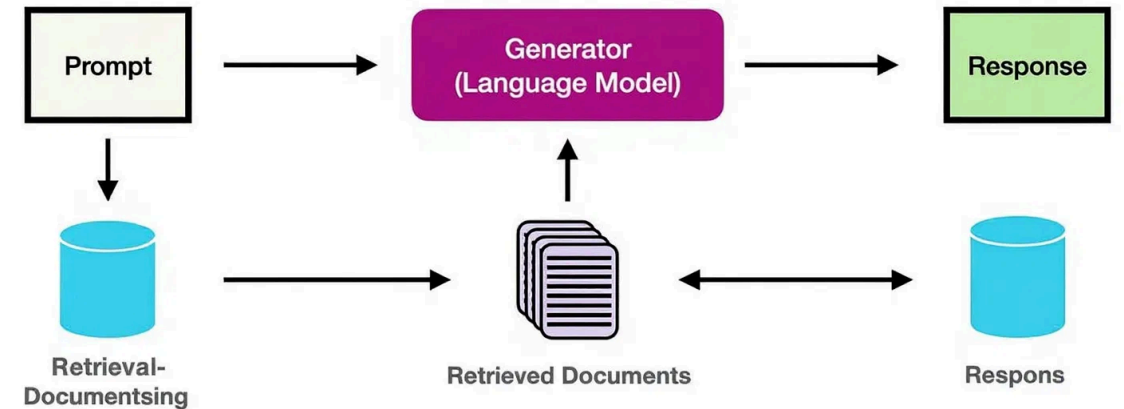


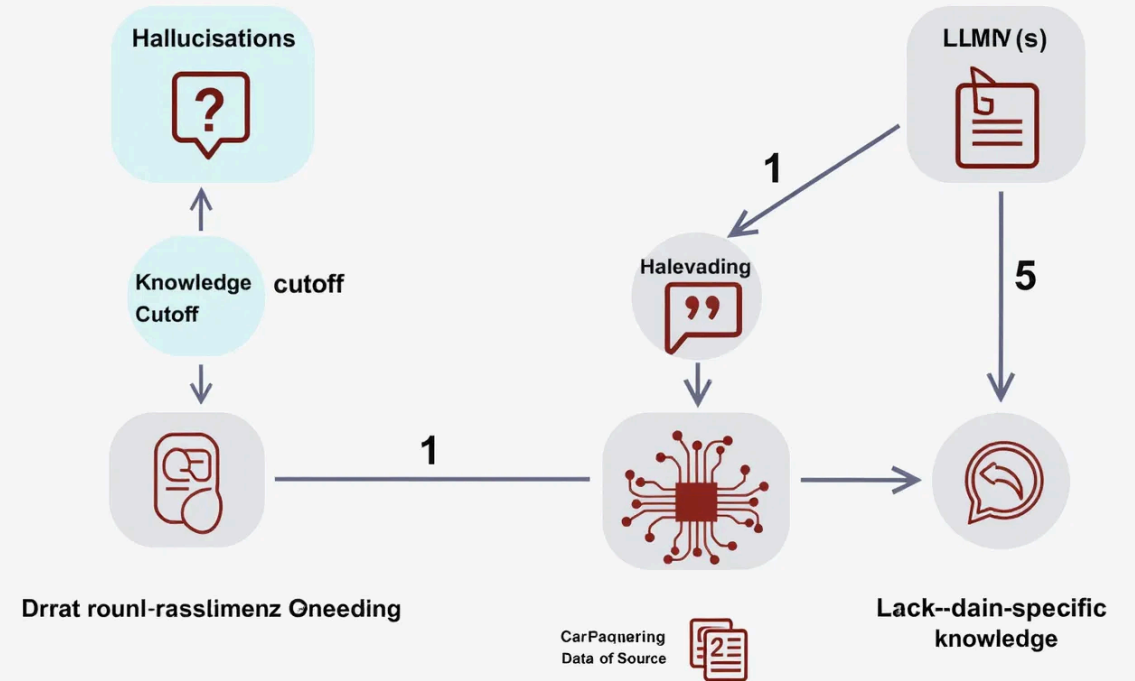
Introduction to RAG

- ✓ **Retrieval-Augmented Generation (RAG)** is an AI framework that enhances LLMs with external knowledge retrieval
- ✓ Combines the strengths of **retrieval-based** and **generation-based** AI approaches
- ✓ Introduced by Meta AI Research in 2020 as a method to improve factual accuracy
- ✓ Enables LLMs to access and leverage information beyond their training data
- ✓ Produces more accurate, up-to-date, and contextually relevant responses








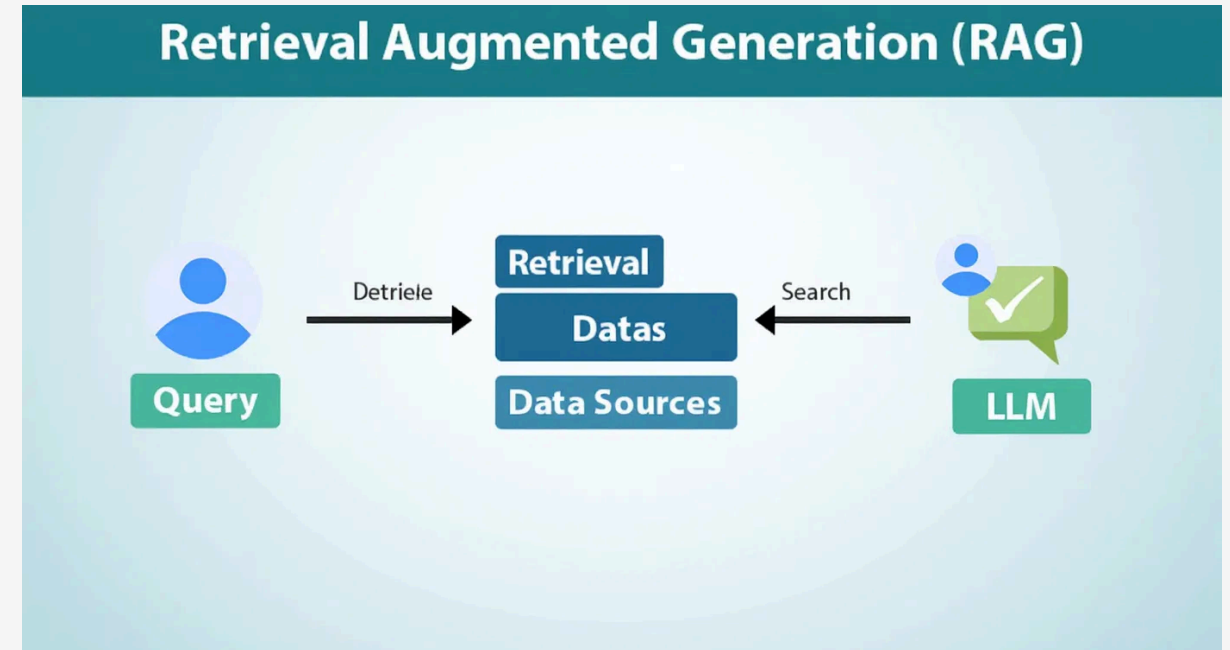
Limitations of Standard LLMs

- ⚠ **Hallucinations:** LLMs can generate plausible but factually incorrect information
- ⚠ **Knowledge Cutoff:** Limited to information available up to their training cutoff date
- ⚠ **Domain Limitations:** Lack specialized knowledge in specific fields or proprietary information
- ⚠ **Source Attribution:** Difficulty in citing sources or providing evidence for generated content
- ⚠ **Context Window:** Limited ability to process and reference large amounts of information








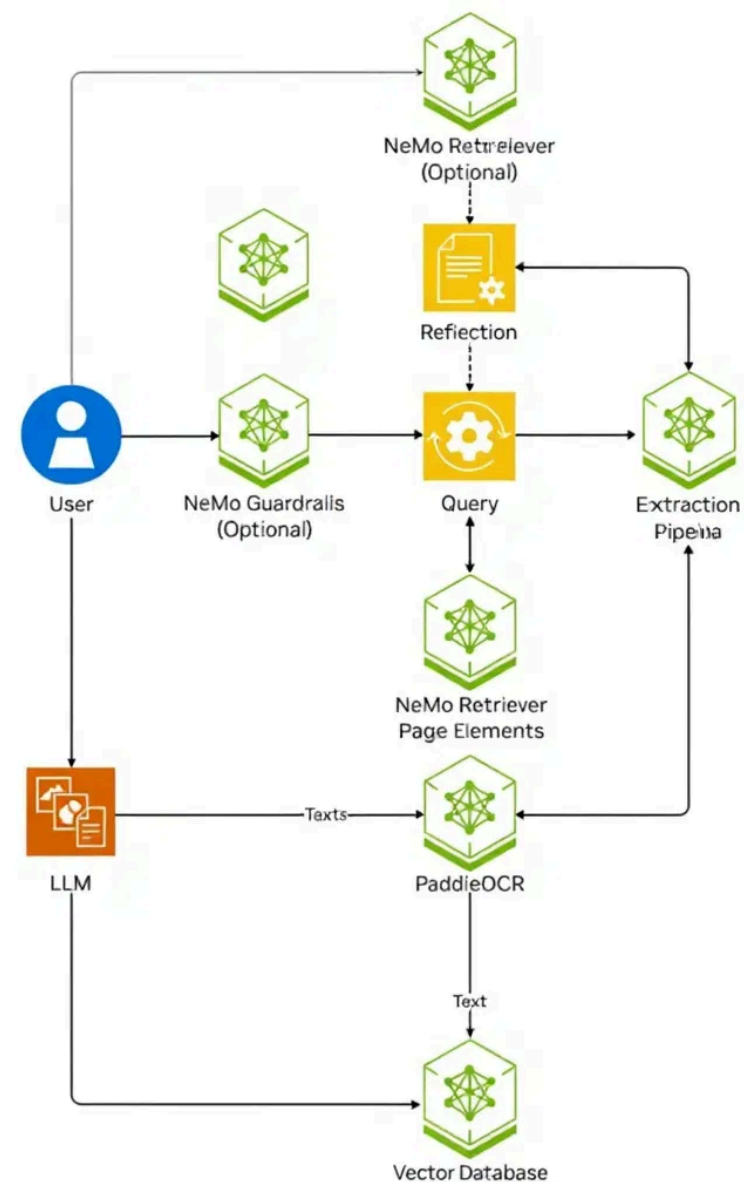
How RAG Works

-  **1. Retriever:** Searches for relevant information from external knowledge sources
-  **2. Knowledge Base:** Contains indexed documents, often using vector embeddings for semantic search
-  **3. Context Augmentation:** Retrieved information is added to the prompt as context
-  **4. Generator:** LLM uses the augmented context to produce accurate, grounded responses
-  **5. End-to-End Flow:** Query → Retrieval → Augmentation → Generation → Response








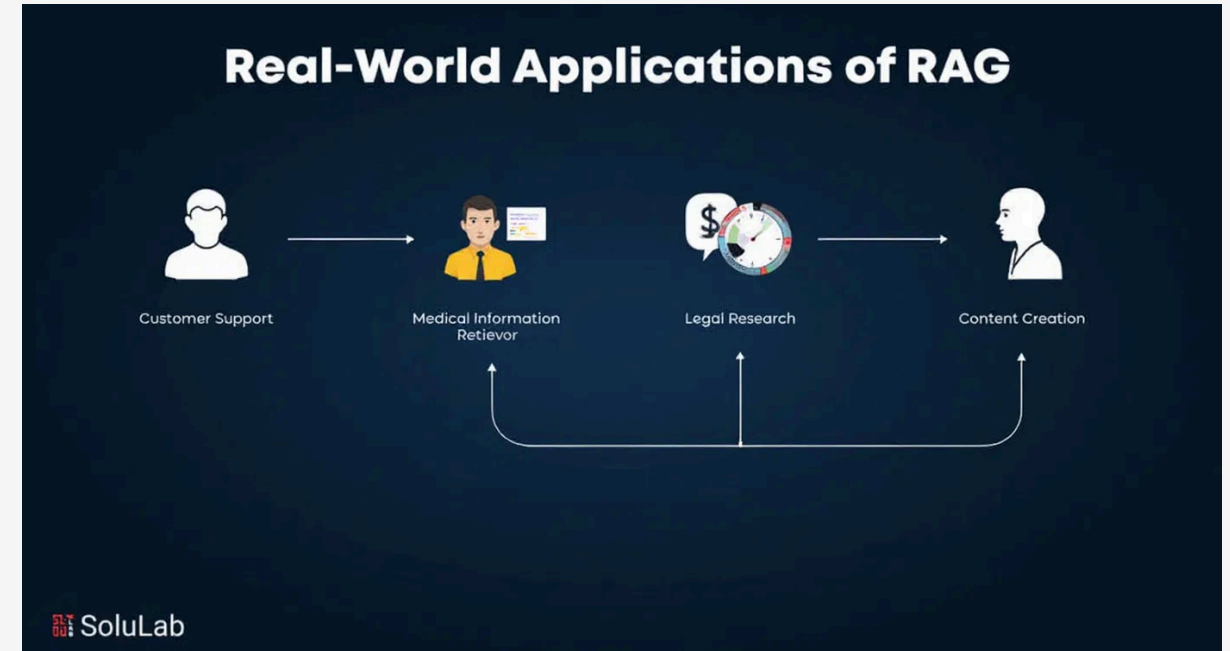
Key Components and Architecture

-  **Knowledge Base:** Document store containing indexed information (PDFs, web pages, databases)
-  **Embedding Model:** Converts text into vector representations for semantic search
-  **Retrieval Pipeline:** Query processing, vector search, and relevance ranking
-  **Integration Layer:** Combines retrieved context with user query for the LLM
-  **Generator:** LLM that produces the final response using augmented context



Real-World Applications

-  **Customer Support:** Chatbots that access company knowledge bases to provide accurate, contextual responses
-  **Healthcare:** Clinical decision support systems that retrieve medical literature and patient records
-  **Legal Research:** Document analysis tools that search case law and legal precedents
-  **Content Creation:** Writing assistants that incorporate factual information and citations
-  **Education:** Personalized tutoring systems that access textbooks and learning materials



Benefits and Challenges

Benefits




- ✓ **Improved Accuracy:** Reduces hallucinations by grounding responses in factual data
- ✓ **Current Information:** Accesses up-to-date knowledge beyond training cutoff
- ✓ **Domain Expertise:** Incorporates specialized knowledge from proprietary sources

Challenges




- ! **System Complexity:** More components to maintain and optimize
- ! **Retrieval Quality:** Results depend on the quality of search and indexing
- ! **Computational Cost:** Additional processing overhead compared to standard LLMs

RAG: Benefits & Challenges

Benefits

-  Improved Accuracy
-  Access to Current Information
-  Reduced Hallucinations

Challenges

-  Complexity
-  Retrieval Quality
-  Computational Cost