

# VECTOR DATABASES

---

## What this lesson is about:

You'll learn what vector databases are, why they're critical for agents, and how they support long-term memory and semantic search in LLM-powered systems like GPT.

## What is a Vector Database?

- A special kind of database made for AI and LLMs
- Stores **vectors** (numeric representations of meaning) instead of raw text
- Finds info based on **context and similarity**, not exact keywords
- Example: Like searching for places similar to your favorite coffee shop, not its exact name

## What are Vectors and Vector Stores?

- A **vector** = a group of numbers that represent meaning
- A **vector store** = a searchable memory where these vectors are saved
- Helps AI compare meanings and find relevant info even with different words

## Why Agents Need Them:

- LLMs need context and memory to generate better responses
- Vectors let agents retrieve **relevant chunks** of information
- Enables **semantic search**: finding data by meaning instead of exact match
- Used in RAG (Retrieval-Augmented Generation) agents and smart assistants

## How the Vectorization Process Works:

1. Start with your data: PDFs, text, JSON
2. **Chunk** the content into small pieces
3. **Embed** each chunk into a vector (using an LLM)
4. **Store** the vectors in a vector database
5. When a user asks something, convert it into a vector
6. Compare it to stored vectors → find best matches
7. Send relevant chunks to GPT to generate a response

# VECTOR DATABASES

---

## Popular Vector Databases Compared:

### 1. Pinecone

- Fully managed and built for semantic search
- Easy to use, plug-and-play setup
- Great for agents, chatbots, recommendations
- Best for: Fast, real-time agents
- We'll use Pinecone later in the course

### 2. Supabase + PGVector

- Not a vector DB by default, but supports vectors via extension
- Combines structured data (like user profiles) with vector search
- Open source, can be self-hosted
- Best for: Hybrid use cases with both SQL and vector needs
- We'll use Supabase in this course

### 3. Qdrant

- Open source, self-hosted vector DB
- Fast and powerful for mid-to-large datasets
- Works well with LangChain and other frameworks
- Best for: Full control, custom agent systems
- Requires more setup than Pinecone

## Summary: Which to Choose?

- Pinecone – Best for fast, managed semantic search
- Supabase – Best for structured + vector data together
- Qdrant – Best for self-hosted, cost-efficient solutions