

AI Agents Masterclass: RAG, Embedding & Vectors

Deep dive into building smart AI agents with your personal data using
Retrieval Augmented Generation

Created by Chinmay Kaitade
MERN Stack Developer | AI Enthusiast



The Data Library Problem

Imagine your personal data as a vast library filled with notes, files, and PDFs. Your chatbot struggles to find specific information within this overwhelming collection.

Traditional search methods fall short when dealing with contextual meaning and semantic understanding.

Created by Chinmay Kaitade

MERN Stack Developer | AI Enthusiast



What is Embedding?

01

Data Input

Your documents, notes, and files are processed

02

Vector Conversion

Text is converted into meaningful numerical codes (vectors)

03

Semantic Understanding

Similar meanings receive similar vector representations

04

Query Matching

Your questions are matched against relevant data vectors

Created by Chinmay Kaitade

MERN Stack Developer | AI Enthusiast



The Magic of Vector Representation

Embedding = Meaningful Number Codes that help RAG fetch the right answer!

Vectors capture context and meaning, enabling AI to understand relationships between concepts rather than just matching keywords.

Created by Chinmay Kaitade

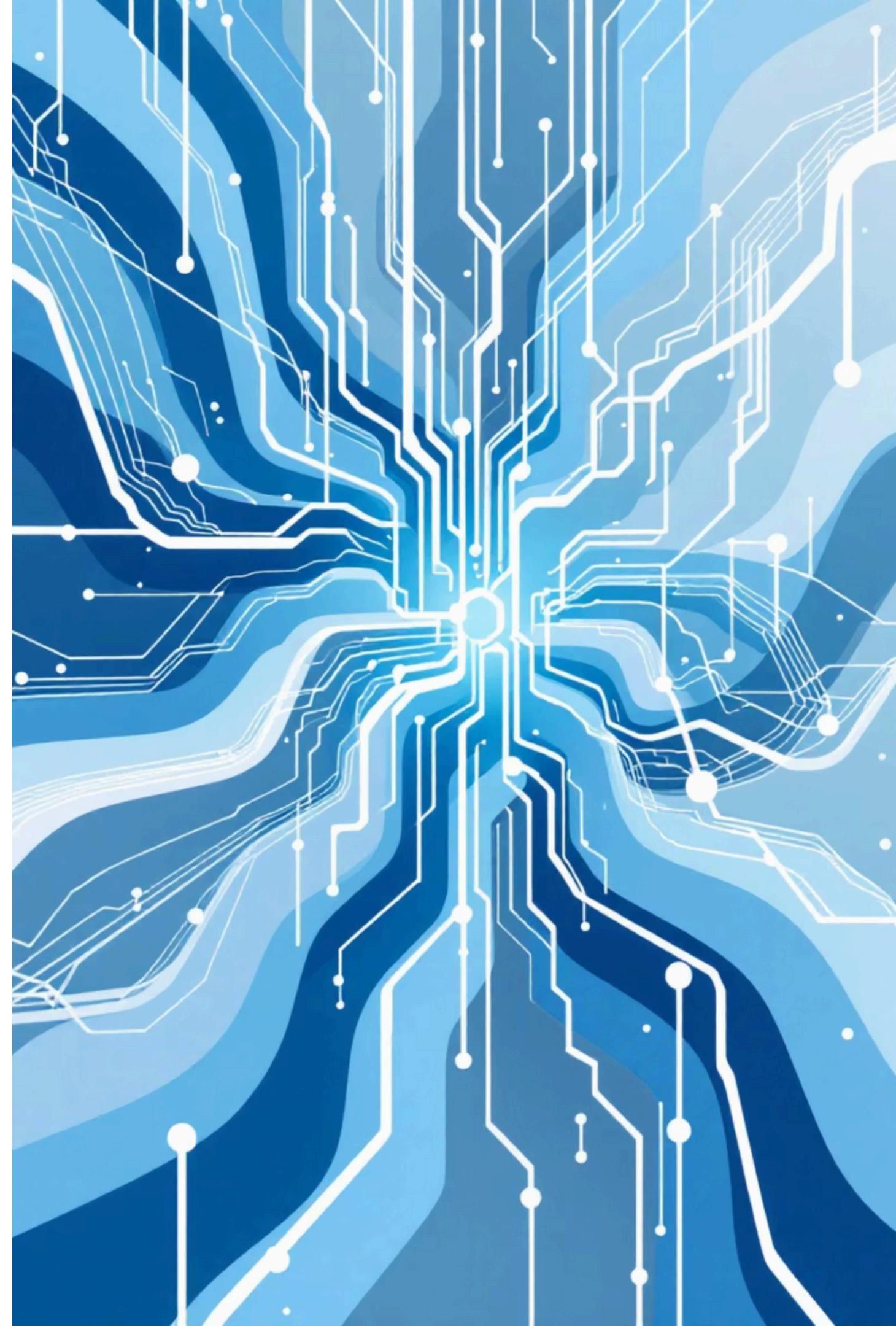
MERN Stack Developer | AI Enthusiast

OpenAI Embedding Models

Industry-leading embedding models with proven effectiveness in large-scale applications

- Pricing Based on Tokens:** Roughly 800 tokens ≈ 1 page of text

Created by Chinmay Kaitade
MERN Stack Developer | AI Enthusiast



Model Comparison

Small Model

~62,500 pages/\$1

Performance: 62.3%

Max Input: 8,192 tokens

 **Recommended** for most projects

Medium Model

~2,000 pages/\$1

Balanced performance

Suitable for mid-level projects

Large Model

~965 pages/\$1

Highest accuracy

Complex projects only

Created by Chinmay Kaitade

MERN Stack Developer | AI Enthusiast

Vector Databases: Storing Smart Data

Specialized databases designed to efficiently store and retrieve high-dimensional vector data

Traditional databases aren't optimized for similarity searches across vector spaces. Vector databases use specialized indexing for lightning-fast semantic searches.

Created by Chinmay Kaitade
MERN Stack Developer | AI Enthusiast



Pinecone: Cloud Vector Database

Pricing Plans

- **Starter:** Free (2GB storage, perfect for testing)
- **Standard:** \$25/month (unlimited storage, multiple projects)
- **Enterprise:** \$500/month (high-scale deployments)

The starter plan covers most initial development and testing needs effectively.



Database Options Comparison

Pinecone

-  Cloud-based solution
-  Easy setup and scaling
-  Generous free tier
-  Recommended for beginners

Qdrant

-  Local/self-hosted option
-  More complex setup
-  Greater control
-  Advanced users only

Created by Chinmay Kaitade

MERN Stack Developer | AI Enthusiast

Connect & Continue Learning

Join my AI learning journey! I share insights, tutorials, and experiences daily.



Instagram

[@chinmaykaitade_hunter](https://www.instagram.com/chinmaykaitade_hunter)



GitHub

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



LinkedIn

[Connect with me](https://www.linkedin.com/in/chinmaykaitade/)



X (Twitter)

[@chinmaydotcom](https://twitter.com/chinmaydotcom)

Created by Chinmay Kaitade

MERN Stack Developer | AI Enthusiast