

INTRODUCTION TO RAG

What is Retrieval-Augmented Generation (RAG)?

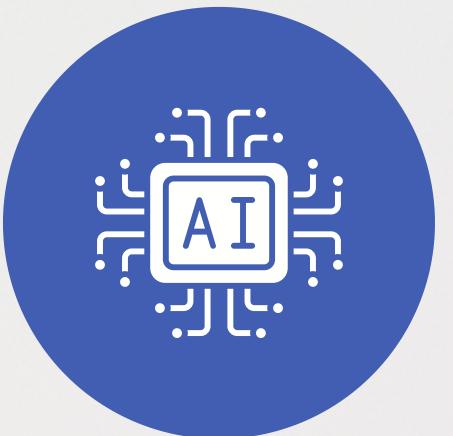
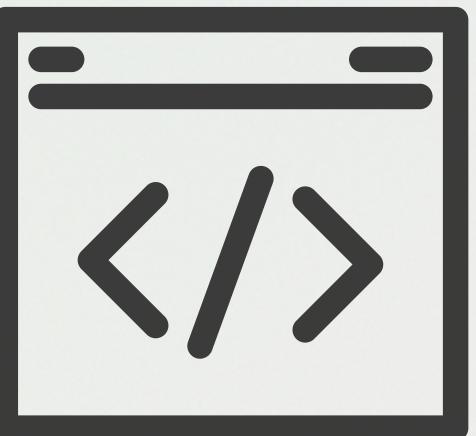
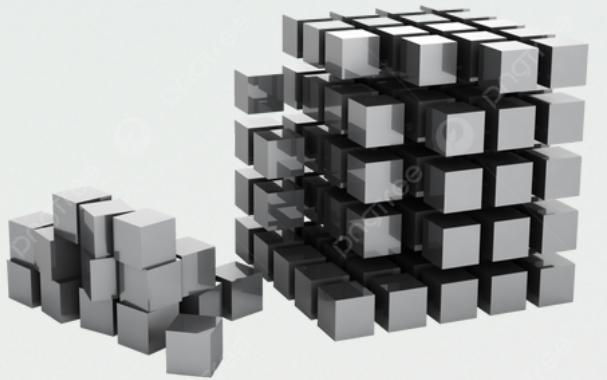
- RAG stands for Retrieval-Augmented Generation
- It helps AI find information before answering
- Gives better and more accurate answers

RAG is a smart way to improve AI. Instead of only using its memory, the AI searches through documents or databases before answering. Just like we use Google to check facts, RAG helps AI do the same!

👉 It's helpful when you want your AI to talk using your own proprietary knowledge base.

Components of RAG

- Chunking
- Embeddings
- Vector Databases
- Semantic Search
- Reranking
- LLM
- User Query
- Response Generation



WHAT IS CHUNKING?

Chunking: Breaking Big Text into Small Parts

- Break big files into smaller sections (chunks)
- Each chunk should have complete meaning
- Helps AI find the right part of text



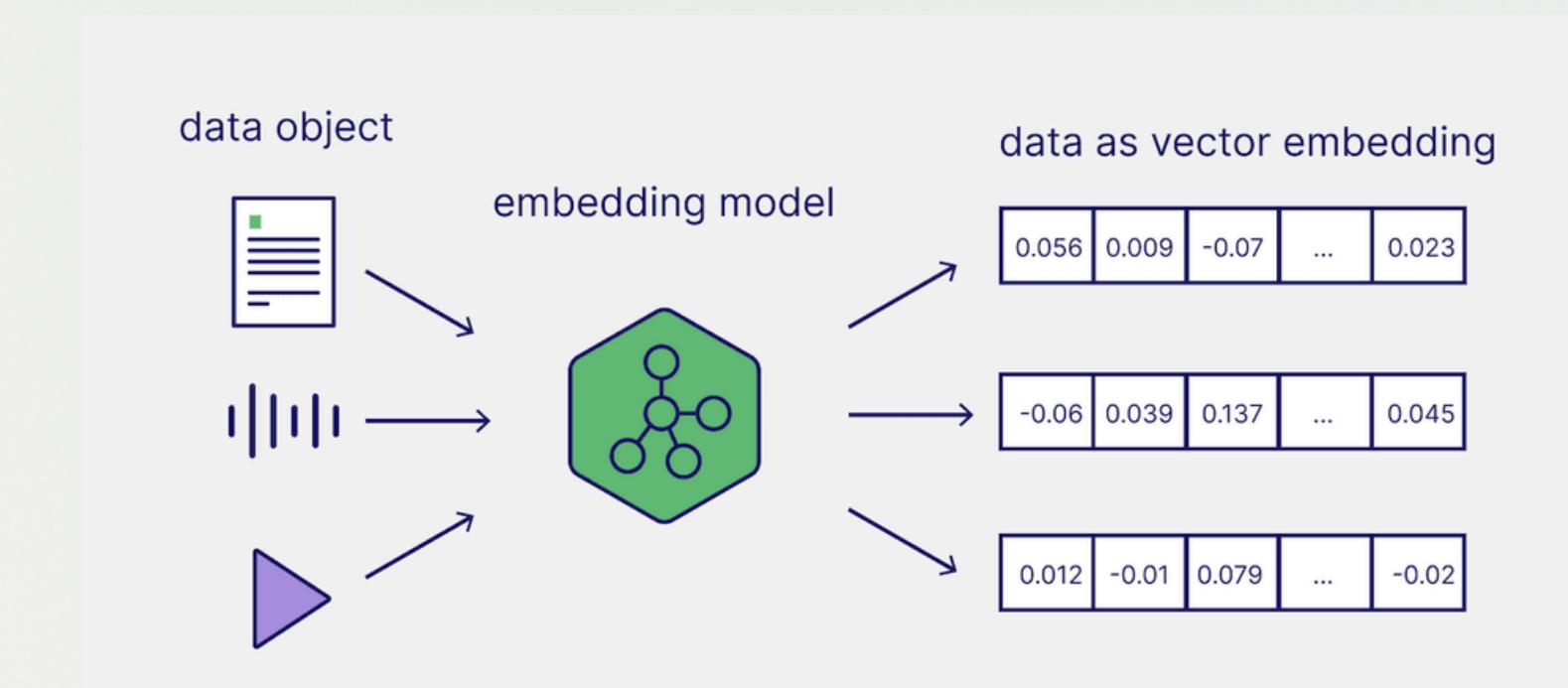
Chunking is like cutting a big book into pages. If we break the content properly, AI can search and understand better. If chunks are too big or too small, AI might miss important details.

👉 A good chunk size is usually between 200 to 500 words depending on the use case.

WHAT IS EMBEDDING?

AI doesn't understand plain text the way humans do – it converts text into numerical representations called embeddings.

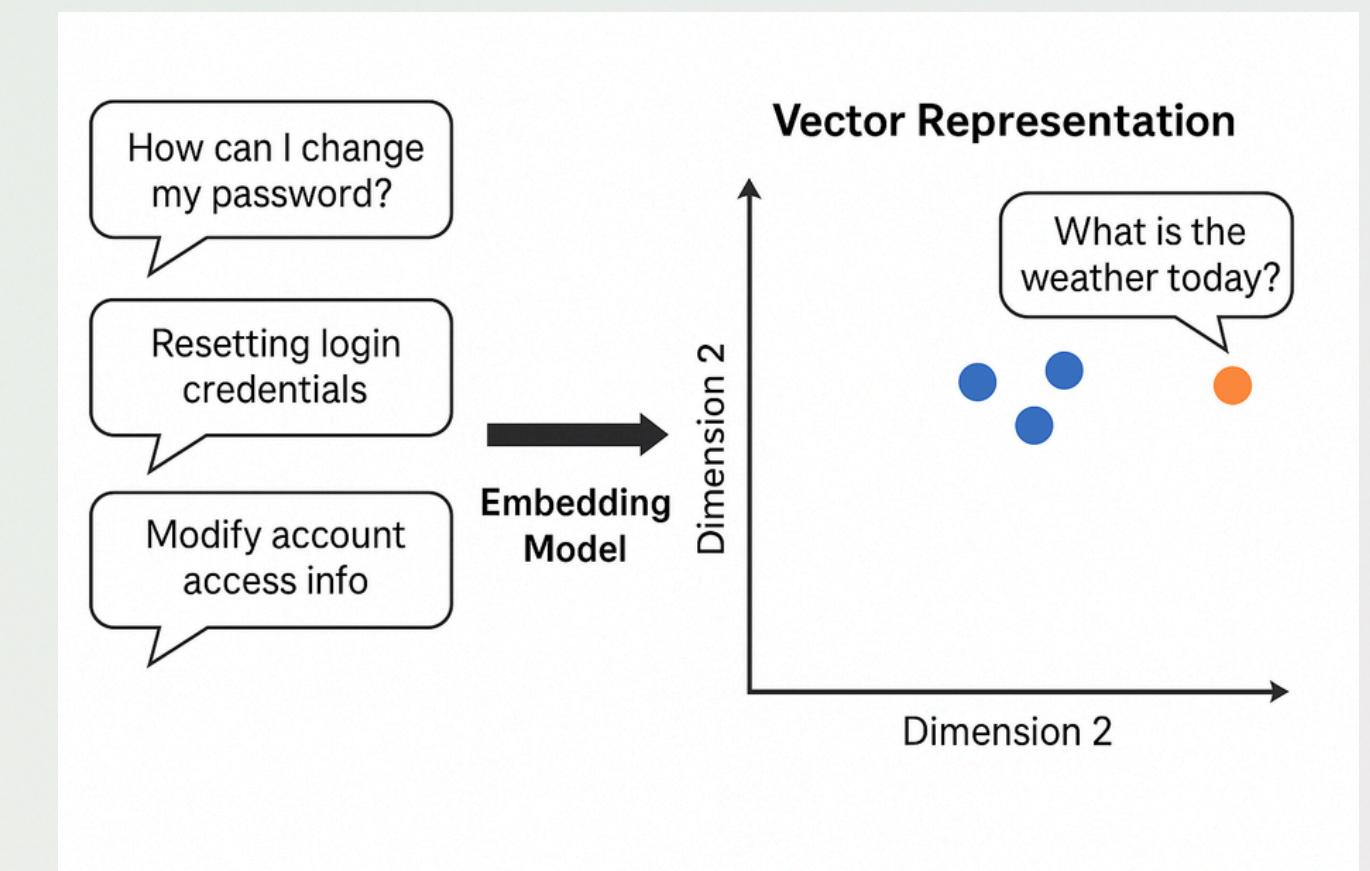
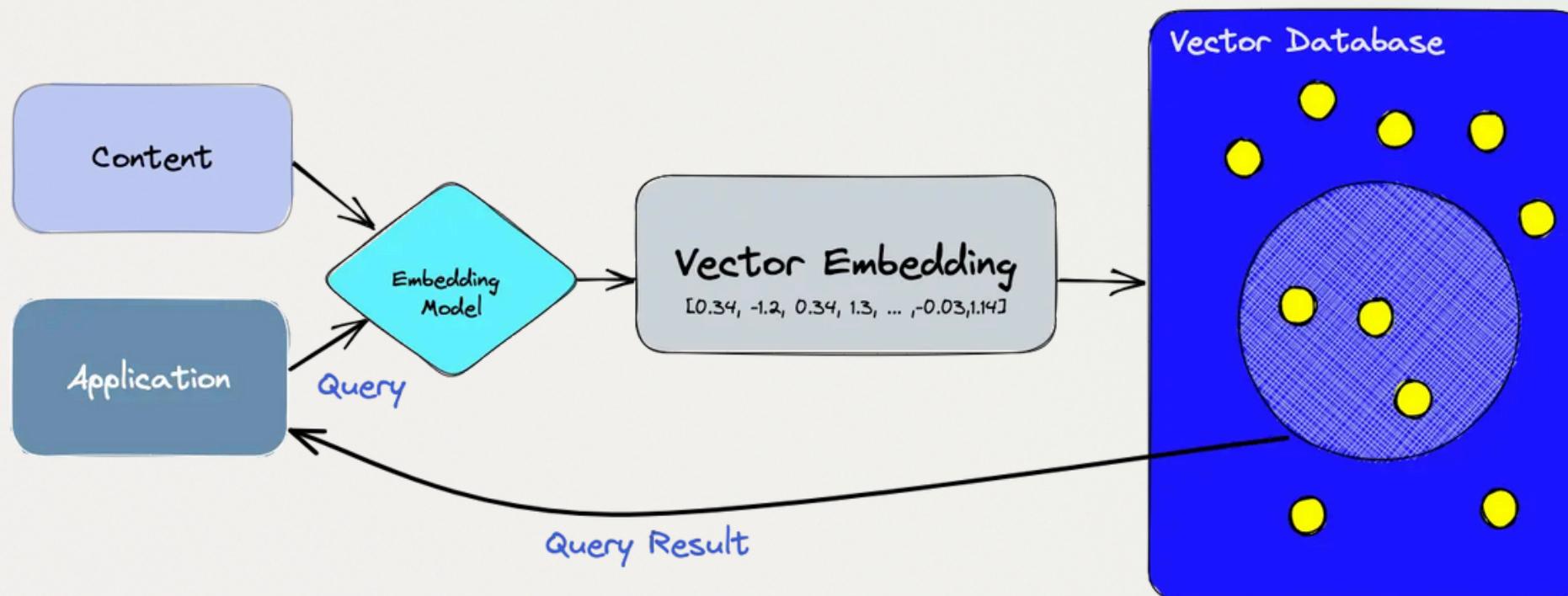
- Embeddings are numerical representations of text
- Similar meanings → similar vectors
- Enables AI to search by meaning, not just keywords
- Good embeddings capture both context and meaning





HOW ARE EMBEDDINGS GENERATED?

Embeddings are numerical representations of text (or other data) generated by AI models that capture the context and meaning.



NEED FOR A VECTOR DATABASE

As AI advances, searching by meaning—not just keywords—is essential. That's where Vector Databases come in.

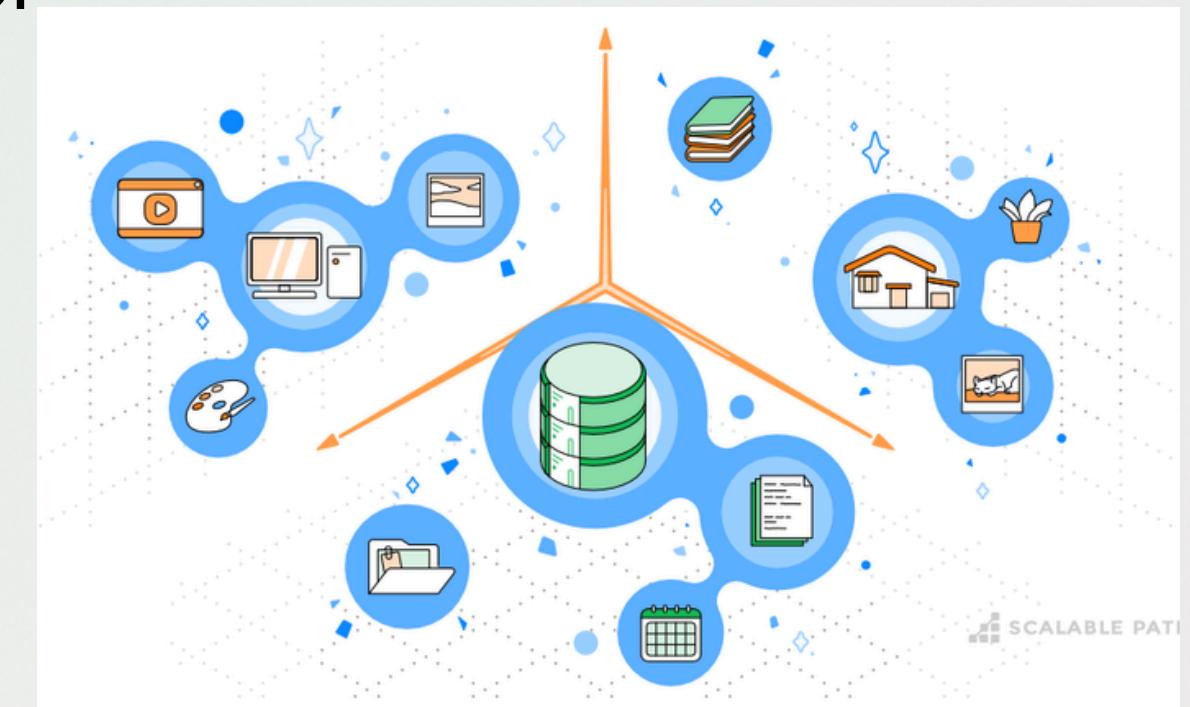
- Traditional DBs can't handle high-dimensional vectors efficiently
- AI models create embeddings that need to be stored and searched
- Essential for recommendations, semantic search, and chatbots
- Vector DBs enable fast, scalable similarity search across millions of embeddings



WHAT IS A VECTOR DATABASE (VECTOR DB)?

A Vector Database stores and manages embeddings (numerical vectors) to enable search by meaning rather than exact matches.

- Stores vector representations of text, images, or other data
- Enables semantic search – finds results based on similarity of meaning
- Ideal for use with AI models to handle recommendations, Q&A, clustering, and more
- Optimized for high-speed similarity search in large datasets



KEY CONCEPTS OF VECTOR DATABASES

Understanding the core ideas behind how Vector DBs work:

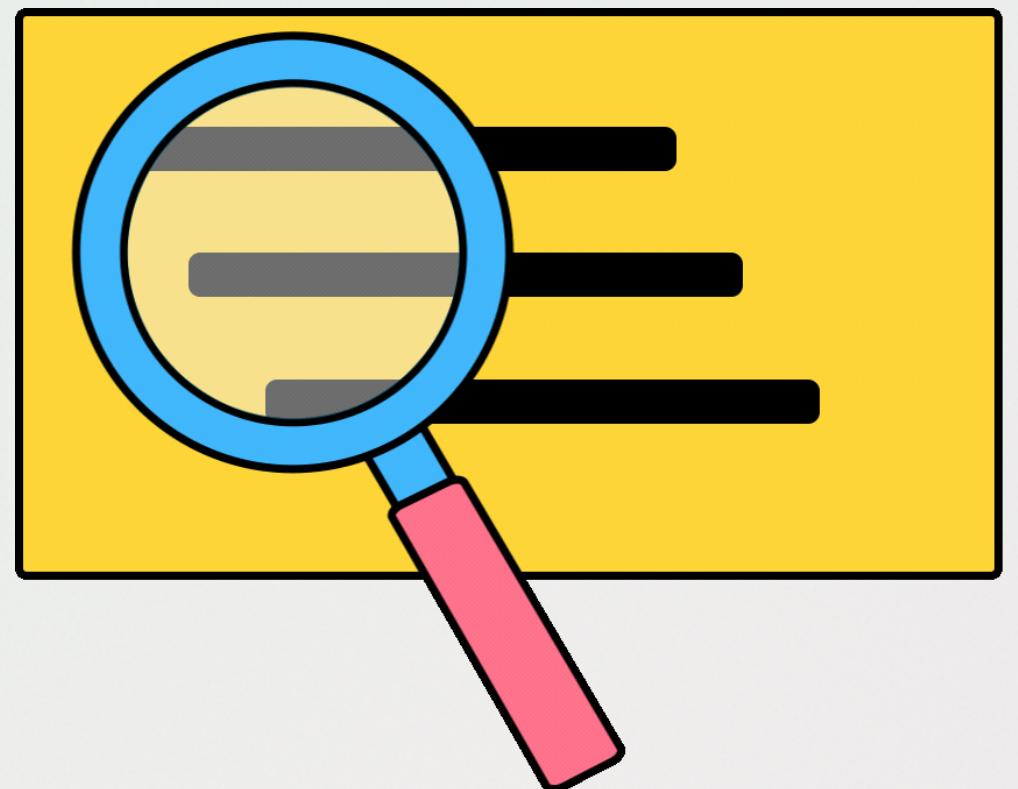
- Embeddings.
- Vector Similarity.
- Indexing.
- Similarity Search (k-NN).
- Scalability.



WHAT IS SEMANTIC SEARCH?

Semantic Search is a search technique that focuses on the meaning behind the words, not just the exact keywords.

- Understands Meaning, Not Just Words
- Goes Beyond Keyword Matching
- Finds Smarter, More Relevant Results
- Powered by AI and Vector Search
- Enhances Search in Smart Applications



HOW SEMANTIC SEARCH WORKS

How It Works:

- Turn Text into Vectors with Embeddings
- Save Those Vectors in a Vector Database
- Convert the User's Question into a Vector
- Search for Similar Vectors (Like Using Cosine Similarity)
- Get Results That Match Meaning, Not Just Words



ASKING A QUESTION (USER QUERY)

- User asks a question
- Question is also converted into a vector
- AI searches for similar chunks in the database

When someone asks a question, the system changes it into a vector (just like chunks). Then it compares this with all saved chunks to find the closest matches.

👉 The quality of the answer depends on how well the question is converted into an embedding.