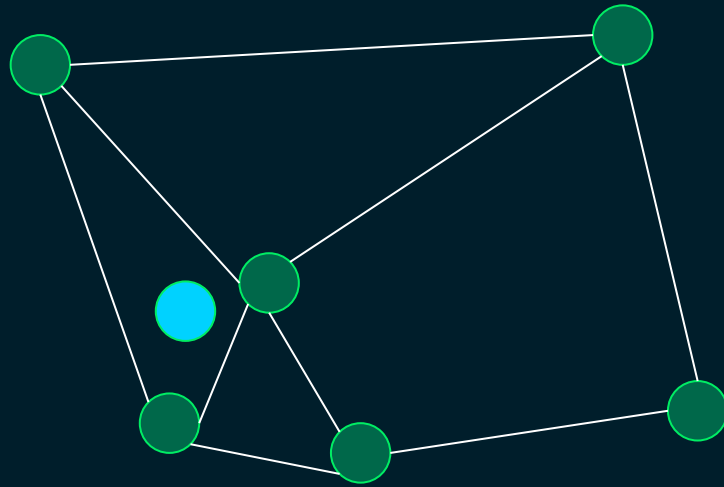
# Skip List



Find 9



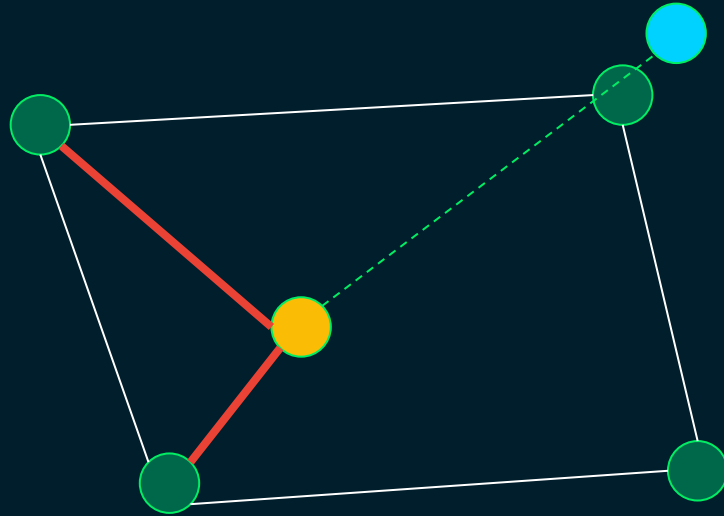
# Navigable Small World Graph



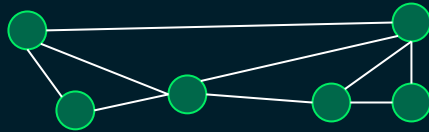




# Navigable Small World Graph



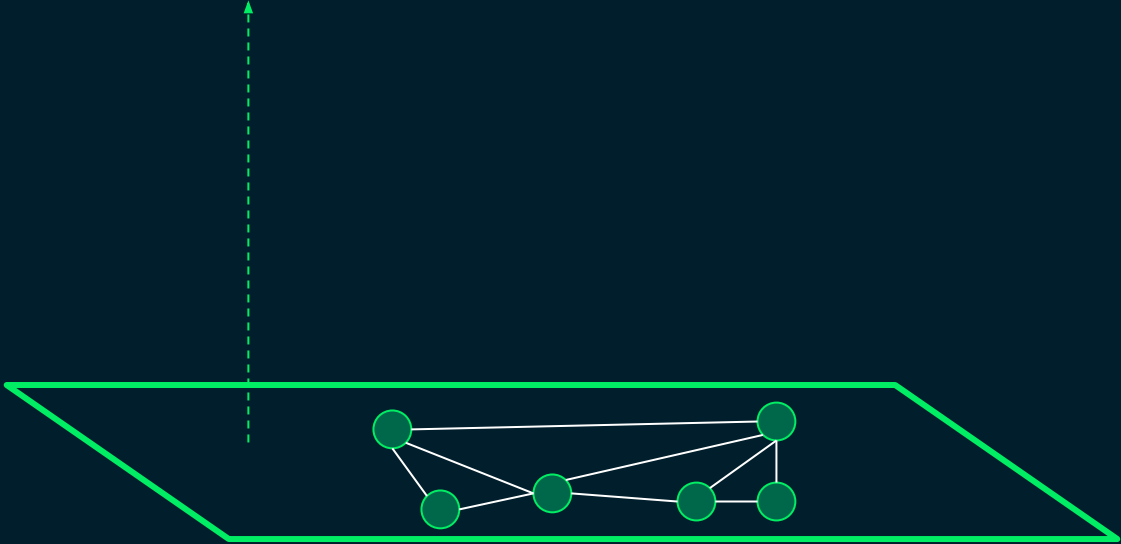
# Hierarchical Navigable Small World Graph



# Hierarchical Navigable Small World Graph



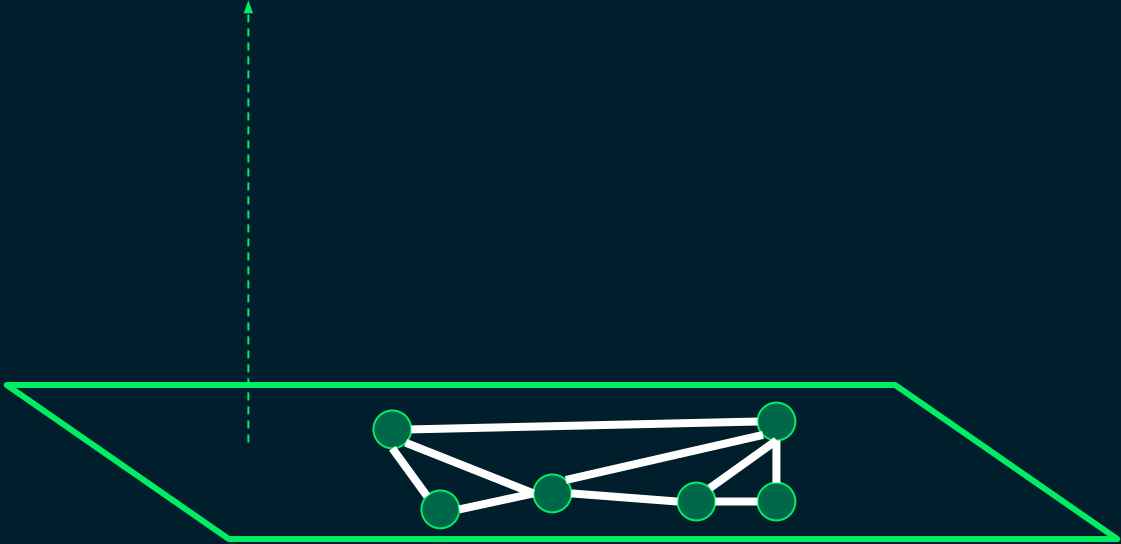
Layer 0



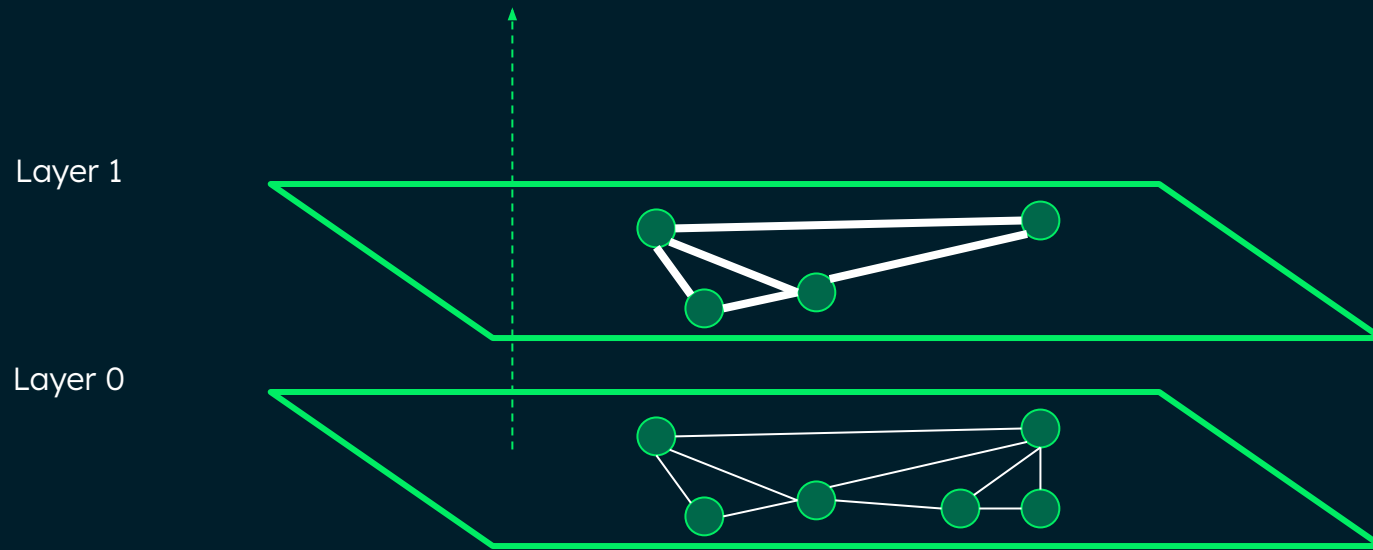
# Hierarchical Navigable Small World Graph



Layer 0



# Hierarchical Navigable Small World Graph



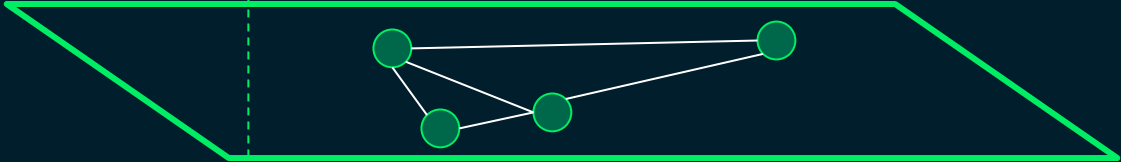
# Hierarchical Navigable Small World Graph



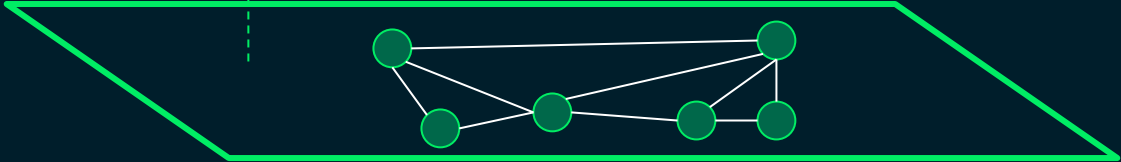
Layer 2



Layer 1



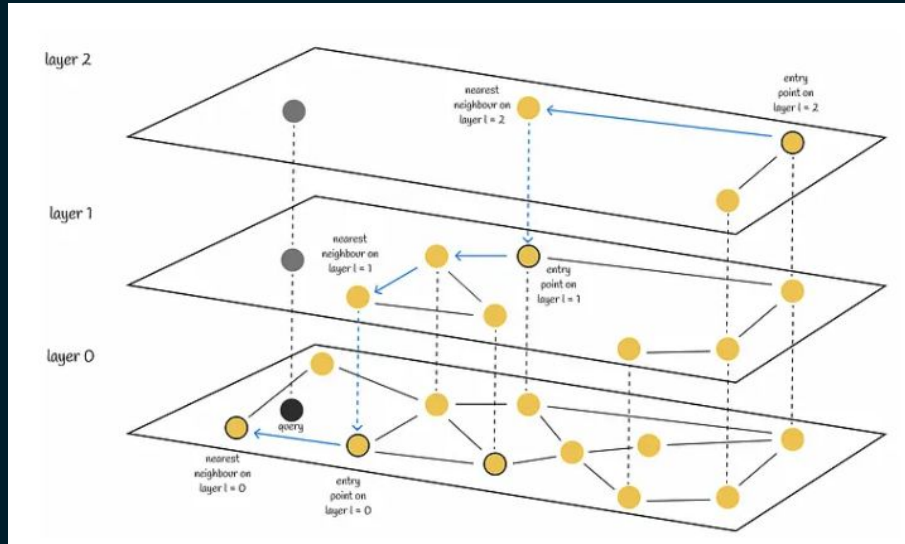
Layer 0





# How vector search works in MongoDB

SKILL



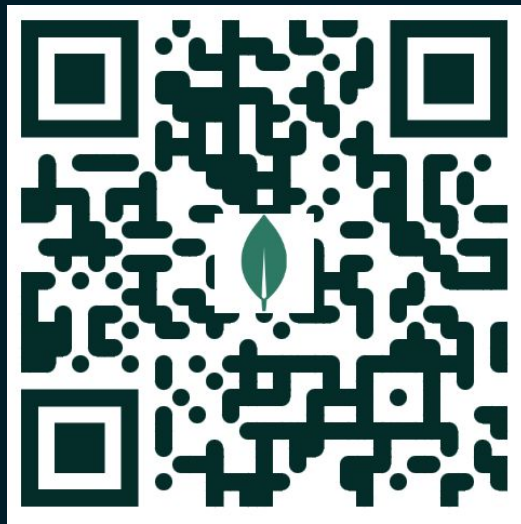
## Hierarchical Navigable Small Worlds

- Creates layered, connected graphs with vectors as nodes, edges created based on distance in vector space
- Coarse search at top layers, refinement at lower layers
- Efficient due to approximate nearest neighbor search (ANN)

Source: Towards Data Science



# HNSW deep-dive

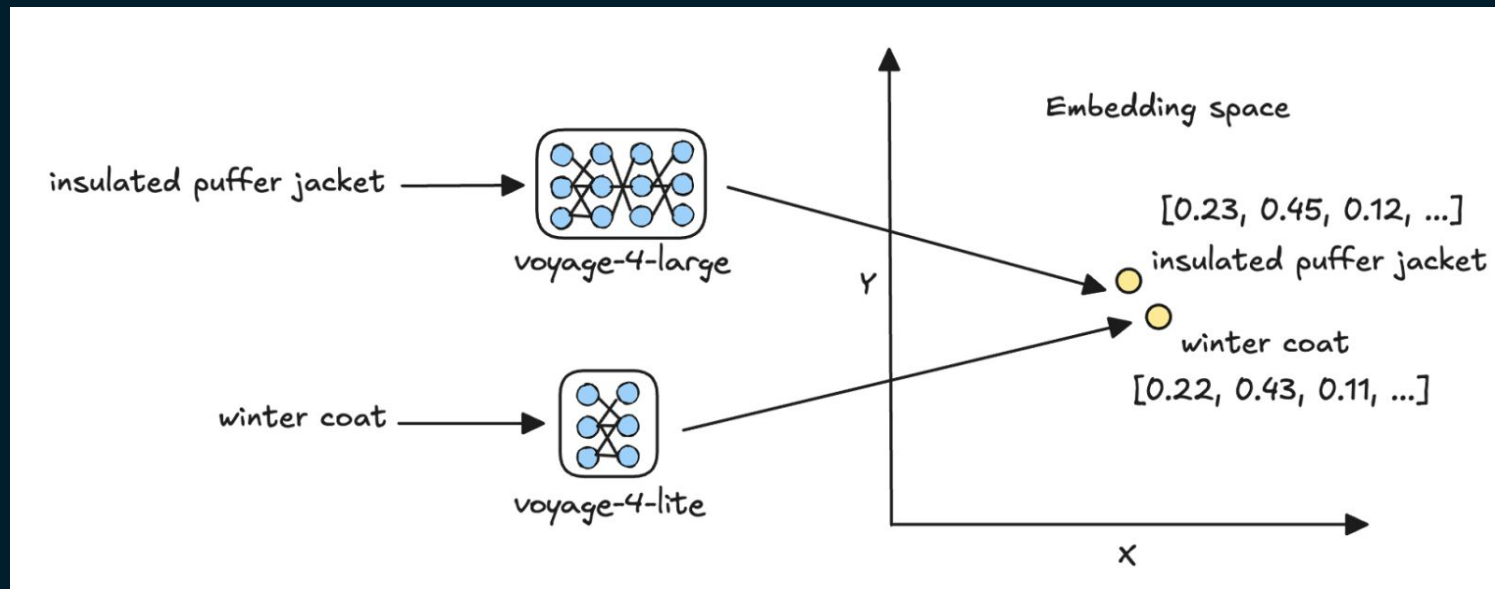


[mdb.link/hnsw-deepdive](https://mdb.link/hnsw-deepdive)

# Asymmetric retrieval using Voyage 4



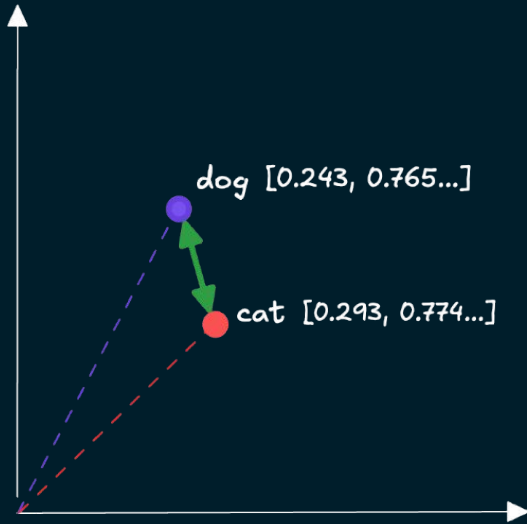
SKILL



Read more: [mdb.link/voyage-4](https://mdb.link/voyage-4)



# Calculating distance in vector space



## Euclidean Distance

Measures absolute distance between vectors



SKILL



# Calculating distance in vector space

$$\begin{array}{ccc} [4, & 0, & 1] \\ \updownarrow * & \updownarrow * & \updownarrow * \\ [3, & 1, & 2] \end{array}$$

## Dot product

Vector multiplication as a measure of alignment

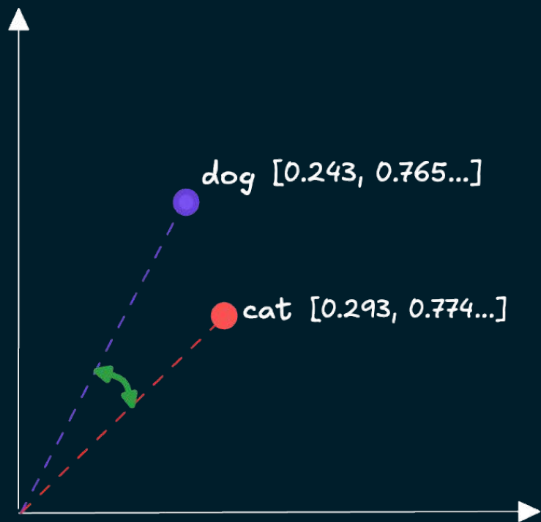
$$12 + 0 + 2 = 14$$



SKILL



# Calculating distance in vector space



## Cosine similarity

Measures the angle between vectors



SKILL




# Recap

- Vector search retrieves documents closest to the query embedding in vector space.
- Use compatible embedding models to embed the data and the user queries.
- Distance in vector space is calculated using mathematical functions.
- Cosine similarity works well with most embedding models.



# Vector search in MongoDB

1. Generate embeddings 
2. Create a vector search index
3. Send a vector search query



# Create a vector search index

## ERAS API

```
{
  "fields": [
    {
      "type": "vector",
      "path": "embedding",
      "numDimensions": 1024,
      "similarity": "cosine"
    },
    ...
  ]
}
```

## Auto-embedding

```
{
  "fields": [
    {
      "type": "autoEmbed",
      "modality": "text",
      "path": "synopsis",
      "model": "voyage-4"
    },
    ...
  ]
}
```



# Send a vector search query

## ERAS API

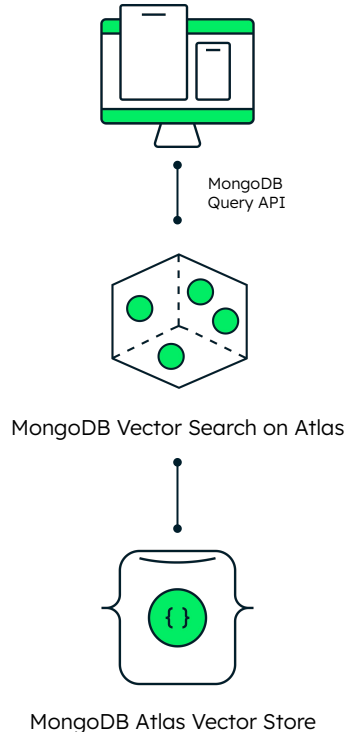
```
[
  {
    "$vectorSearch": {
      "index": "vector_index",
      "path": "embedding",
      "queryVector": [0.02421053, ...],
      "numCandidates": 200,
      "limit": 10
    }
  },
  {
    "$project": {
      "_id": 0,
      "title": 1,
      "score": {"$meta": "vectorSearchScore"}
    }
  }
]
```

## Auto-embedding

```
[
  {
    "$vectorSearch": {
      "index": "vector_index",
      "path": "synopsis",
      "query": {"text": "A boy in a..."},
      "model": "voyage-4",
      "numCandidates": 200,
      "limit": 10
    }
  },
  {
    "$project": {
      "_id": 0,
      "title": 1,
      "score": {"$meta": "vectorSearchScore"}
    }
  }
]
```



# MongoDB Vector Search on Atlas



**Integrated platform** that simplifies your application architecture

- Data is **automatically synchronized** between the database and vector index
- Developers work with database and vector search via the **unified MongoDB Query API**
- **Fully managed** for you so you can focus on your application
- Search nodes **scale your search workloads** independent of the operational database

Vector search simplified

Avoid the tax synchronization

Remove operational heavy lifting



SKILL



# Recap

- To perform vector search in MongoDB, you need to generate embeddings, create a vector search index and send a query.
- The number of dimensions in the embedding vector depends on the embedding model used.
- Vector search should always be the first stage in a vector search aggregation pipeline.

# In this exercise, you will:

- Create a vector search index.
- Perform vector search queries.

## **TODO:**

- Navigate to [mdb.link/devrel-vector-search](https://mdb.link/devrel-vector-search)
- Complete **Step 5-6** in the notebook



# Optimizing Your Vector Search Queries



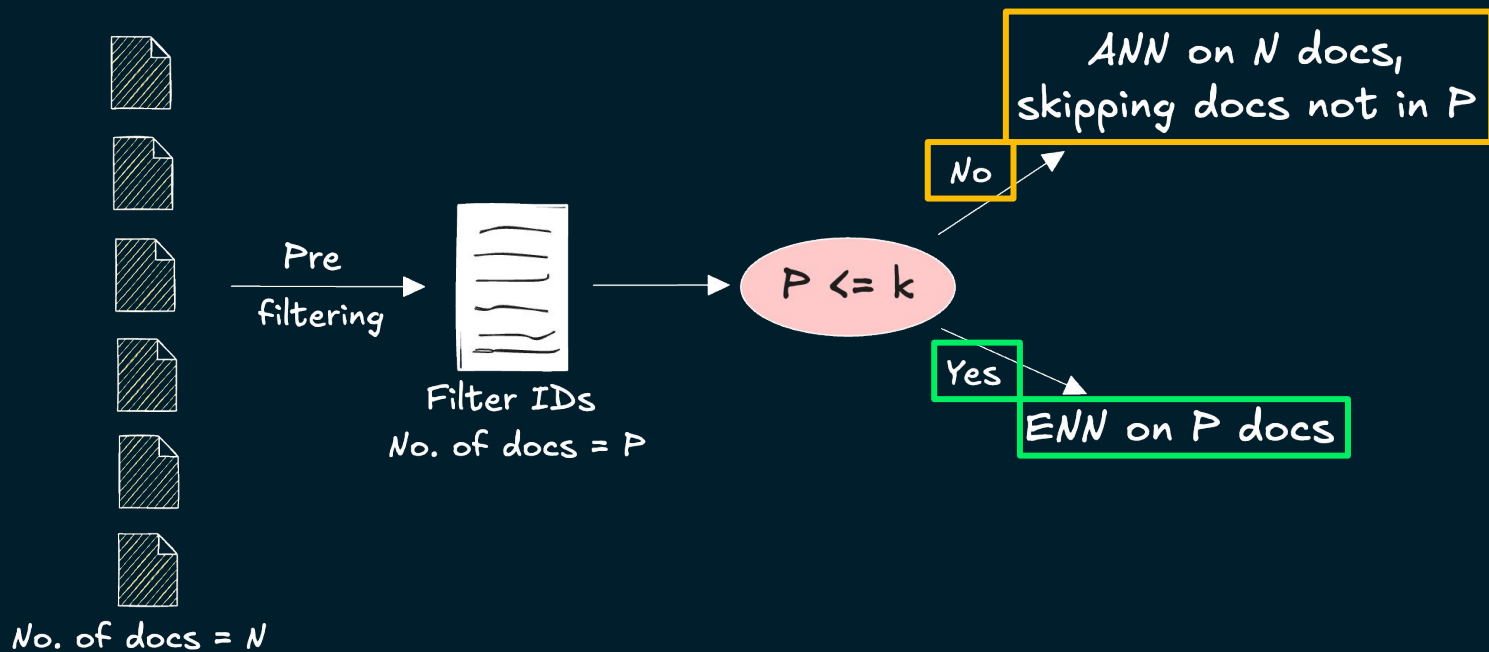
# Pre-filtering

# Pre-filtering

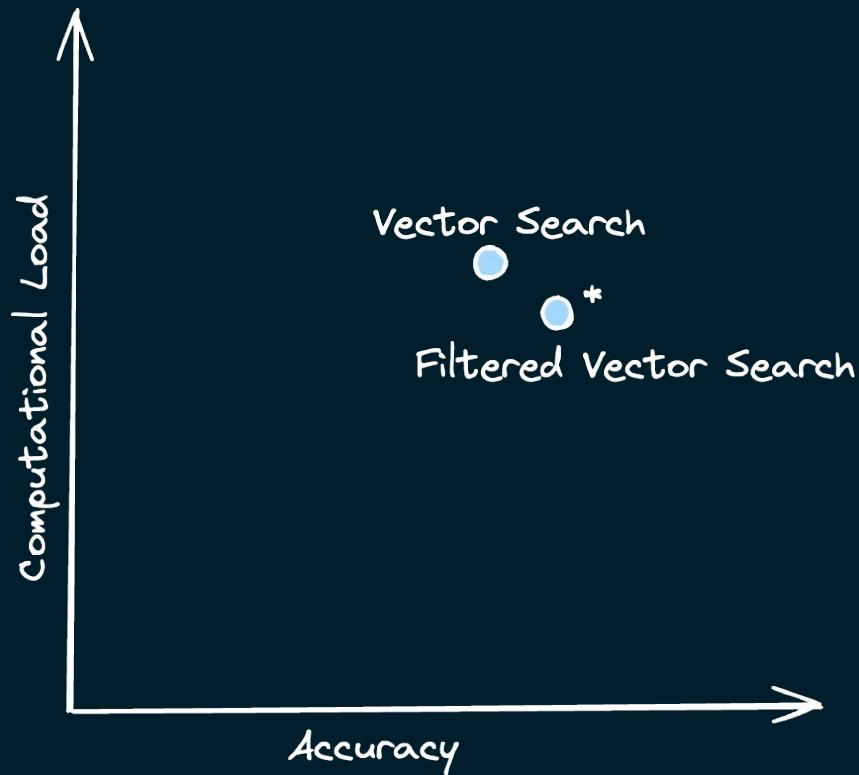
SKILL



$k$  = No. of results to return



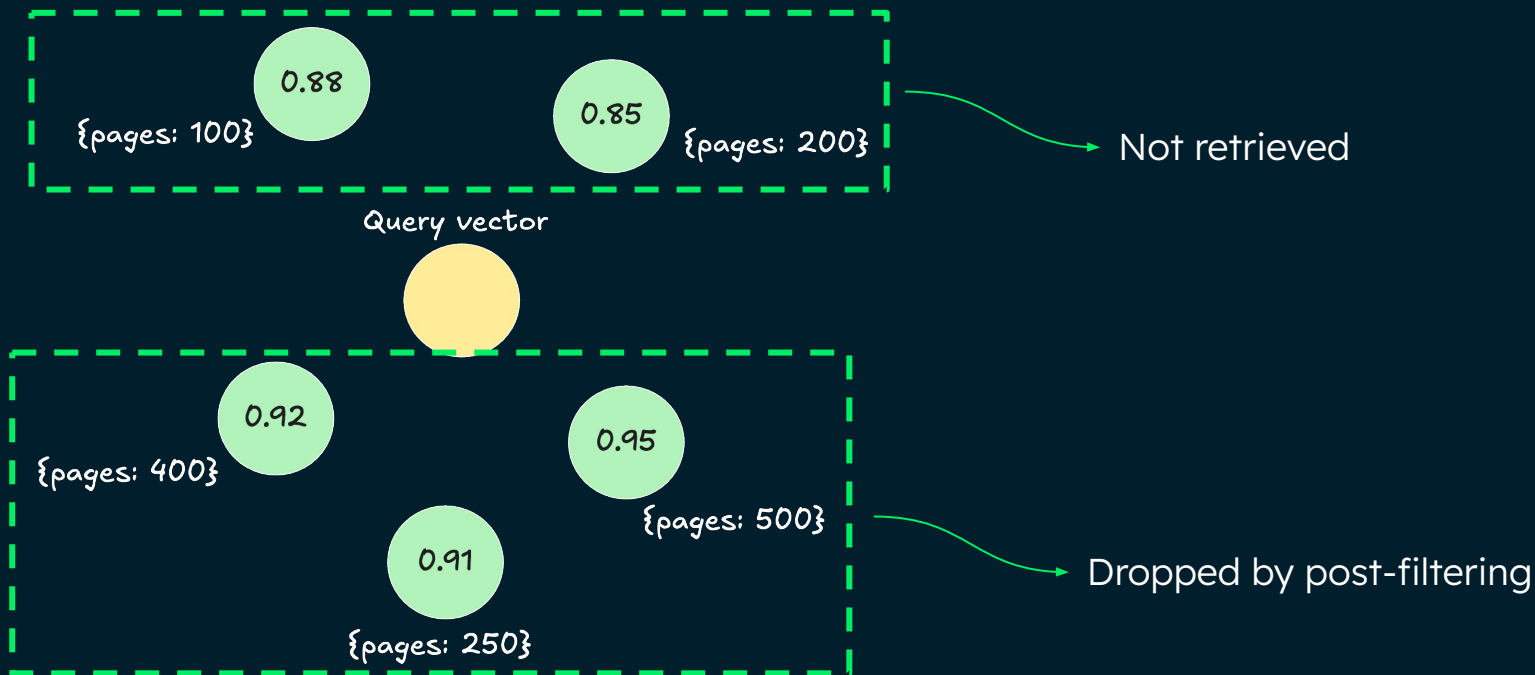
# Pre-filtering





# Pre-filtering >>> post-filtering

**filter** : pages <= 200    **k** : 3





# Adding pre-filters to vector search

```
{
  "fields": [
    {
      "type": "vector",
      "path": "embedding",
      "numDimensions": 1024,
      "similarity": "cosine"
    },
    {
      "type": "filter",
      "path": "pages"
    },
    ...
  ]
}
```



# Adding pre-filters to vector search

```
pipeline = [  
  {  
    "$vectorSearch": {  
      "index": "vector_index",  
      "path": "embedding",  
      "filter": {"pages": { "$lte": 200 } },  
      "queryVector": [0.02421053, -0.022372592, ...],  
      "numCandidates": 20,  
      "limit": 10  
    }  
  },  
  {  
    "$project": {  
      "_id": 0,  
      "Content": 1,  
      "score": {"$meta": "vectorSearchScore"}  
    }  
  }  
]
```

# Recap



SKILL



- Use pre-filters to filter vector search results based on certain criteria.
- Pre-filtering can slow down the vector search if the filters are too restrictive.
- Post-filtering is generally not recommended.

## In this exercise, you will:

- Combine pre-filtering with vector search.


### **TODO:**

- Navigate to [mdb.link/devrel-vector-search](https://mdb.link/devrel-vector-search)
- Complete **Step 7** in the notebook

A decorative green line starts from the top left, goes down, then curves right and then down again, ending near the text. A light purple rounded shape is on the left side, partially overlapping the green line.

# Vector Quantization



 Train ▾

 Deploy ▾

 Use this model ▾

Downloads last month  
**4,924,653**



 Safetensors ⓘ

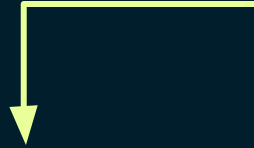
Model size | 33.4M params

Tensor type | 164 · FP16





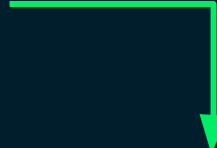
# FP32



## Data type

FP/float: Floating point

int: Integer



## No. of bits required to store each number

32 bits = 4 bytes



**$20 * 1024 * 4 \cong 81920$  bytes  $\cong$  82 KB**

What if you had 100 million embeddings?

**$100000000 * 1024 * 4 \cong 410$  GB**



# Vector quantization

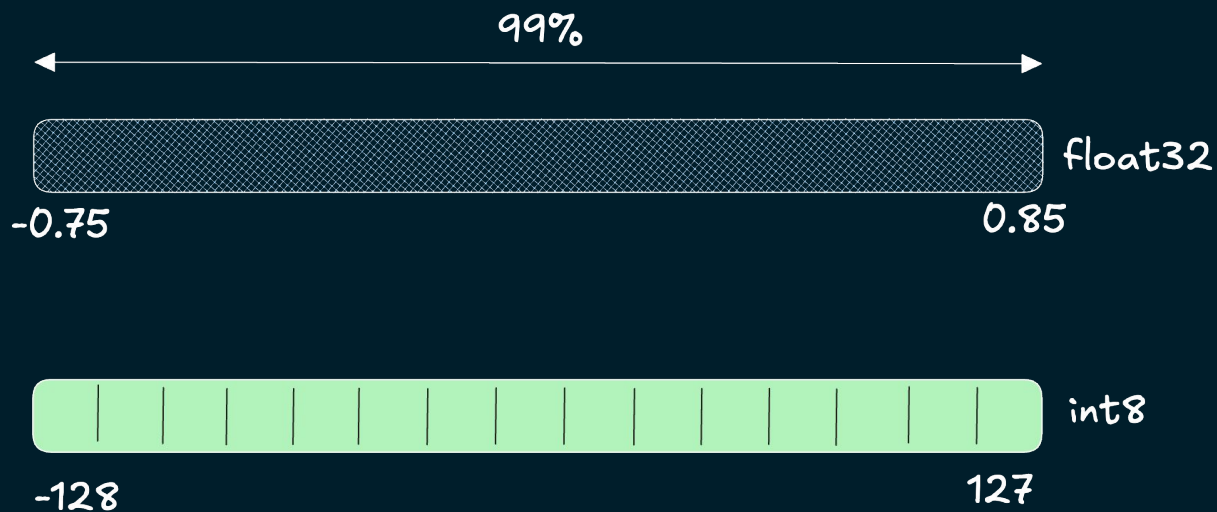
Quantization is the process of shrinking full-fidelity vectors into fewer bits

- **Scalar**
- **Binary**



# Scalar quantization

Takes each vector dimension and buckets it into a smaller set of discrete integers





# Binary quantization

Sets each vector dimension to a binary value based on a threshold

$$q(x) \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$



Data size: 3M

<b>Test</b>	<b>Memory and index savings</b>	<b>% of default performance</b>
Full-fidelity	1x	100%
Scalar quantization	3.75x	97%
Binary quantization	24x	95%

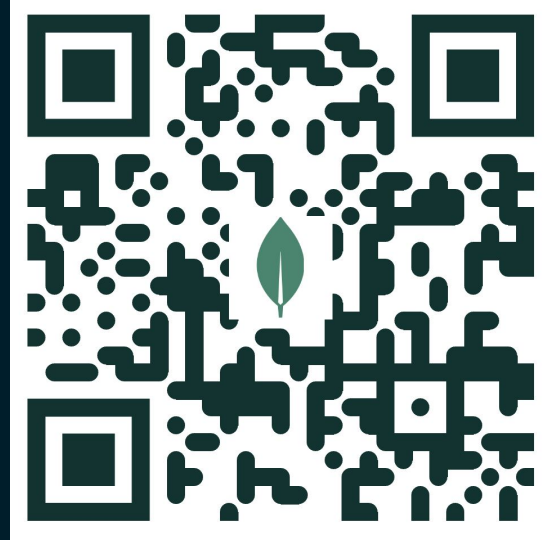


# Auto-quantization in MongoDB

```
{
  "fields": [
    {
      "type": "vector",
      "path": "embedding",
      "numDimensions": 1024,
      "similarity": "cosine",
      "quantization": "scalar"
    },
    ...
  ]
}
```



# Read more about quantization



[mdb.link/quantization](https://mdb.link/quantization)

# Recap



SKILL



- Vector quantization is the process of shrinking full-fidelity vectors into fewer bits.
- Quantization reduces retrieval latency and storage costs for a small dip in accuracy.
- Scalar quantization buckets each vector dimension into a smaller set of discrete integers.
- Binary quantization sets each vector dimension to a 0 or 1.

# In this exercise, you will:

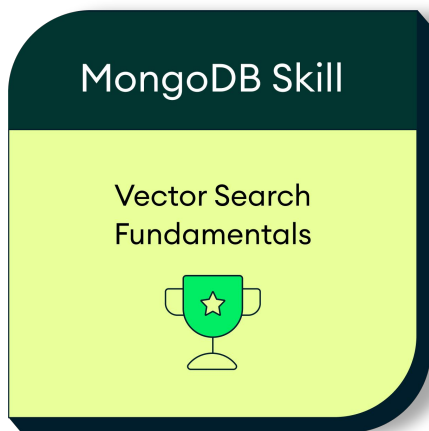
- Enable auto-quantization in MongoDB.

## **TODO:**

- Navigate to [mongodb.link/devrel-vector-search](https://mongodb.link/devrel-vector-search)
- Complete **Step 8** in the notebook

## SKILL CHECK

# Vector Search Fundamentals



- Understand Semantic Search
- Store Embeddings for Your Data
- Perform a Vector Search



<https://mdb.link/devday-chicago-26-vector-search>

## Questions?



Ask your neighbor



Raise your hand

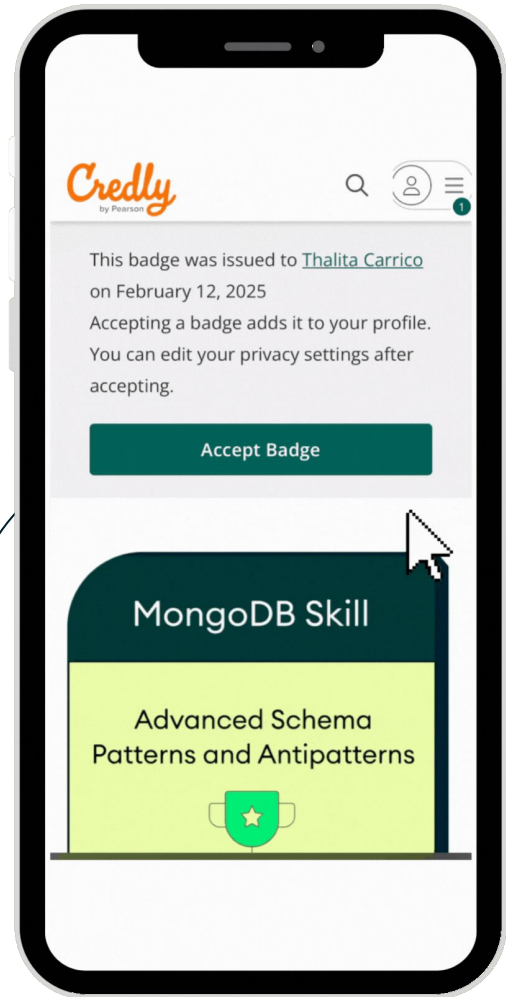


Let's do it together!



# Share Your Badge

Use **#MDBDAYCHICAGO**

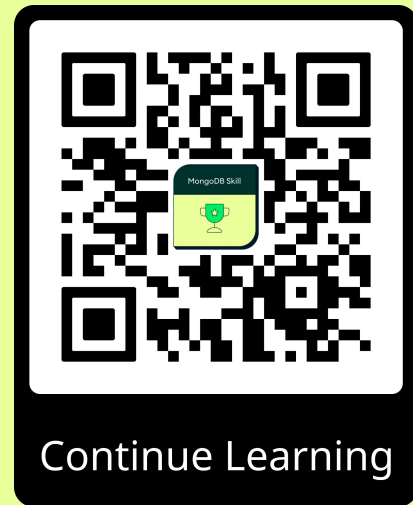


# Keep Learning

Every badge you earn elevates your skillset, empowering you and your team to accelerate and optimize application development and operations.

## Explore topics like:

- Data modeling
- Generative AI
- Aggregation
- Querying
- Security
- Indexing
- Sharding
- Monitoring, Tuning, and Automation
- Search
- Performance at Scale





# Check out our GenAI Examples repo!



[mdb.link/genai-examples](https://mdb.link/genai-examples)



# Check out our AI Learning Hub!



[mdb.link/ai-learning-hub](https://mdb.link/ai-learning-hub)



# Q & A

A decorative graphic on the left side of the slide consists of a light purple rounded rectangular shape at the bottom, with a thin green line that curves upwards from its top edge, then turns 90 degrees to run vertically along the left margin.

**Thank You!**