# Deploying and Maintaining RAG Systems

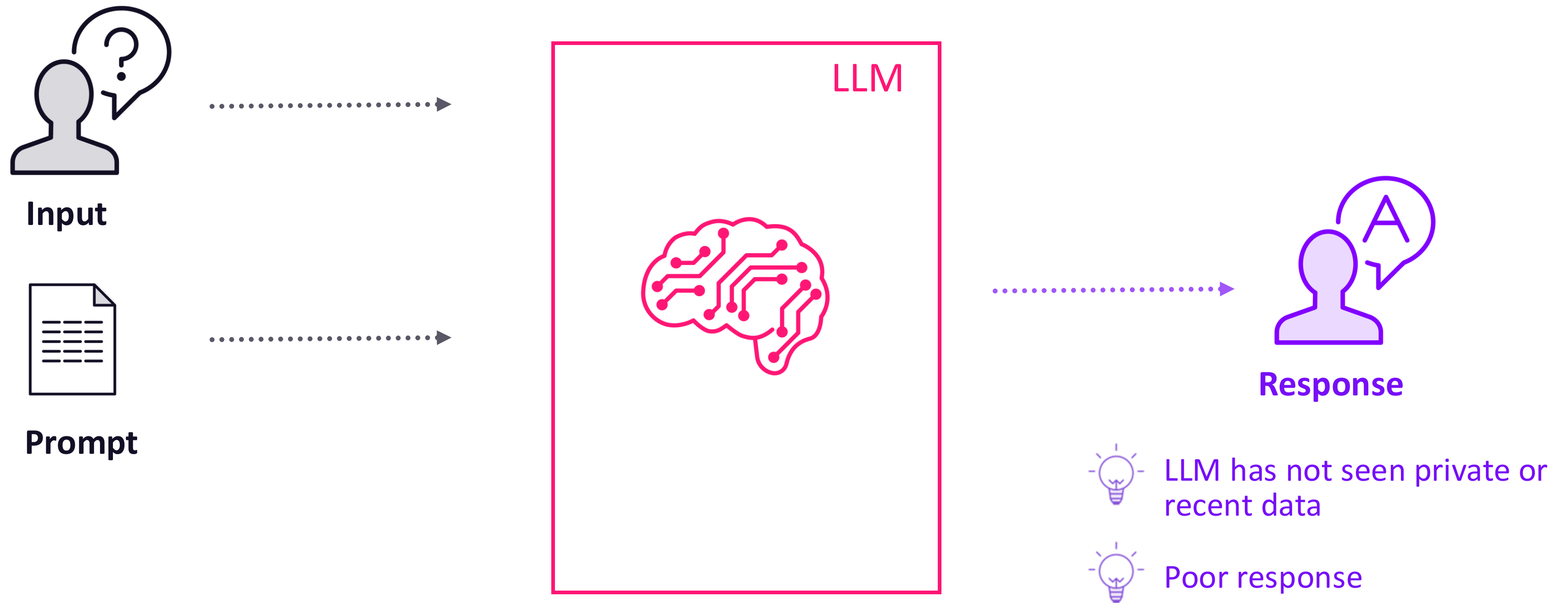**Building and Deploying RAG in Production**

**Abhishek Kumar**

Data Scientist | Author | Speaker
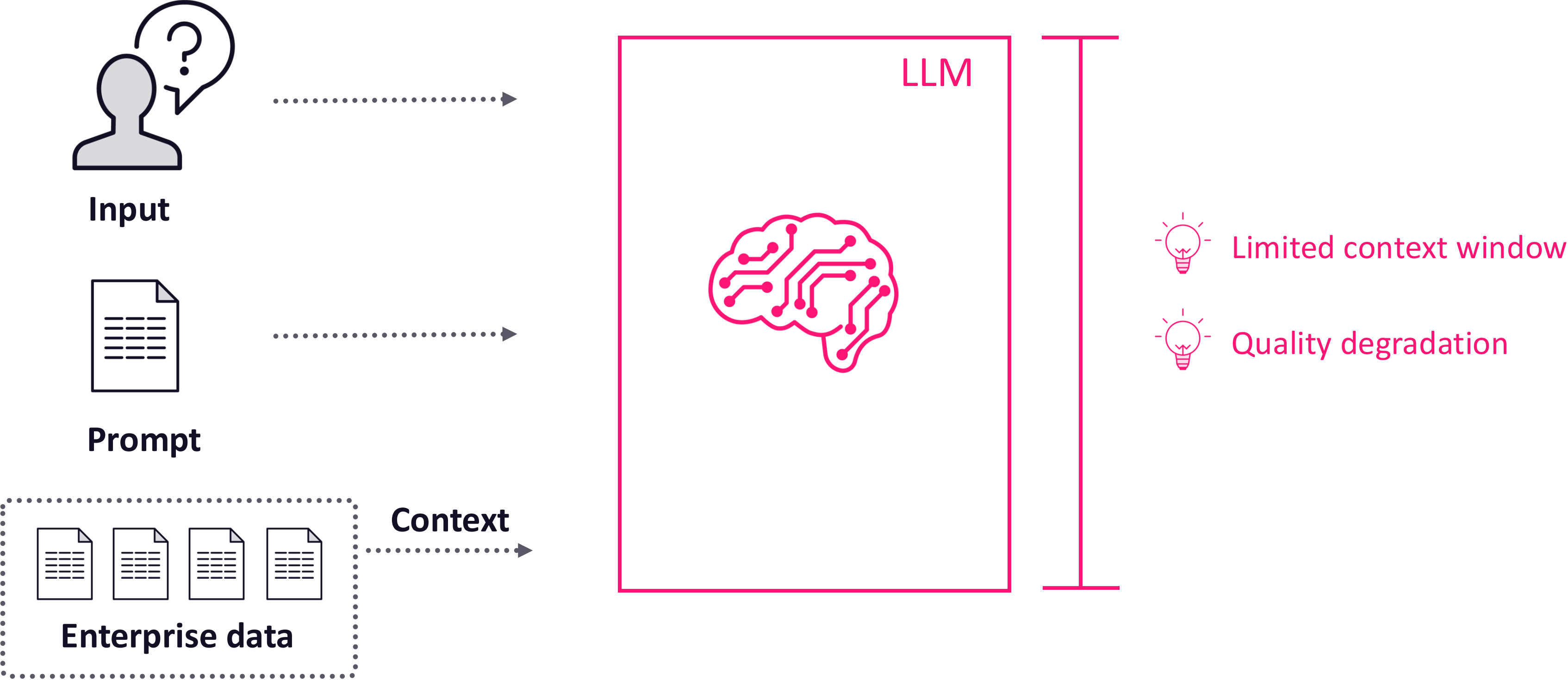
@meabhishekkumar
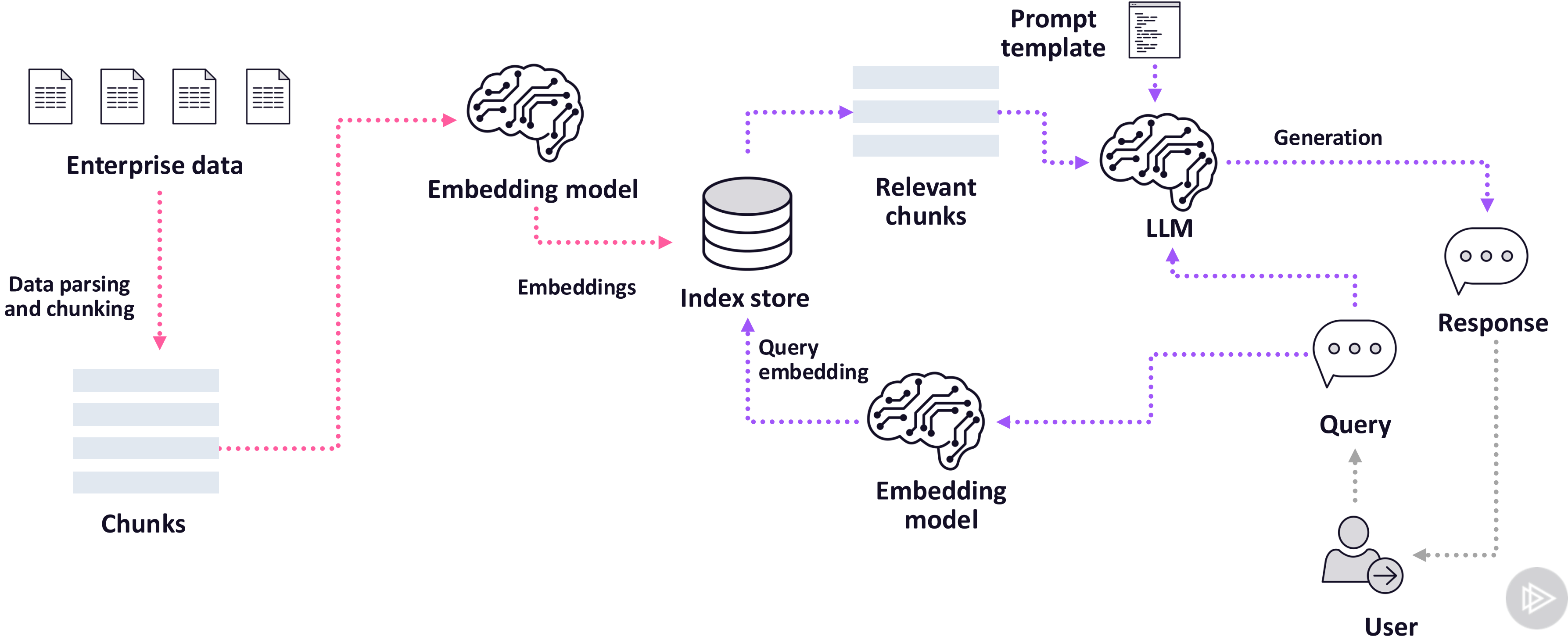
# Prompt Engineering Alone Is Not Sufficient

# Can We Not Stuff Everything in Prompt?

Input

Prompt

Context

Enterprise data

LLM

Limited context window

Quality degradation

# Retrieval Augmented Generation (RAG) to the rescue

# RAG System

# Demo

**Creating embeddings**

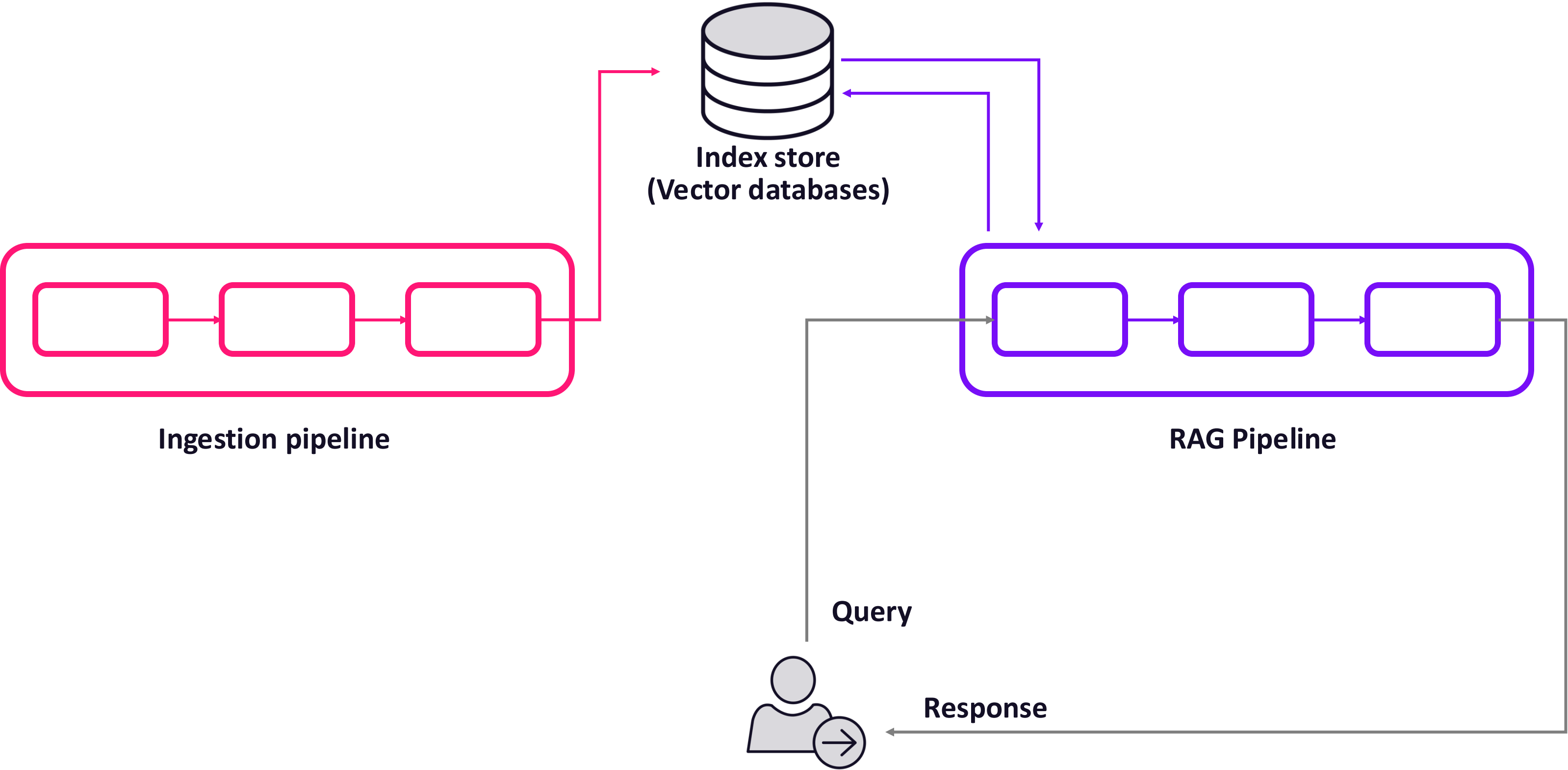# Benefits of RAG Systems

Include your private and recent data

Low LLM cost based on usage

Grounding and reducing hallucination

Source citations and attribution

# RAG System

# Many Choices to Be Made

| | | |
|---|---|---|
| **Parsing and chunking strategy** | **Embedding model and configuration** | **Retrieval techniques** |
| **Top N chunks** | **LLM and associated parameters** | **Prompt template** |

Building all RAG system components from scratch can be tedious.

# RAG Integration Frameworks

Fasttrack RAG system development

Modular components and pre-configured chains

Adapt as per your needs

Open source
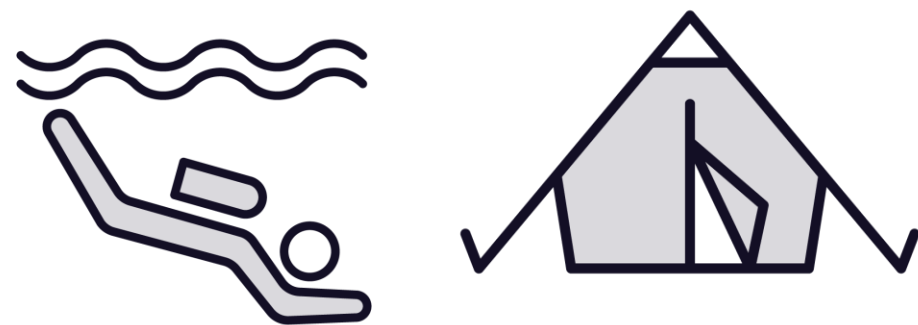- LlamaIndex
- Langchain

# Demo

**Building simple RAG system**

- Using LlamaIndex framework

# Online Retailer
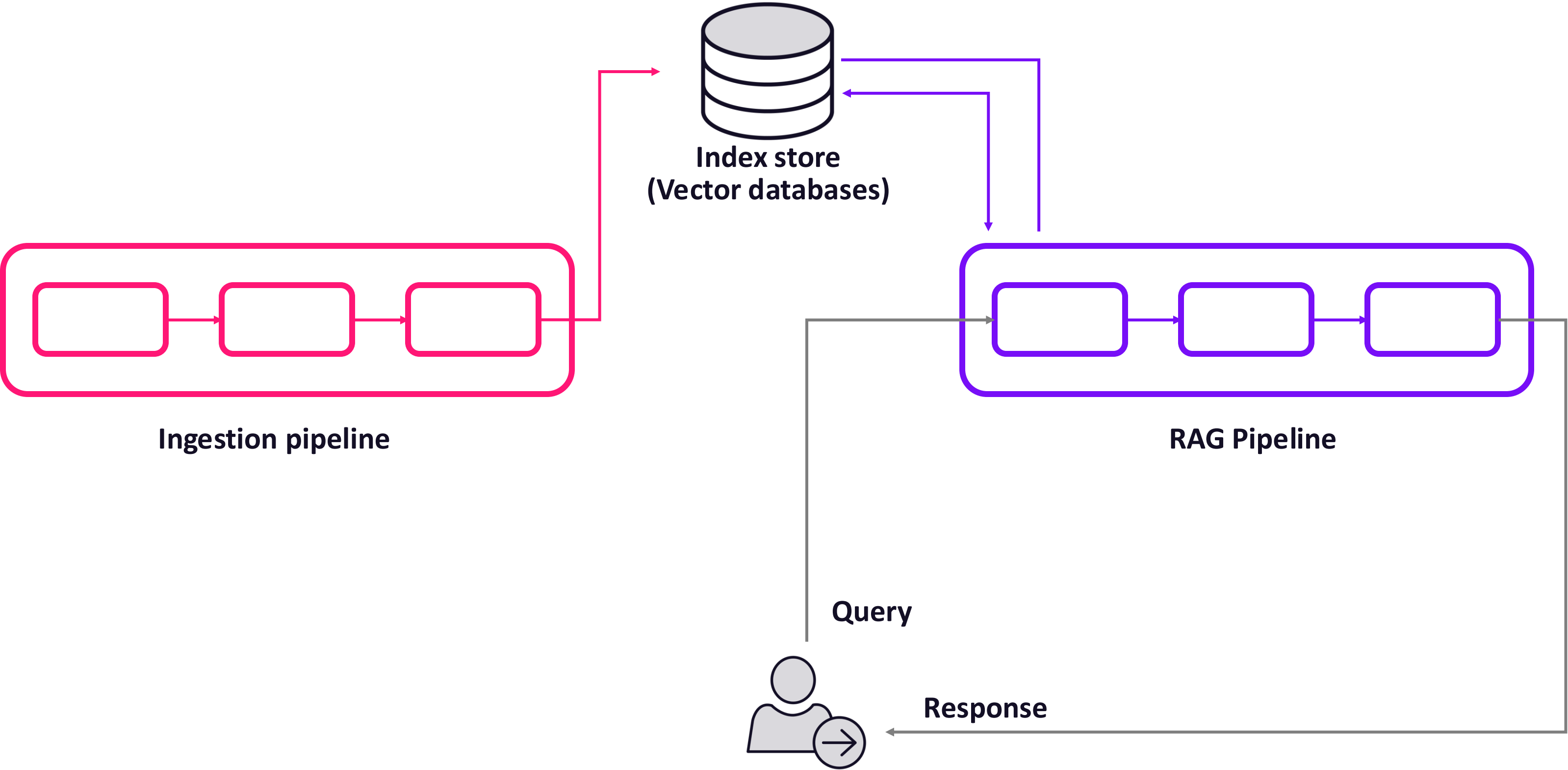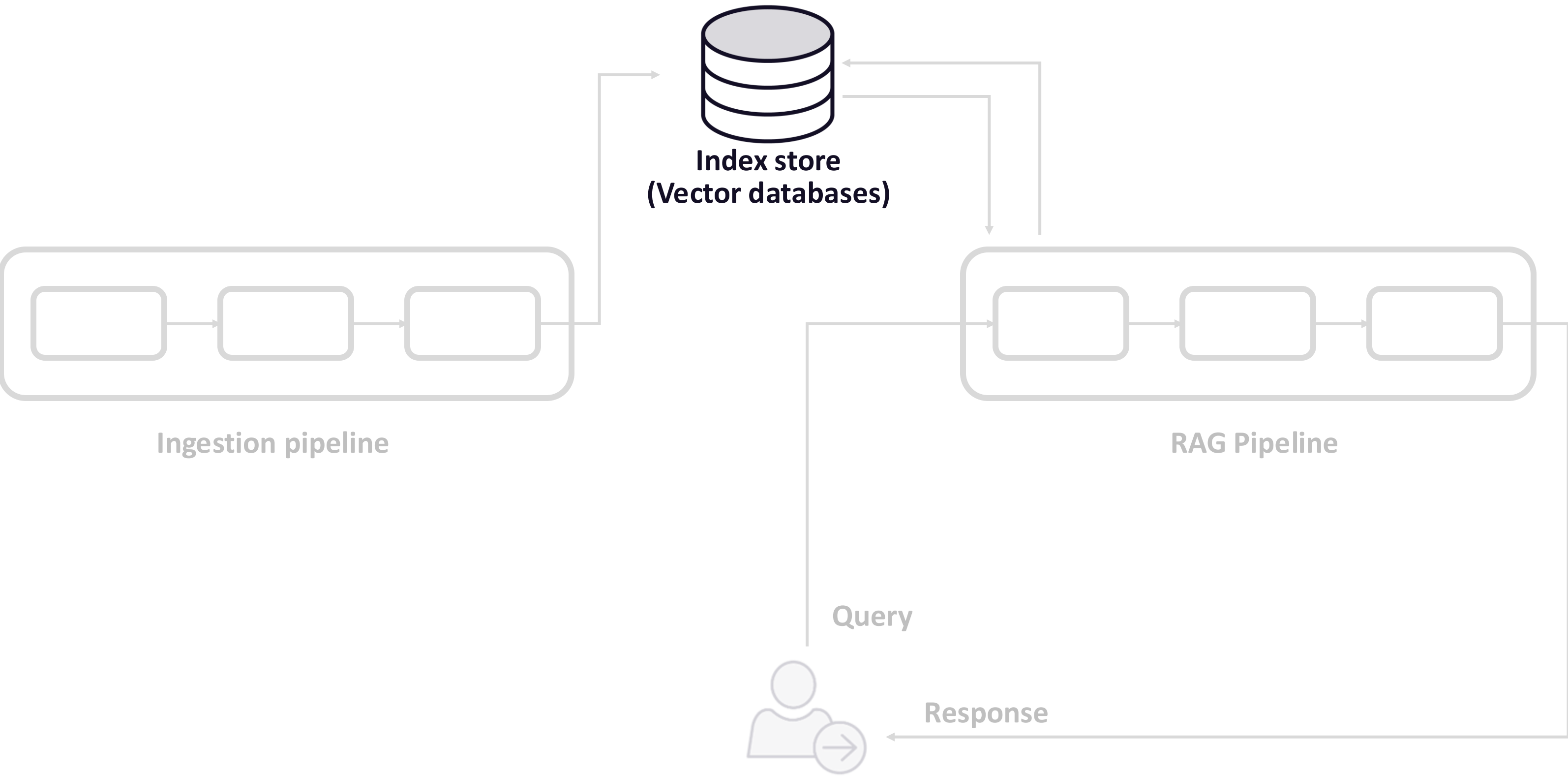
**Foo**
**adventure gear company**

**Bar**

**Qux**

# RAG System

# Deployment Blueprint - Vector Databases



**Index store
(Vector databases)**

Ingestion pipeline

RAG Pipeline

Query

Response

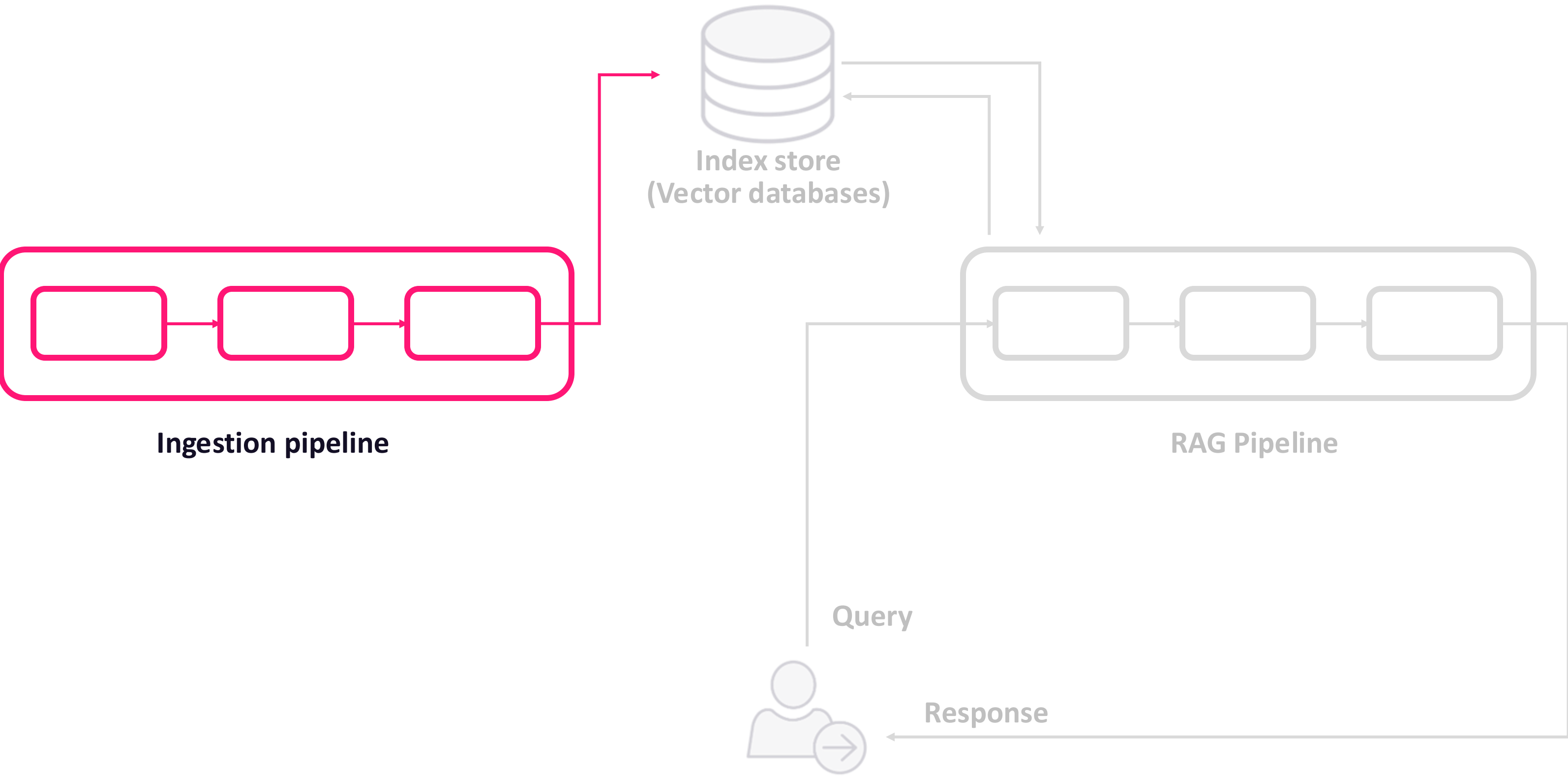# Deployment Blueprint - Vector Databases

Managed databases

- Weaviate, Pinecone

- Cloud provider's offerings such as GCP vector search

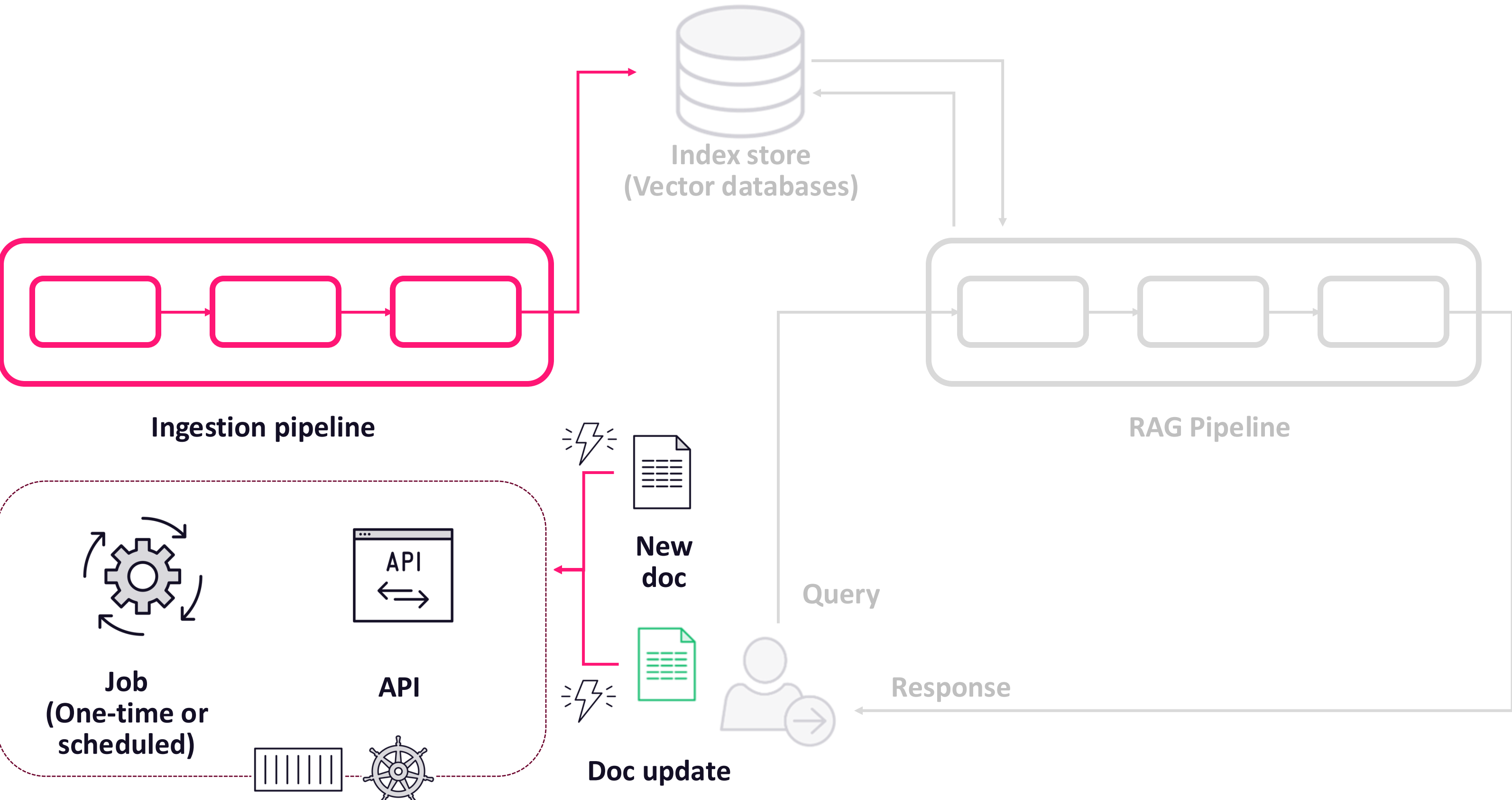Self-hosted and managed database

- Open-source database such as Chroma, Milvus, Elasticsearch

- Containers with persistent volumes

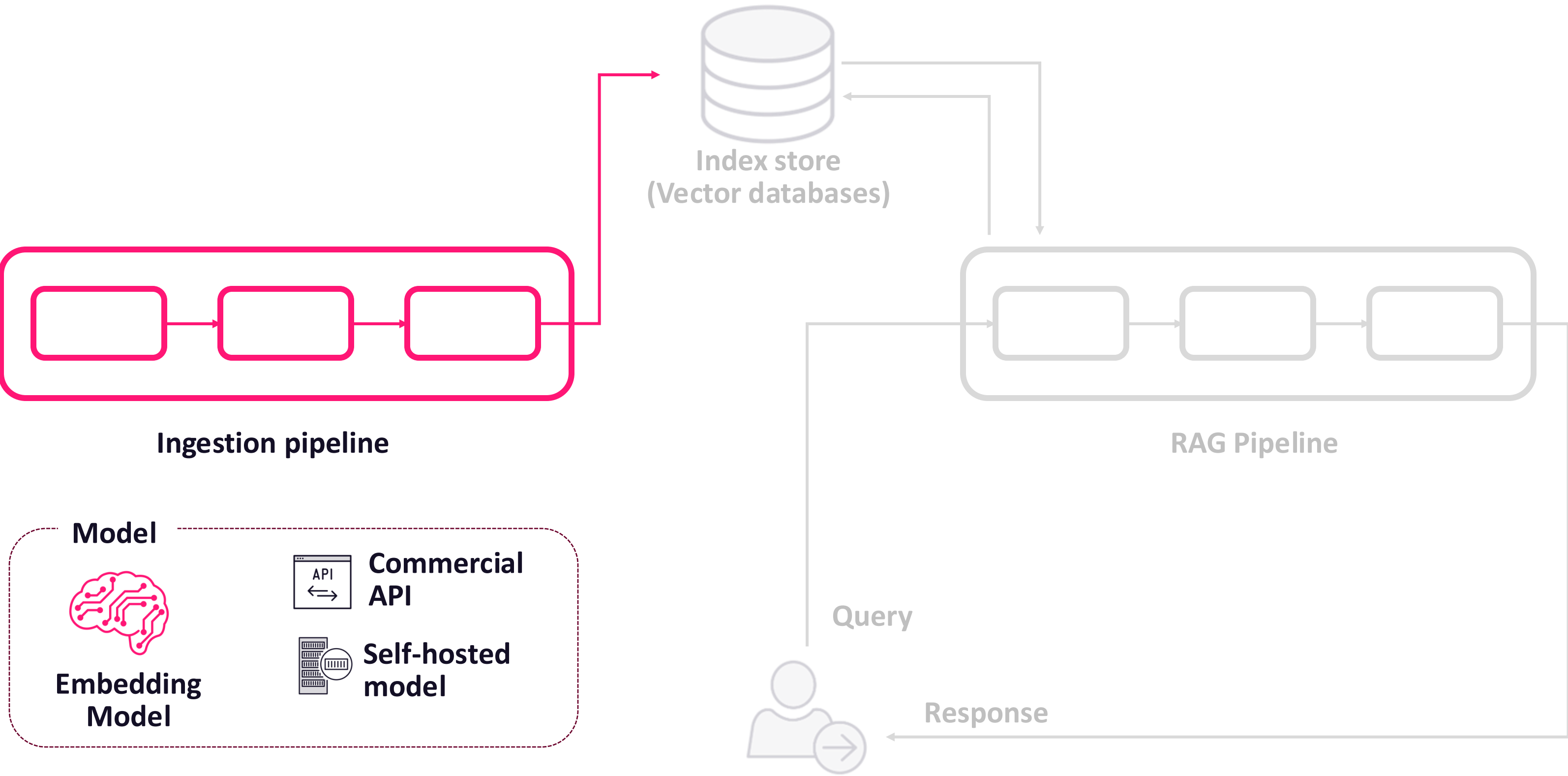- Container orchestration with Kubernetes
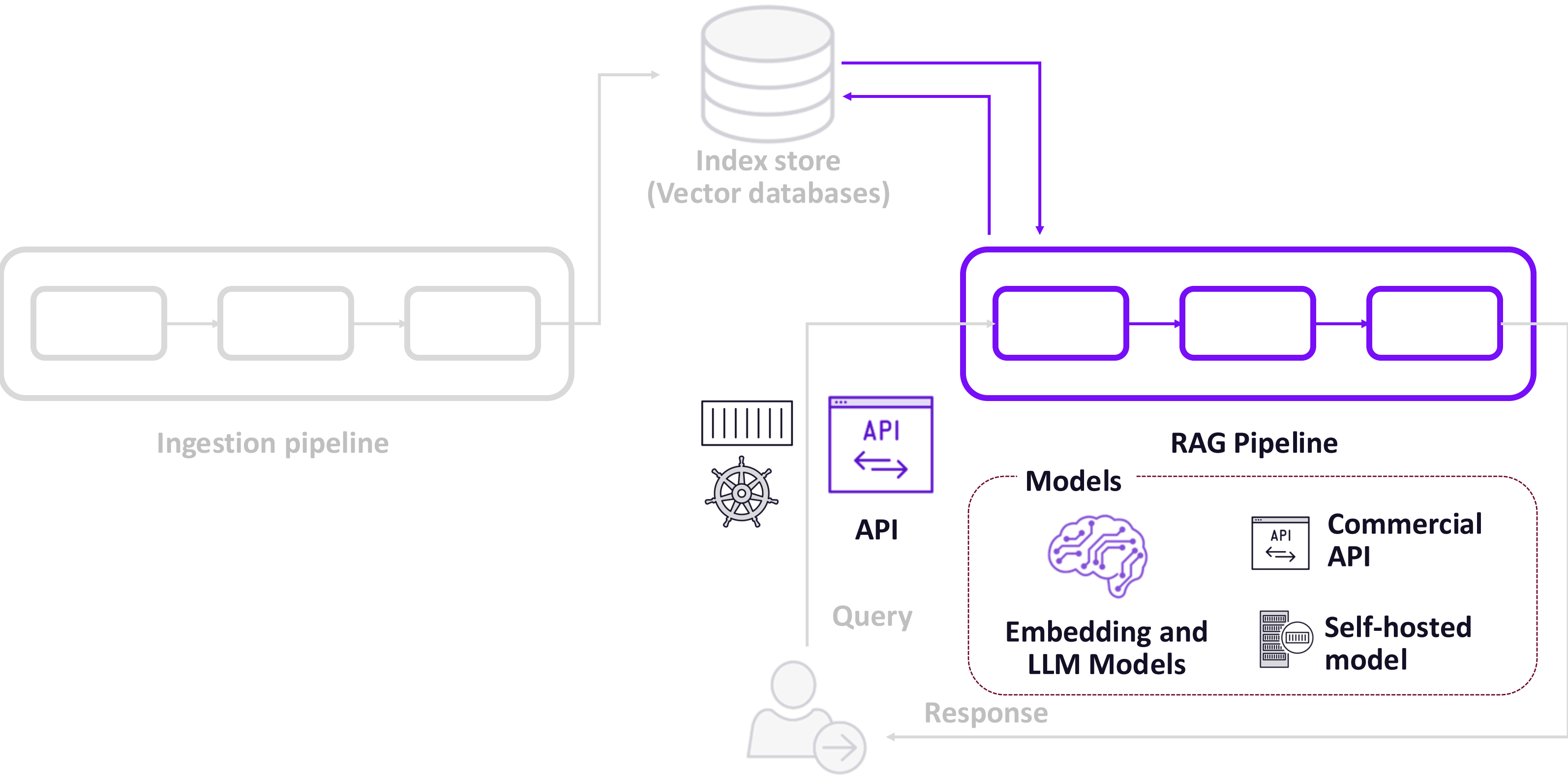
# RAG System - Ingestion Pipeline



Index store
(Vector databases)

Ingestion pipeline

RAG Pipeline

Query

Response

# Deployment Blueprint - Ingestion Pipeline



Index store
(Vector databases)

Ingestion pipeline

RAG Pipeline

Job
(One-time or
scheduled)

API

New
doc

Doc update

Query

Response

# Deployment Blueprint - Ingestion Pipeline

# RAG System - RAG Pipeline

Up Next:

# Managing RAG Systems