# Building Scalable RAG Pipelines

Enhancing AI applications with Retrieval-Augmented Generation creates more accurate, updatable systems. These pipelines power fact-aware LLM applications that respond with real-time information.

**T** **by The XYZ Company**

# What Is Retrieval-Augmented Generation (RAG)?

### Knowledge Integration

Combines information retrieval with generative AI capabilities. Creates systems that reference external knowledge before responding.
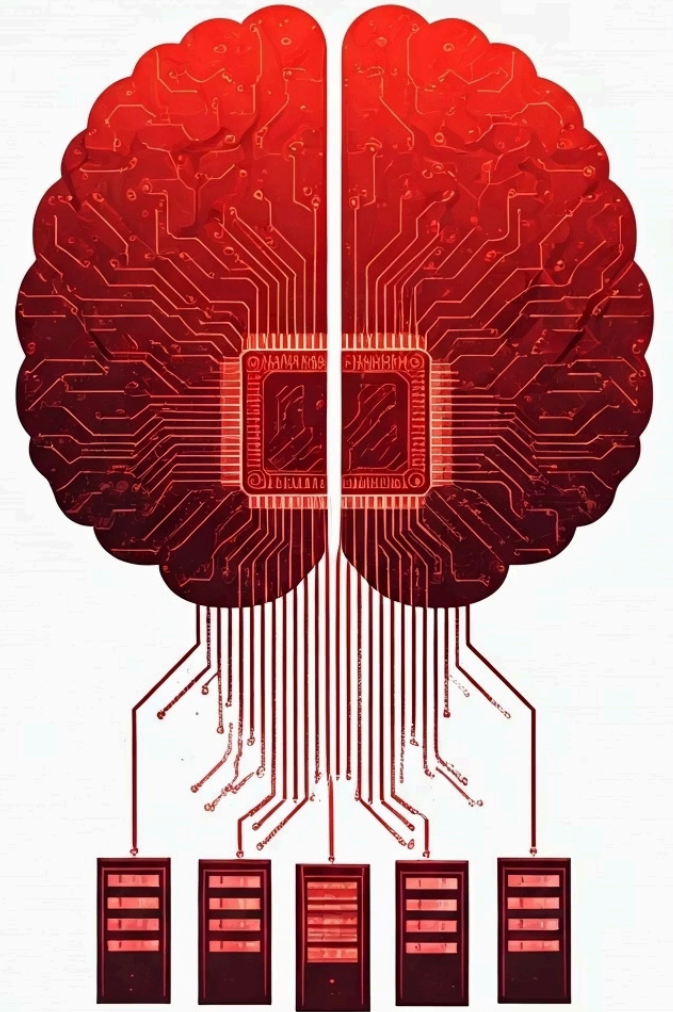
### Fact Grounding

Fetches specific facts from knowledge bases to ground LLM responses in verifiable information.
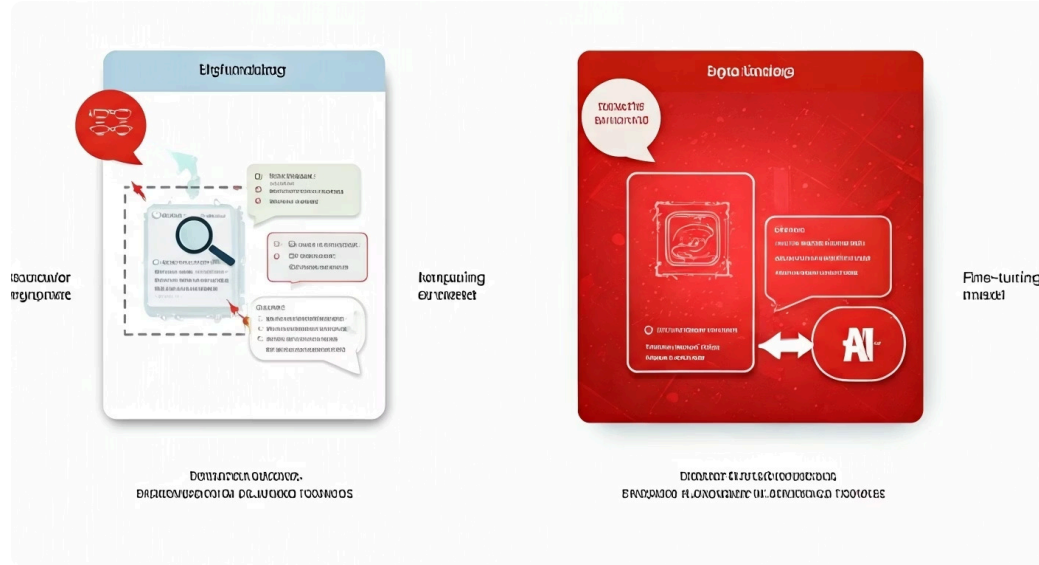
### Improved Accuracy

Minimizes hallucination by referencing actual data. Enables up-to-date responses without model retraining.

# RAG vs. Traditional Fine-Tuning



## Instant Updates

RAG refreshes knowledge through data updates. Fine-tuning needs complete model retraining.

## Cost Efficiency

Lower compute and storage requirements compared to full model retraining cycles.

## Real-Time Data

Ideal for domains with frequent information changes or need for current information.

# Key Components of a RAG Pipeline

## Document Ingestion

Importing and continuously updating various data sources into the system.

## Chunking

Breaking text into manageable units optimized for retrieval and context.

## Embedding

Converting text chunks into vector representations that capture semantic meaning.

## Retrieval

Finding relevant chunks via vector similarity to match user queries.

## Generation

Language model produces final answers using retrieved context information.

# Document Ingestion & Chunking Strategies

## Format Support

- PDFs and documents
- Web pages and HTML
- Structured records (JSON, CSV)
- APIs and databases

## Chunking Approaches

**1 Size Optimization**

Balance between comprehensive context and precise search results.

**2 Overlap Strategy**

Overlapping chunks maintain context across segment boundaries.

# Vector Databases: Core to RAG

### Embedding Storage

Store vector representations of text for semantic search capabilities. Enable dimension reduction for efficiency.
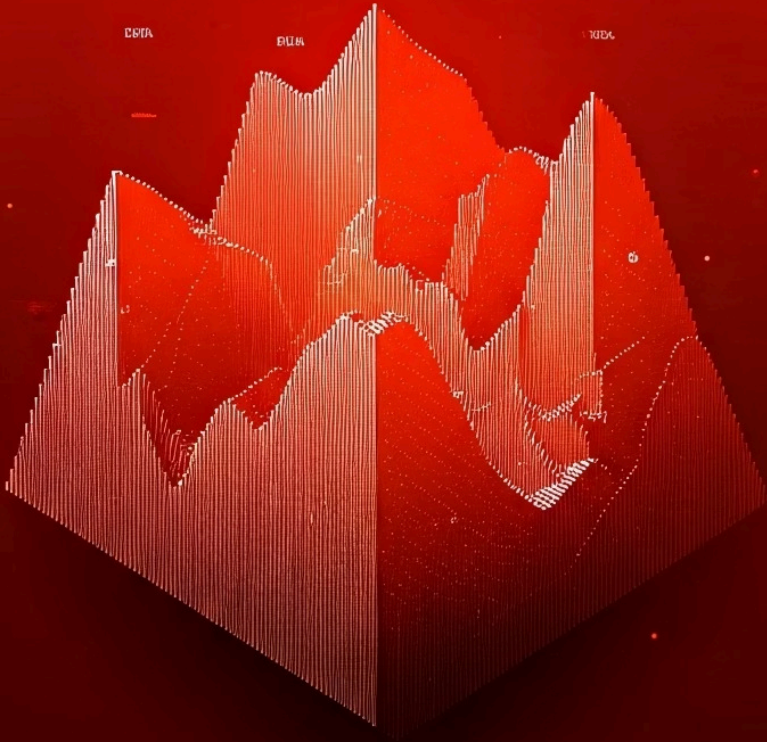
### Similarity Search

Find semantically similar content in milliseconds. Scale to millions of chunks without performance degradation.
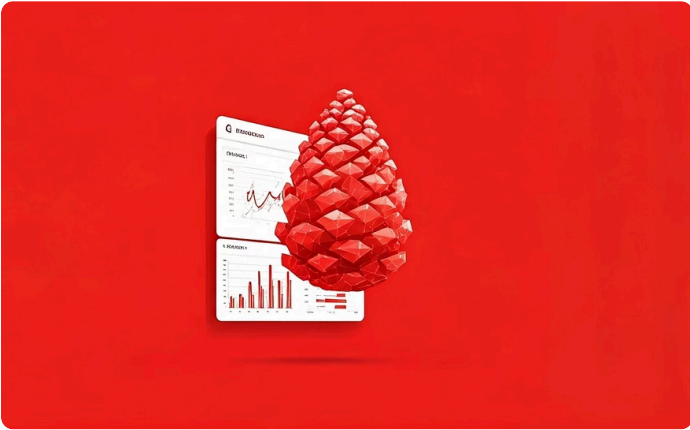
### Performance Optimization

Designed for high-throughput, low-latency AI workloads. Support for index sharding and distributed queries.

# Popular Vector Database Solutions

### Pinecone

Fully managed service optimized for production-grade RAG applications. Excels at scale with minimal operational overhead.

### Weaviate

Open-source with hybrid search capabilities. Features schema-based organization and multimodal search options.

### ChromaDB

Developer-friendly open-source option. Easily embeddable with Python-first API design for rapid prototyping.

# Summing Up: Building for Scale & Accuracy

## Key Benefits

- Up-to-date, fact-founded responses

- Lower cost than continuous fine-tuning

- Flexible knowledge updates

## Best Practices

- Automate data refresh pipelines

- Tune chunk size for your domain

- Monitor retrieval quality continuously

- Implement feedback loops for improvement

**Start Building Your RAG Pipeline**　　**Learn Advanced Techniques**