

Retrieval-Augmented Generation Systems

Retrieval-Augmented Generation Systems

Retrieval-Augmented Generation (RAG) is an advanced AI architecture that combines the strengths of retrieval-based and generation-based models. This approach enhances large language models by providing them with relevant contextual information retrieved from a knowledge base.

How RAG Works:

1. Document Ingestion: Documents are processed and stored in a vector database
2. Embedding Generation: Text is converted into numerical vector representations
3. Query Processing: User queries are embedded using the same model
4. Retrieval: Similar documents are found using vector similarity search
5. Generation: The LLM generates responses using the retrieved context

Benefits of RAG Systems:

- Reduces hallucination by grounding responses in actual documents
- Enables up-to-date information without retraining models
- Provides source attribution and citations
- More cost-effective than fine-tuning large models
- Can work with private or domain-specific knowledge

Key Components:

- Vector Database: ChromaDB, Pinecone, Weaviate, or Milvus
- Embedding Models: OpenAI embeddings, Sentence Transformers, or Cohere
- Large Language Models: GPT-4, Llama, Claude, or Palm
- Document Processors: PyPDF2, Unstructured, or LangChain

RAG systems are particularly useful for question-answering systems, customer support chatbots, research assistants, and any application requiring accurate, source-backed responses from a specific knowledge base.