# What is a RAG System?

## Instructor

### Dipanjan Sarkar

Head of Community & Principal AI Scientist at Analytics Vidhya

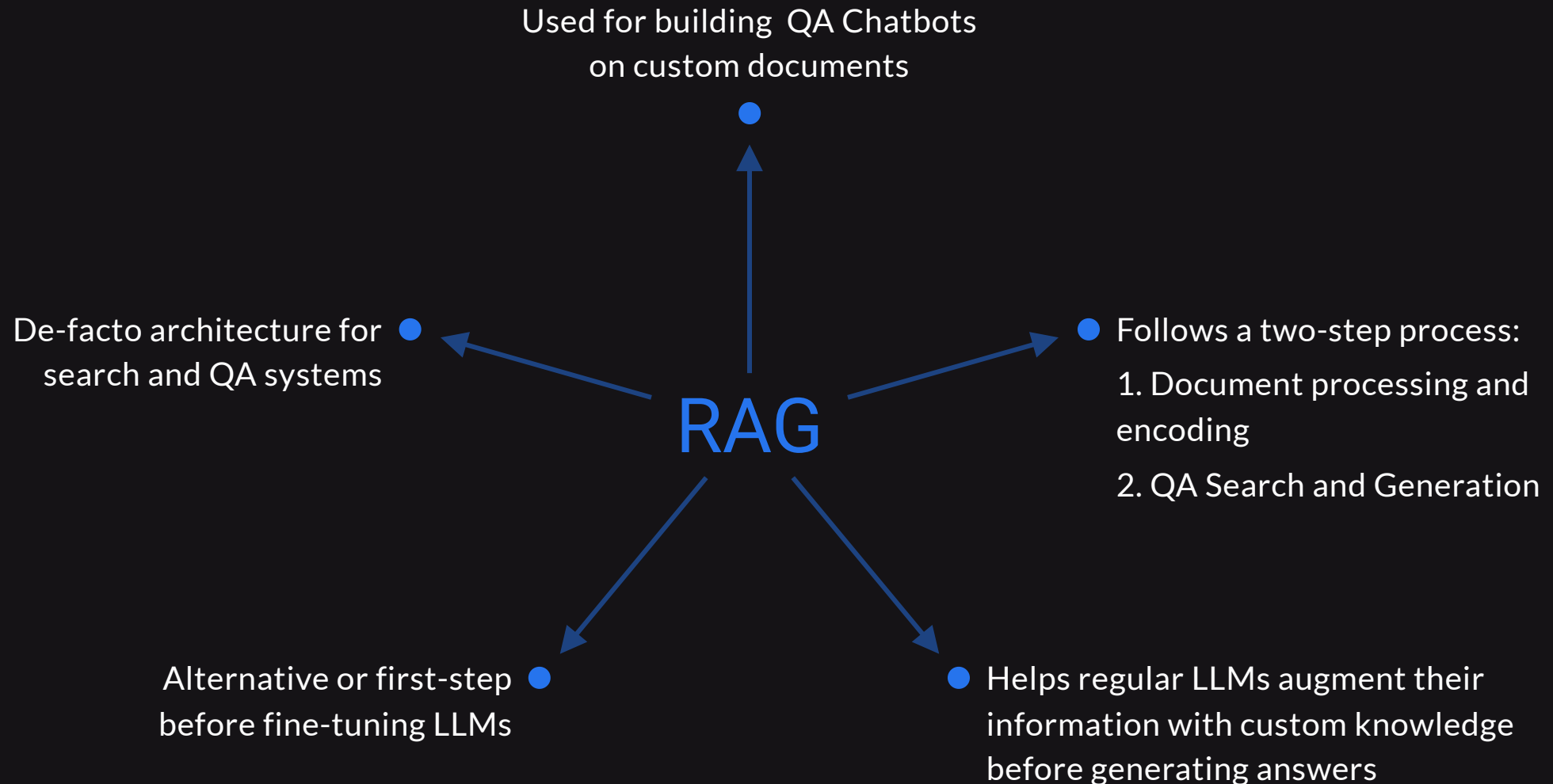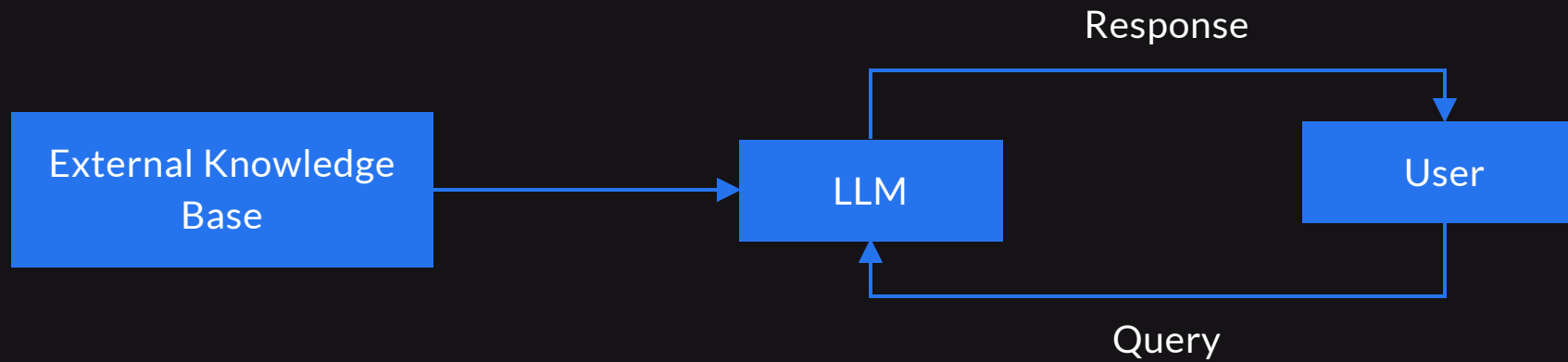Google Developer Expert - ML & Cloud Champion Innovator

Published Author

# Outline

- What is a RAG System

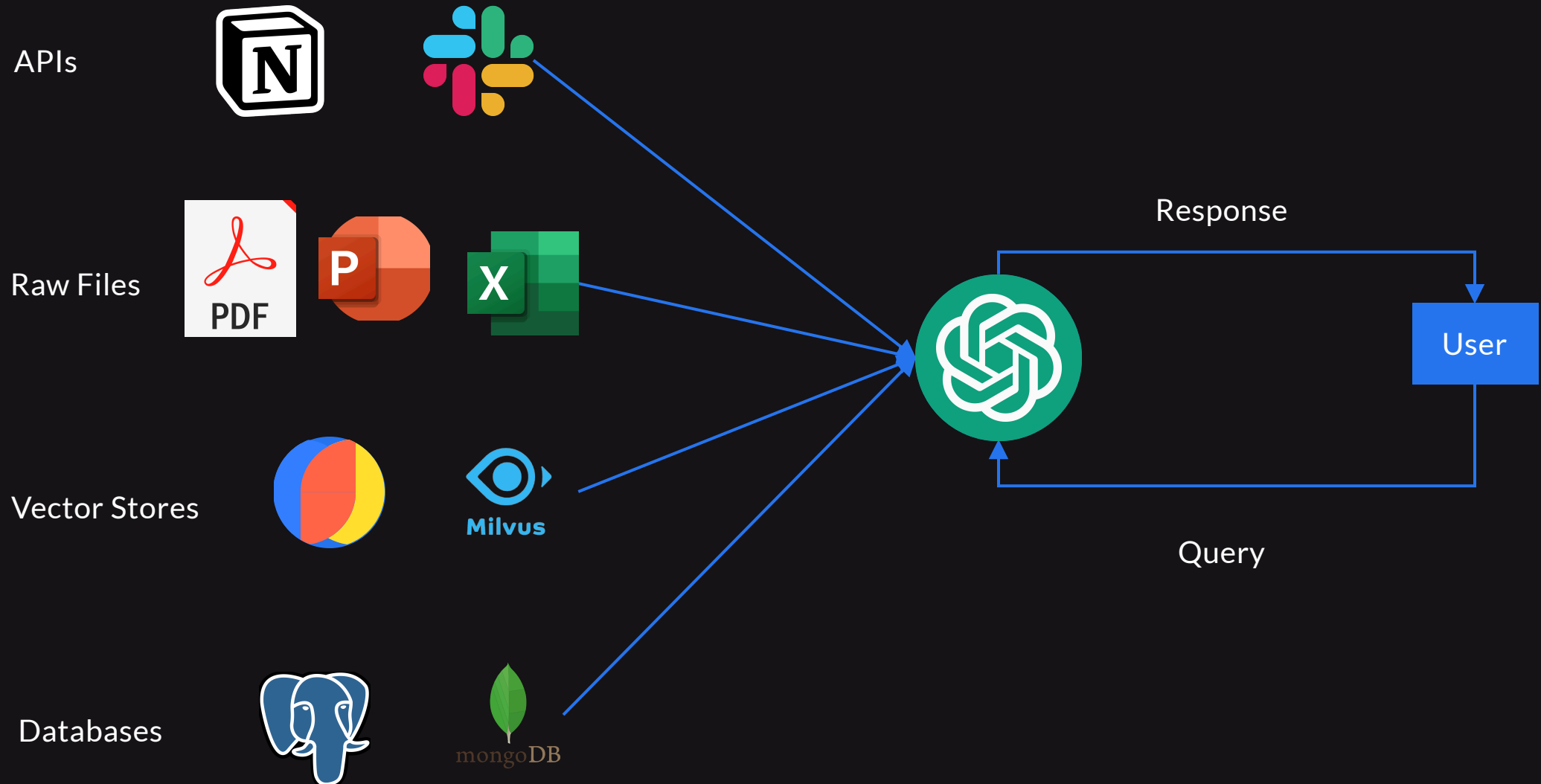- RAG Workflow

- RAG System Architecture

# What is a RAG System



External Knowledge Base → LLM

LLM → User: Response

User → LLM: Query

# What is a RAG System



APIs

Raw Files

Vector Stores

Databases

Response

Query

User

Analytics Vidhya

# What is a RAG System

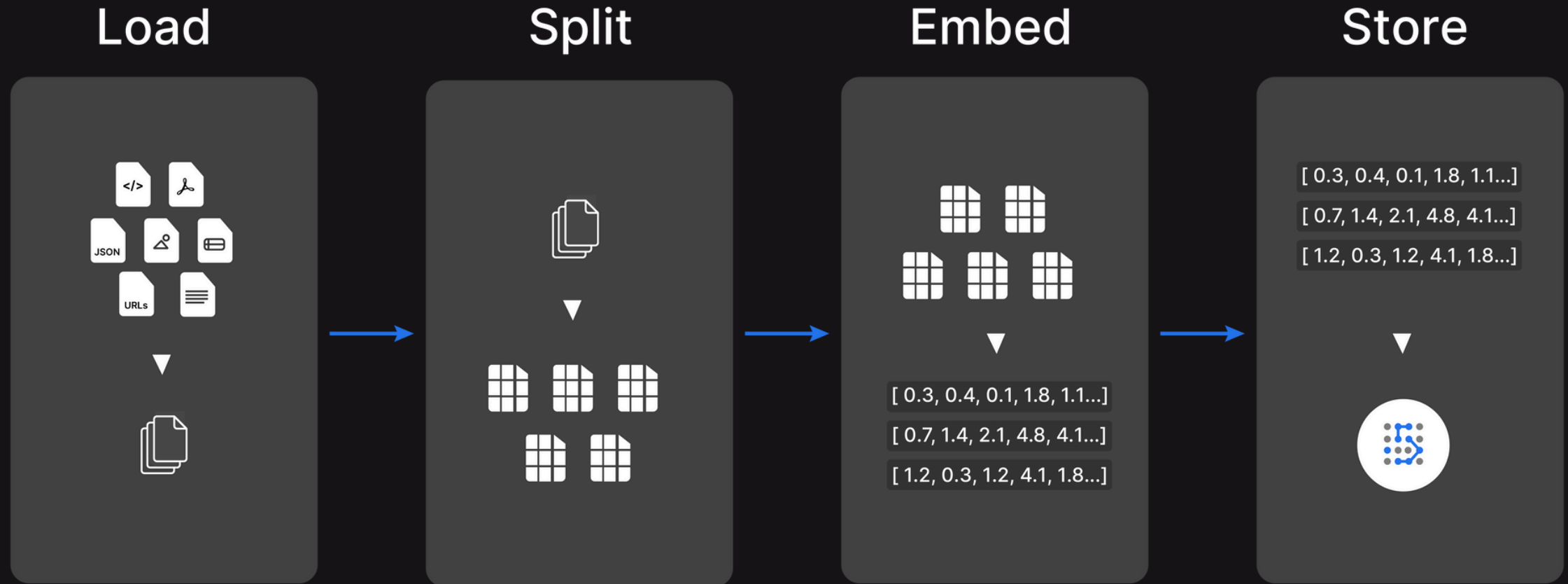| Step 1 | Step 2 |
|---|---|
| Data Processing & Indexing | Retrieval & Response Generation |

- Custom documents are processed and chunked
- These chunks are converted to embeddings with a LLM (transformer)
- Chunks and embeddings are stored in a Vector Database index (along with metadata)

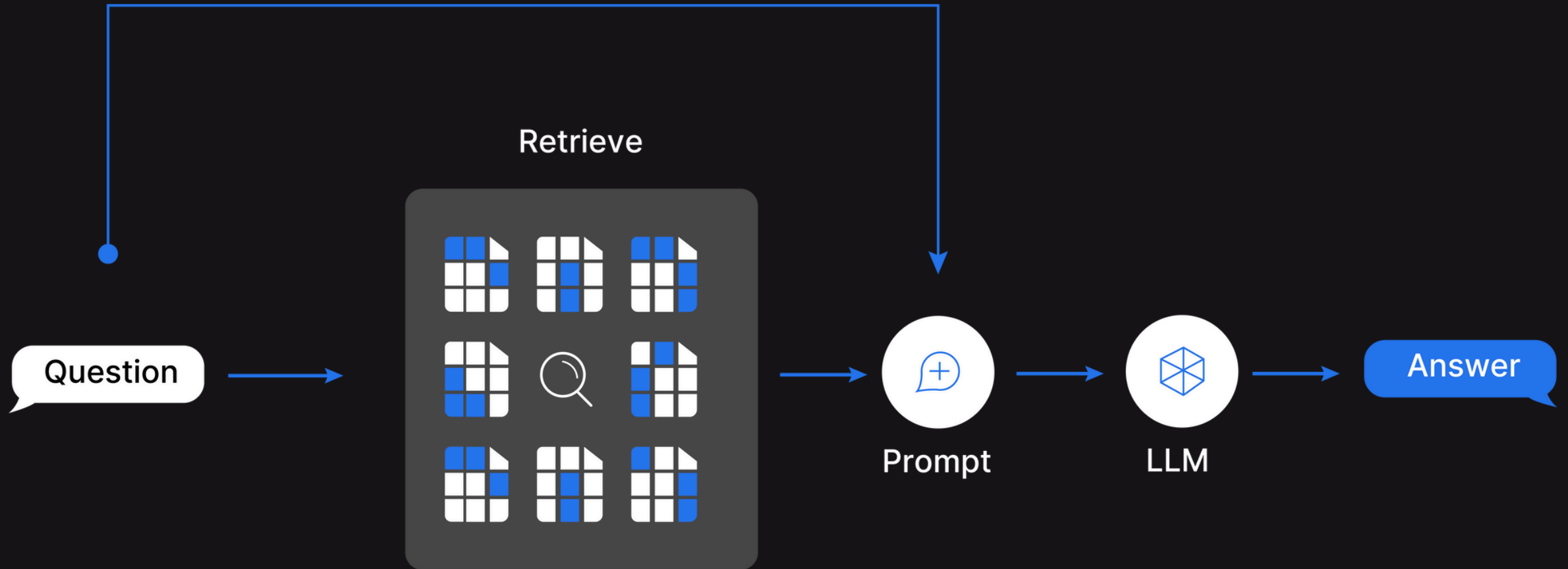- Based on user query the system retrieves relevant document chunks
- Passes these chunks to the LLM along with the query to augment its knowledge
- The LLM generates a human-like response to the user query based on this information

Analytics
Vidhya

# RAG Workflow - Step 1 - Data processing and Indexing



**Load**

**Split**

**Embed**

[ 0.3, 0.4, 0.1, 1.8, 1.1...]
[ 0.7, 1.4, 2.1, 4.8, 4.1...]
[ 1.2, 0.3, 1.2, 4.1, 1.8...]

**Store**

[ 0.3, 0.4, 0.1, 1.8, 1.1...]
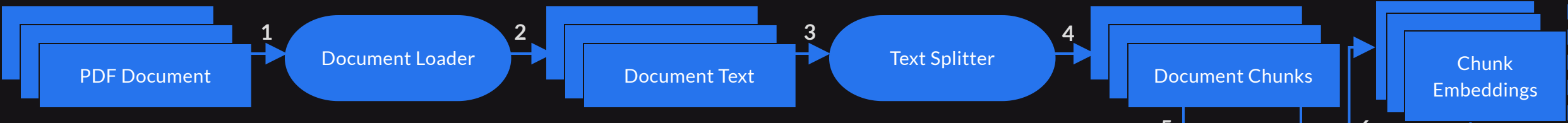[ 0.7, 1.4, 2.1, 4.8, 4.1...]
[ 1.2, 0.3, 1.2, 4.1, 1.8...]

Analytics Vidhya

RAG Workflow - Step 2 - Retrieval and Response Generation

# RAG System Architecture



**Step 1 - Data processing & Indexing**

PDF Document → (1) → Document Loader → (2) → Document Text → (3) → Text Splitter → (4) → Document Chunks → (6) → Chunk Embeddings

Document Chunks → (5) → LLM Embedder → (7) → Vector DB

Chunk Embeddings → (7) → Vector DB

**Step 2 - Retrieval & Response Generation**

User → (8) → User Query → (9) → LLM Embedder

User Query → (10) → Query Embedding → (11) → DB Retriever

Vector DB → (12) → DB Retriever → (13) → Similar Document Chunks

User Query → (14) → Prompt Template

Similar Document Chunks → (14) → Prompt Template → (15) → LLM Prompt → (16) → LLM → (17) → Generated Response → (18) → User

Analytics Vidhya