# Agenda

1. RAG Concepts
2. Why MongoDB
3. Models in Google
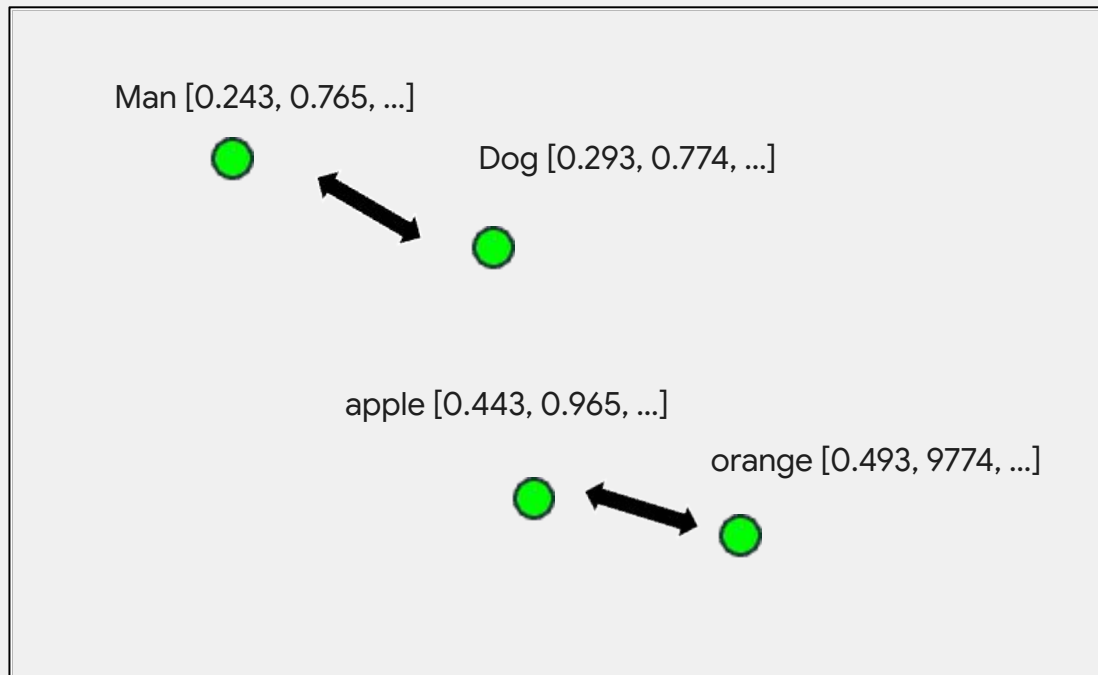4. Why Google Models
5. Workshop

# **Vectors** in ML are **lists of numbers** that represent attributes.

[0.743, 0.720, −0.325, 0.195, 0.835, −0.945, …]
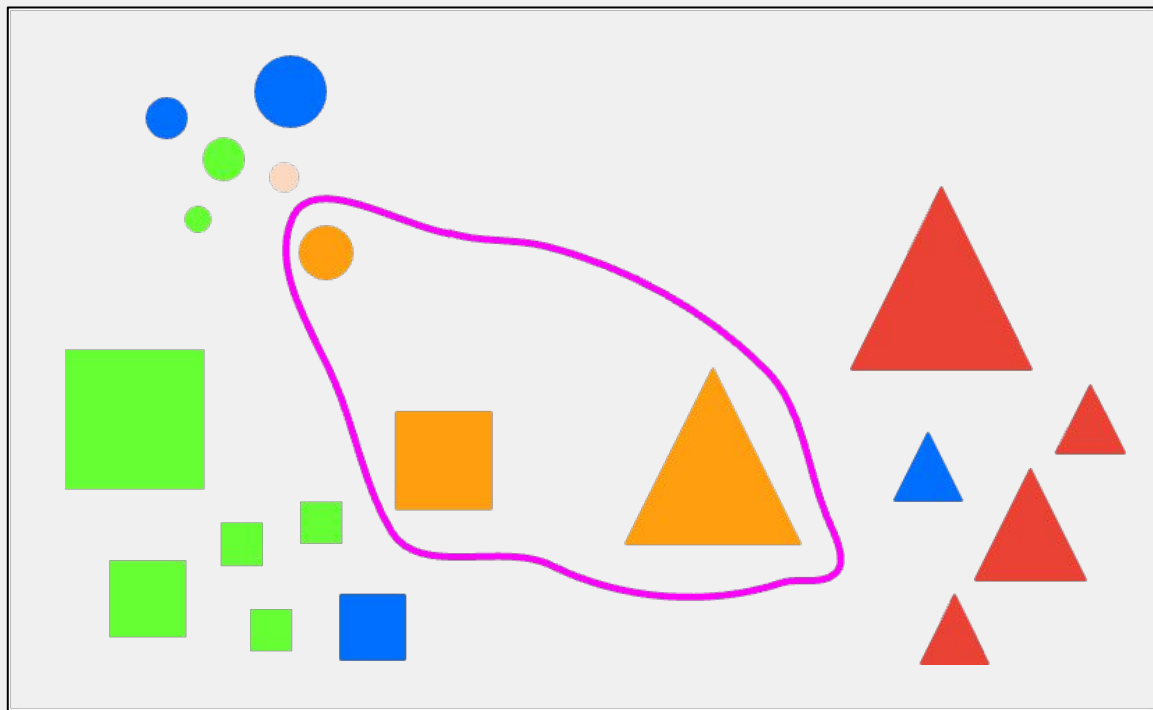
Each number represents a feature
(or a property of a data object)

# Similar vectors plotted in space will be near one another



Man [0.243, 0.765, …]

Dog [0.293, 0.774, …]

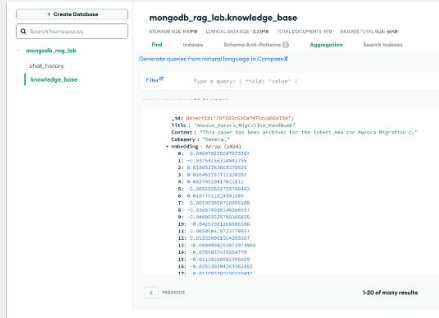apple [0.443, 0.965, …]

orange [0.493, 9774, …]

We can then take **search queries** and use **algorithms** to find **clusters** in high-dimensional space

# Why RAG in **MongoDB**

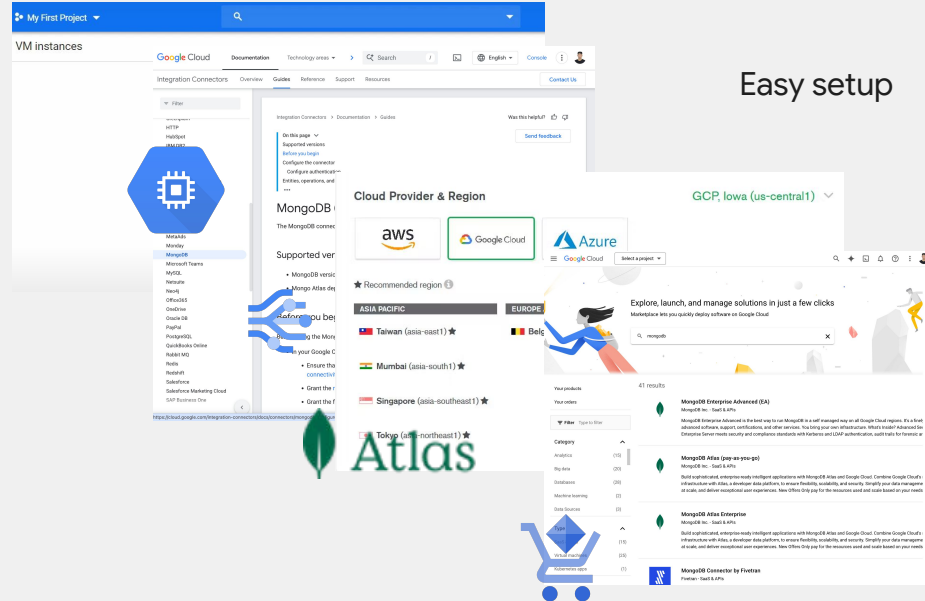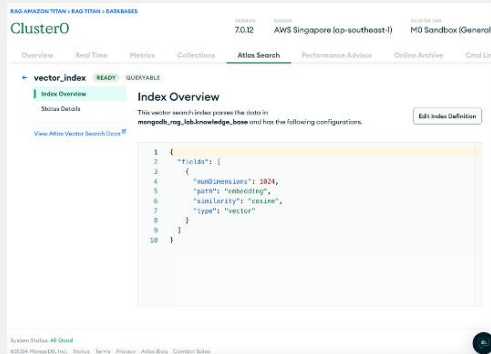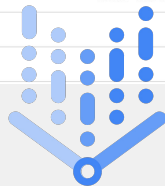## Embedding as Attributes



## Vector Search Index



## Easy setup

## Chunking

Vertex AI Agent Builder

## Embedding

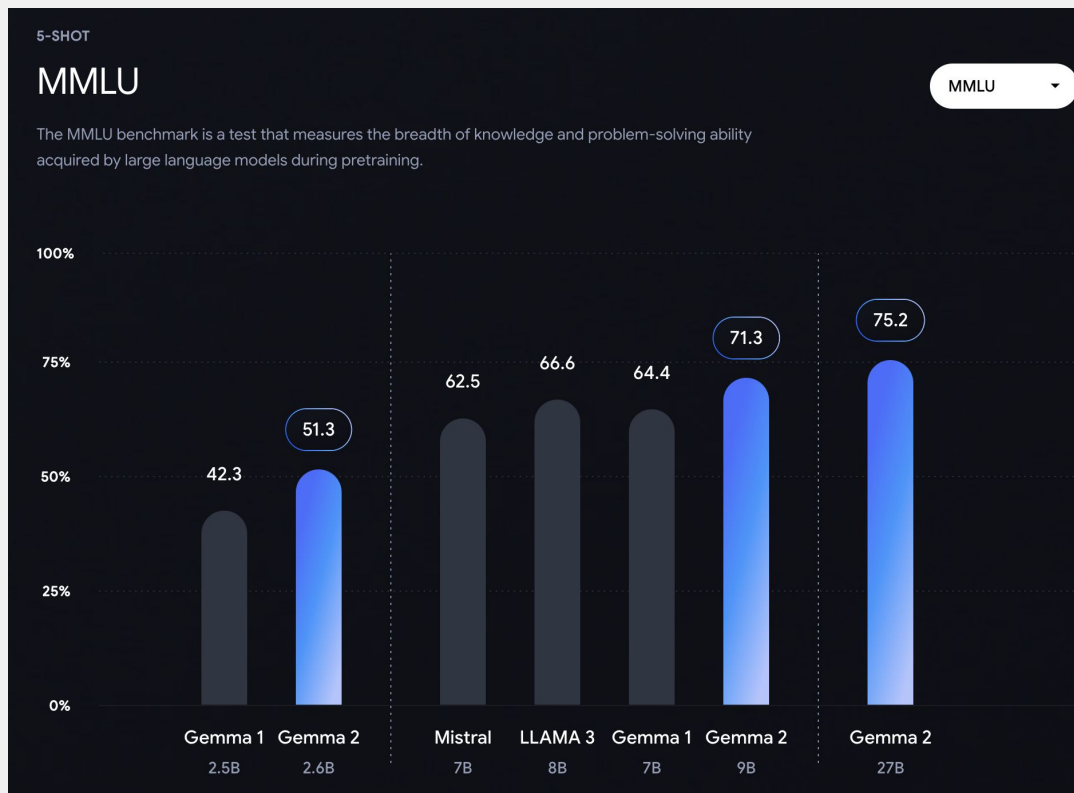| English models | Multilingual models |
|---|---|
| textembedding-gecko@001 | textembedding-gecko-multilingual@001 |
| textembedding-gecko@003 | text-multilingual-embedding-002 |
| text-embedding-004 | |
| text-embedding-005 | |

Vertex AI API

## Generation

| Model | Inputs | Outputs | Use case | Try the model |
|---|---|---|---|---|
| Gemini 1.5 Flash | Text, code, images, audio, video, video with audio, PDF | Text | Provides speed and efficien... volume... effectiv... | Try the Gemini 1.5 Flash model |
| Gemini 1.5 Pro | Text, code, images, audio, video, video with audio, PDF | Text | Suppor... prompt... code re... Suppor... unders... the ma... token li... | |
| Gemini 1.0 Pro | Text | Text | The be... model... of text... | |
| Gemini 1.0 Pro Vision | Text, images, audio, video, video with audio, PDF | Text | The be... image... unders... to hand... range o... | |

Gemini / Gemma

# Why **Gemma**



**5-SHOT**

## MMLU

The MMLU benchmark is a test that measures the breadth of knowledge and problem-solving ability acquired by large language models during pretraining.

MMLU ▾

| Model | | Score |
| --- | --- | --- |
| Gemma 1 2.5B | | 42.3 |
| Gemma 2 2.6B | | 51.3 |
| Mistral 7B | | 62.5 |
| LLAMA 3 8B | | 66.6 |
| Gemma 1 7B | | 64.4 |
| Gemma 2 9B | | 71.3 |
| Gemma 2 27B | | 75.2 |

| | |
| --- | --- |
| Context Length | 8192 |
| Strength | - Sliding window attention (low memory low time)<br>- Knowledge Distillation (Learning from sensei model)<br>- Model Merging |
| General Performance | 75.2 % MMLU Benchmark |

# **RAG** Flow