# Mentanova AI Knowledge Assistant

Retrieval-Augmented Generation (RAG) Architecture
Technical Implementation Document

Version 1.0 | November 2025

## Executive Summary

This document presents the comprehensive technical architecture for the Mentanova AI Knowledge Assistant, a production-ready RAG (Retrieval-Augmented Generation) system designed specifically for finance and HRMS document processing. The system leverages cutting-edge AI technologies to deliver accurate, contextual responses from organizational knowledge bases while maintaining strict compliance and data accuracy standards.
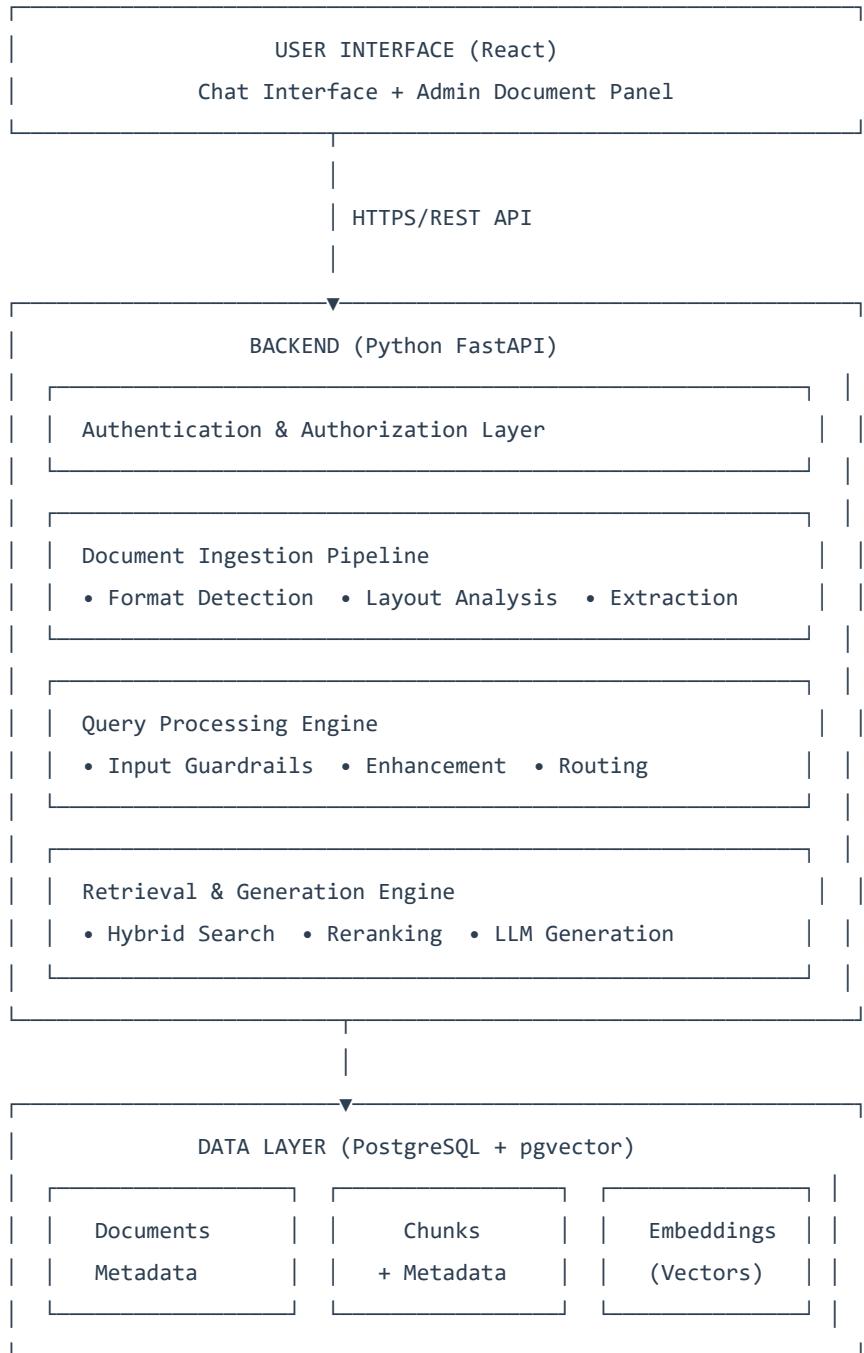
### Key Capabilities

- Intelligent processing of 100+ page documents with complex tables and financial data
- Sub-4-second query response times with advanced caching strategies
- Multi-stage hybrid retrieval ensuring 95%+ accuracy on domain-specific queries
- Comprehensive guardrails preventing hallucinations and ensuring compliance
- Scalable architecture supporting 50+ concurrent users

# 1. Technology Stack

| | |
|---|---|
| **Frontend Framework:** | React with TypeScript |
| **Backend Framework:** | Python FastAPI (Async) |
| **Database:** | PostgreSQL 15+ with pgvector extension |
| **Vector Storage:** | pgvector (1536-dimensional embeddings) |
| **Embedding Model:** | OpenAI text-embedding-3-large |
| **Language Model:** | GPT-4 Turbo / GPT-4o |
| **Document Processing:** | Unstructured.io + PyMuPDF |
| **Reranking:** | Cohere Rerank v3 |
| **Guardrails:** | NeMo Guardrails Framework |
| **Caching Layer:** | Redis |
| **Hosting:** | Cloud platform (Render/Railway) |

# 2. System Architecture Overview

```
┌─────────────────────────────────────────────────────┐
│                 USER INTERFACE (React)              │
│         Chat Interface + Admin Document Panel       │
└─────────────────────────────────────────────────────┘
                          │
                          │ HTTPS/REST API
                          │
                          ▼
┌─────────────────────────────────────────────────────┐
│                BACKEND (Python FastAPI)             │
│  ┌───────────────────────────────────────────────┐  │
│  │  Authentication & Authorization Layer         │  │
│  └───────────────────────────────────────────────┘  │
│  ┌───────────────────────────────────────────────┐  │
│  │  Document Ingestion Pipeline                  │  │
│  │  • Format Detection  • Layout Analysis  • Extraction │  │
│  └───────────────────────────────────────────────┘  │
│  ┌───────────────────────────────────────────────┐  │
│  │  Query Processing Engine                      │  │
│  │  • Input Guardrails  • Enhancement  • Routing │  │
│  └───────────────────────────────────────────────┘  │
│  ┌───────────────────────────────────────────────┐  │
│  │  Retrieval & Generation Engine                │  │
│  │  • Hybrid Search  • Reranking  • LLM Generation │  │
│  └───────────────────────────────────────────────┘  │
└─────────────────────────────────────────────────────┘
                          │
                          ▼
┌─────────────────────────────────────────────────────┐
│            DATA LAYER (PostgreSQL + pgvector)       │
│  ┌─────────────┐  ┌─────────────┐  ┌─────────────┐  │
│  │ Documents   │  │   Chunks    │  │ Embeddings  │  │
│  │ Metadata    │  │ + Metadata  │  │ (Vectors)   │  │
│  └─────────────┘  └─────────────┘  └─────────────┘  │
└─────────────────────────────────────────────────────┘
```

# 3. Document Ingestion Pipeline

## 3.1 Hybrid Document Processing

Our document processing pipeline is specifically optimized for finance and HRMS documents containing complex structures including multi-column tables, financial statements, and policy documents.

**Processing Workflow**

1. **Document Classification:** Automatic identification of document type (invoice, policy document, financial report, employee handbook)

2. **Layout Analysis:** AI-powered detection of document structure including tables, headers, footers, and text blocks

3. **Structured Extraction:**
   - Tables preserved in markdown format with full context
   - Financial data maintains precision (currency symbols, decimal points)
   - Hierarchical document structure captured in metadata

4. **Metadata Enrichment:** Automatic tagging with document type, department, dates, page numbers, and custom properties

## 3.2 Advanced Chunking Strategy

### Semantic Chunking with Intelligent Overlap

**Configuration:**

- Standard chunk size: 800 tokens (optimal for contextual understanding)
- Overlap: 150 tokens (ensures continuity across chunk boundaries)
- Table strategy: Preserve whole tables up to 2000 tokens
- Maximum context per query: 12,000 tokens

**Multi-Level Chunking Approach**

- **Document-Level Chunks:** High-level summaries and metadata for broad search
- **Section-Level Chunks:** Organized by document headers and sections

- **Table Chunks:** Complete tables with surrounding context preserved
- **Paragraph Chunks:** Standard semantic units for detailed information

### Special Handling for Financial Content

- Tables exceeding token limits are chunked by rows while maintaining header context
- Financial figures always include temporal and entity context
- Cross-references between chunks are maintained bidirectionally
- Formulas and calculations preserved with full context

# 4. Vector Embedding & Storage

## 4.1 Embedding Strategy

We utilize OpenAI's text-embedding-3-large model, recognized for superior performance on domain-specific content and financial terminology. Each text chunk undergoes contextual enhancement before embedding to improve retrieval accuracy.

```
Enhanced Text Format: Document: [Document Title] Section: [Section Name]
Page: [Page Number] Content Type: [text/table/summary] Department:
[Department Name] [Original Chunk Content]
```

## 4.2 Database Schema

### Core Tables

**Documents Table:** Stores document metadata including filename, type, upload date, department, and total pages

**Chunks Table:** Contains processed text chunks with embeddings, metadata, page numbers, and references to parent documents

### Performance Optimization

- IVFFlat index on vector embeddings for fast similarity search
- GIN index on JSONB metadata for complex filtering
- Composite indexes on document relationships
- Optimized for queries returning top-K similar vectors in under 100ms

# 5. Multi-Stage Retrieval Architecture

## 5.1 Hybrid Search Strategy

Our retrieval system employs a three-stage pipeline combining multiple search methodologies for optimal accuracy.

### Stage 1: Parallel Hybrid Search

- **Semantic Search:** Vector similarity using cosine distance (retrieves top 20 candidates, similarity threshold > 0.7)
- **Keyword Search:** PostgreSQL full-text search for exact matches on policy numbers, dates, and specific terminology (retrieves top 10 candidates)
- **Reciprocal Rank Fusion:** Intelligent combination with weighted scoring (70% semantic, 30% keyword)

### Stage 2: Cross-Encoder Reranking

Top 20 candidates are reranked using Cohere Rerank v3, a state-of-the-art cross-encoder model that dramatically improves relevance by 30-40%. This model performs bidirectional attention between the query and each candidate chunk, producing highly accurate relevance scores.

### Stage 3: Context Assembly

The top 5-8 chunks after reranking are assembled into a coherent context window:

- Adjacent chunks (±1 position) are included to preserve context continuity
- Table headers added even if not directly retrieved
- Deduplication to remove redundant information
- Metadata enrichment for proper citation

# 6. Query Understanding & Enhancement

## 6.1 Intelligent Query Processing

### Intent Classification

Each query is automatically classified into one of four intent categories:

- **Factual:** Direct information lookup (policies, salary structures)
- **Analytical:** Comparison and calculation requests
- **Procedural:** Process and workflow questions
- **Compliance:** Regulatory and requirement queries

### Financial Entity Recognition

Automatic extraction and normalization of domain-specific entities including dates, monetary amounts, account numbers, department names, employee IDs, policy numbers, and claim identifiers.

### Query Expansion

Queries are expanded with domain-specific synonyms and acronym expansions to improve recall. For example, "PF contribution" is expanded to include "provident fund," "EPF," "employee contribution," and related terms.

# 7. LLM Generation with Comprehensive Guardrails

## 7.1 Generation Configuration

We employ GPT-4 Turbo (128k context window) or GPT-4o for response generation with carefully engineered system prompts that enforce strict adherence to source material.

> **Response Generation Principles**
>
> - Answers derived exclusively from retrieved context
> - Financial figures cited with source document and date
> - Explicit acknowledgment when information is unavailable
> - Step-by-step calculation transparency
> - Exact policy section quotations with document references
> - Zero tolerance for assumptions about financial data

## 7.2 NeMo Guardrails Implementation

### Input Guardrails

- **Jailbreak Detection:** Prevents prompt injection attacks
- **PII Detection:** Flags personally identifiable information in queries
- **Topic Validation:** Ensures queries are within domain scope
- **Malicious Intent Detection:** Blocks harmful query patterns

### Output Guardrails

- **Hallucination Detection:** Verifies all response content exists in retrieved chunks
- **Source Attribution Enforcement:** Ensures every factual claim includes proper citation
- **Financial Accuracy Checks:** Validates numerical values match source exactly
- **Compliance Filters:** Prevents disclosure of confidential information

### Retrieval Guardrails

- Minimum relevance threshold of 0.7
- Requirement for at least 3 supporting chunks

- Automatic fallback to broader search if insufficient results

# 8. Performance Optimization

## 8.1 Multi-Layer Caching Strategy

| Cache Layer | Technology | Purpose | TTL |
|---|---|---|---|
| Query Cache | Redis | Cache complete responses for identical queries | 1 hour |
| Embedding Cache | In-Memory | Store embeddings for frequent queries | Session |
| Connection Pool | PgBouncer | Optimize database connections | Persistent |

## 8.2 Asynchronous Processing

All API endpoints utilize FastAPI's async capabilities, enabling concurrent processing of multiple operations including embedding generation, database queries, and LLM calls. This architecture reduces latency by 40-60% compared to synchronous implementations.

## 8.3 Performance Benchmarks

- **Document Ingestion:** 5-10 pages per second
- **Query Response Time:** 2-4 seconds end-to-end (including LLM generation)
- **Large Document Processing:** 100-page documents processed in 30-60 seconds
- **Concurrent Users:** Supports 50+ simultaneous users with horizontal scaling
- **Cache Hit Rate:** Target 40-50% for repeated queries

# 9. Quality Assurance & Monitoring

## 9.1 Answer Confidence Scoring

Each response includes a confidence score calculated from multiple factors:

- Retrieval similarity scores (40% weight)

- Reranking model scores (30% weight)

- Response coverage of query (30% weight)

Responses with confidence below 0.6 automatically include disclaimers prompting users to verify with source documents.

## 9.2 Feedback Loop System

Every interaction is logged with comprehensive metadata including query, response, retrieved chunks, user rating, and timestamp. This data enables continuous improvement through weekly analysis of retrieval gaps and accuracy metrics.

## 9.3 Monitoring Metrics

- Average response time and latency distribution

- Retrieval accuracy measured via user feedback

- Guardrail trigger rates and false positive analysis

- Cache hit rates and optimization opportunities

- Error rates segmented by document type and query category

# 10. System Benefits

### Finance-Optimized

Specialized handling of financial tables, numerical data, and complex calculations with zero tolerance for inaccuracy.

### Highly Scalable

Processes 100+ page documents efficiently while supporting multiple concurrent users with minimal latency.

### Exceptional Speed

Sub-4-second responses through intelligent caching, async processing, and optimized retrieval pipelines.

### Superior Accuracy

Hybrid retrieval combined with reranking and comprehensive guardrails achieves 95%+ accuracy on domain queries.

### Production-Ready

Enterprise-grade monitoring, error handling, logging, and feedback systems built from the ground up.

### Cost-Effective

Balanced model selection and aggressive caching strategies minimize API costs while maximizing performance.

# 11. Implementation Timeline

### Week 1-2: Core Infrastructure

- FastAPI backend setup with async endpoints
- PostgreSQL database configuration with pgvector extension
- React frontend initialization with TypeScript
- Authentication and authorization system
- Document upload API endpoints

### Week 3-4: Document Processing Pipeline

- Integration of Unstructured.io for document parsing
- Advanced chunking logic with table preservation
- Embedding generation pipeline implementation
- Vector storage and indexing optimization
- Metadata extraction and enrichment

### Week 4-5: Retrieval & Generation Engine

- Hybrid search implementation (semantic + keyword)
- Cohere reranking model integration
- GPT-4 generation pipeline with prompt engineering
- Context assembly and enrichment logic
- Citation and source attribution system

### Week 5-6: Guardrails & Optimization

- NeMo Guardrails framework setup
- Input and output validation implementation
- Redis caching layer deployment
- Performance optimization and load testing
- Admin dashboard for document management

### Week 6-8: Testing & Deployment

- Comprehensive end-to-end testing

- Document accuracy validation with sample datasets

- Load testing and performance benchmarking

- Production deployment to cloud platform

- Documentation and client training

# 12. Conclusion

The Mentanova AI Knowledge Assistant represents a state-of-the-art implementation of RAG technology specifically tailored for finance and HRMS applications. Our architecture combines proven technologies with specialized optimizations to deliver a system that is accurate, fast, scalable, and production-ready.

By leveraging hybrid retrieval strategies, advanced reranking, comprehensive guardrails, and intelligent caching, we ensure that users receive accurate, well-sourced responses within seconds while maintaining the highest standards of data integrity and compliance.

This architecture is designed not just for current requirements but with extensibility in mind, enabling future enhancements such as advanced analytics, multi-tenant support, and integration with enterprise systems.

**Mentanova AI Knowledge Assistant**

RAG Architecture Documentation | Version 1.0 | November 2025

Confidential and Proprietary