# RAG

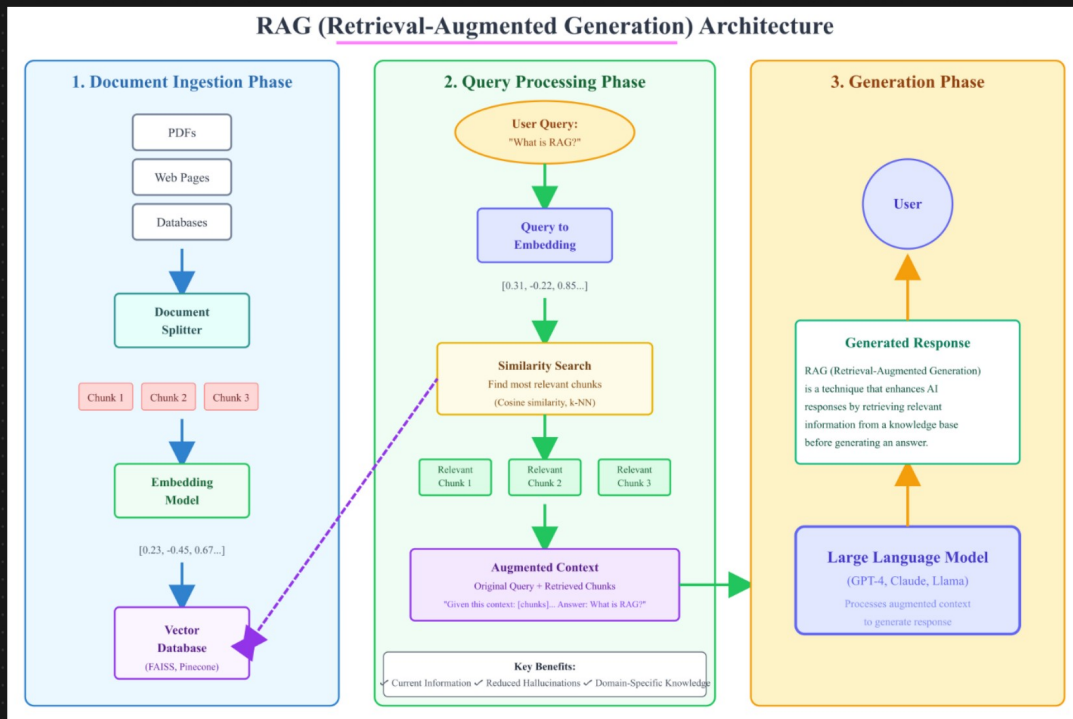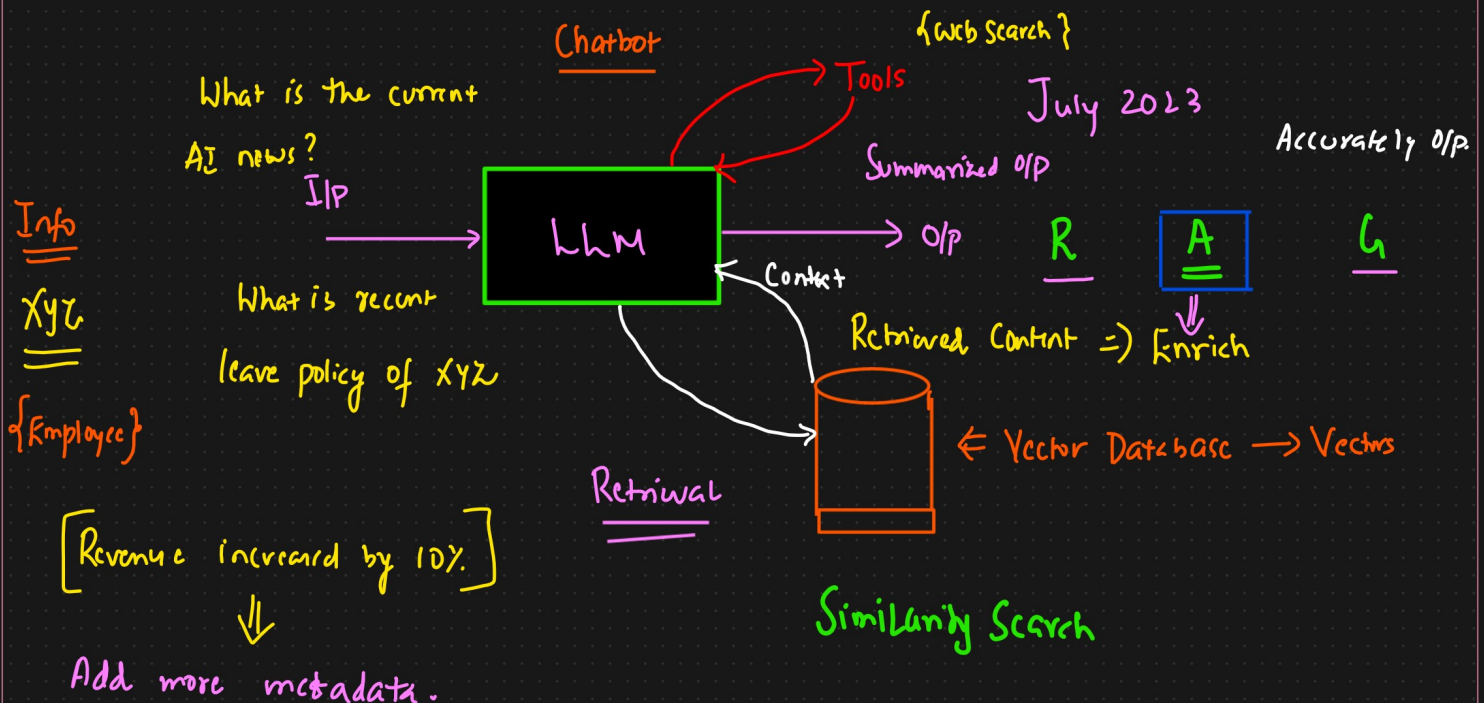**80%**

**RAG**



**RAG (Retrieval-Augmented Generation)** is a powerful technique that enhances AI language models by combining their generation capabilities with external knowledge retrieval.

**What is RAG?**
RAG is like giving an AI assistant access to a library while it's answering questions. Instead of relying solely on what it learned during training, the AI can now look up specific, current, or specialized information from external sources before generating its response.
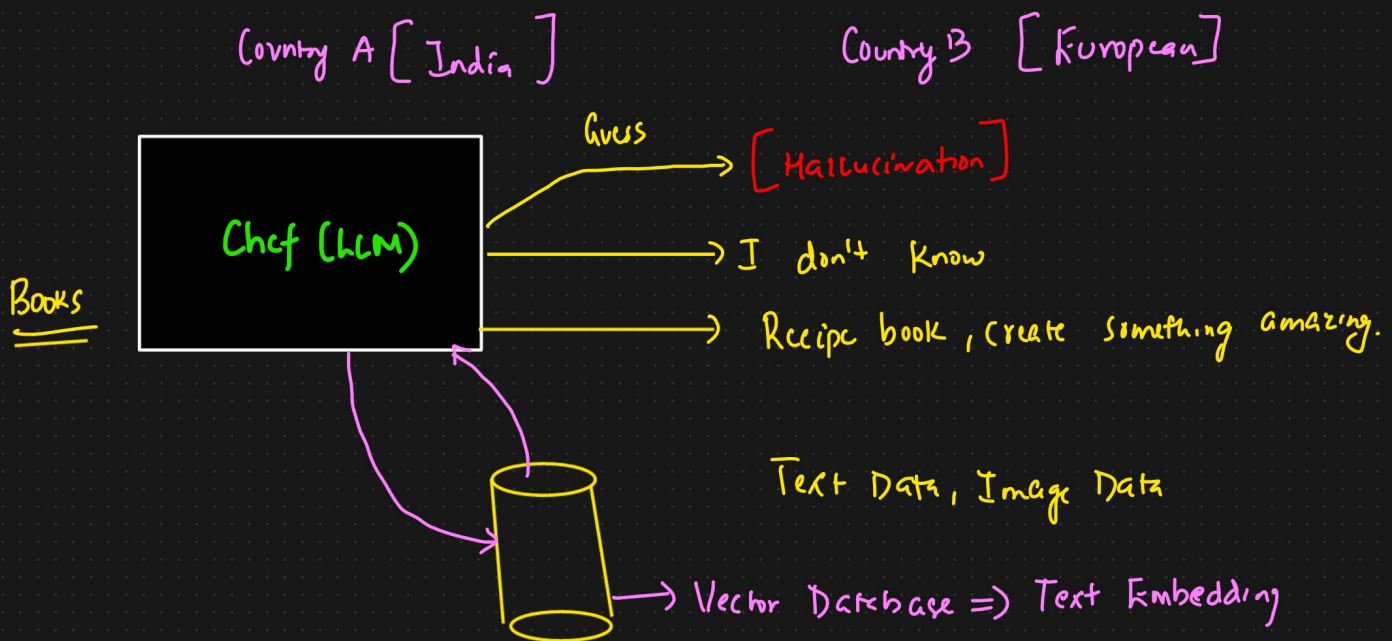Think of it this way: Traditional language models are like students taking a closed-book exam - they can only use what they memorized. RAG-enabled models are like students in an open-book exam - they can reference materials to provide more accurate, detailed, and up-to-date answers.

Chatbot

{Web Search}

July 2023

Accurately O/P

What is the current AI news?

I/P

Tools

Summarized O/P

Info

XYZ

{Employee}

What is recent

leave policy of XYZ

LLM

Context

O/P

R      A      G

Retrieved Content =) Enrich

Retrieval

Vector Database → Vectors

[ Revenue increased by 10% ]

⇓

Add more metadata.

Similarity Search

Augmented = { 
    text: "Revenue increased by 10%.
    "Source:" Tesla Annual Report 2023
    "date"} : —

## 2) RAG

① Retrieval — Finding Relevant Info [Vector Databases]

② Augmentation — Enhancing Context with Metadata.

③ Generation — Producing the Answer.

Country A [India]　　　　　　　Country B [European]

Guess → [Hallucination]

Chef (LLM) →  I don't Know

Books

→ Recipe book, create something amazing.

Text Data, Image Data

→ Vector Database => Text Embedding

## Customer Support

LLM => Customer Support

### Without RAG

Customer: "What's your return policy for items bought during the Black Friday sale?"
AI: "Generally, most companies offer 30-day returns, but policies may vary..."
[Generic, unhelpful response]

Customer: "What's your return policy for items bought during the Black Friday sale?"
[RAG searches company policy database]
AI: "According to our current policy (Policy Doc v3.2, updated Nov 2024), Black Friday purchases have an extended 60-day return window until January 31st. Items must be unused with original tags. Electronics have a 15-day return period due to rapid depreciation. Would you like me to start a return for a specific item?"

**Cost Savings**: JPMorgan saved $150M annually by implementing RAG for research analysis instead of fine-tuning models monthly
**Accuracy**: Microsoft reported 94% reduction in AI hallucinations after implementing RAG in their Copilot products
**Flexibility**: Bloomberg updates their financial AI assistant hourly with new market data - impossible with traditional LLMs
**Compliance**: Healthcare companies use RAG to ensure AI responses always cite approved medical sources"