

AI Research Assistant: Detailed Technical Report

Dasari Santhan Reddy

December 3, 2025

Contents

1	Project Overview	4
2	System Architecture	4
2.1	Frontend Architecture	4
2.2	Backend Architecture	4
2.3	AI/ML Components	5
3	Detailed Component Breakdown	5
3.1	Document Processing Pipeline	5
3.1.1	File Upload Handler	5
3.1.2	Text Extraction	5
3.1.3	Content Processing	6
3.2	Search System	6
3.2.1	Full-Text Search	6
3.2.2	Semantic Search	6
3.3	User Management System	6
3.3.1	Authentication	6
3.3.2	Authorization	6
4	Database Schema	6
4.1	Users Table	7
4.2	Documents Table	7
4.3	Document Sections Table	7
5	API Endpoints	7
5.1	Authentication APIs	7
5.2	Document APIs	7

5.3	Search APIs	7
6	AI Integration	7
6.1	Document Analysis	7
6.2	Question Answering	8
7	Deployment	8
7.1	Containerization	8
7.2	Cloud Deployment	8
7.3	CI/CD Pipeline	8
8	Security Measures	8
8.1	Data Security	8
8.2	Application Security	8
9	Performance Optimization	9
9.1	Caching	9
9.2	Database Optimization	9
10	Monitoring and Logging	9
10.1	Monitoring	9
10.2	Logging	9
11	Testing Strategy	9
11.1	Unit Testing	9
11.2	Integration Testing	9
12	Future Enhancements	9
12.1	AI Features	9
12.2	Collaboration	10
13	Conclusion	10

Abstract

The AI Research Assistant is an intelligent research support platform developed to automate document processing, semantic search, and question answering using advanced Artificial Intelligence techniques. The system integrates a scalable web architecture with Large Language Models and vector-based search to enhance academic research efficiency. This report presents a comprehensive technical description of system architecture, implementation, security, performance optimization, testing strategy, and deployment. The proposed system demonstrates how modern AI can significantly reduce manual research effort while improving accuracy, accessibility, and collaboration.

1 Project Overview

The AI Research Assistant is a full-stack, AI-driven platform developed to streamline the research process by providing automated document ingestion, semantic retrieval, summarization, and intelligent question answering. Traditional research workflows often involve manual searching, repetitive reading, and lack of contextual retrieval. This system addresses these challenges by integrating Natural Language Processing (NLP), vector databases, and Large Language Models (LLMs).

The platform allows users to upload research documents in multiple formats, automatically extract and index their content, and perform semantic searches across large document collections. AI-powered reasoning enables users to ask natural language questions and obtain precise, context-aware answers with proper source citations. This drastically reduces the time required for literature review and knowledge extraction.

2 System Architecture

The system follows a three-tier modular architecture composed of:

1. Frontend (Presentation Layer)
2. Backend (Application Layer)
3. AI/ML Processing Layer

This separation of concerns ensures independent scaling, easy maintenance, and technology upgrades without affecting other components.

2.1 Frontend Architecture

The frontend is built using React.js (v18+), providing a fast, responsive, and interactive user interface. The component-based design improves reusability and maintainability. State management is handled using React Context API with useReducer for predictable state transitions. Axios is used to communicate securely with backend REST APIs.

JWT-based authentication ensures secure access to protected resources. React Router v6 enables seamless navigation between pages. The frontend architecture supports dynamic dashboards, document upload interfaces, advanced search panels, and real-time AI response visualization.

2.2 Backend Architecture

The backend is implemented using FastAPI (Python 3.9+) due to its asynchronous support, high performance, and auto-generated API documentation using OpenAPI.

SQLAlchemy ORM is used for database interactions, ensuring database abstraction and secure queries.

PostgreSQL is used as the production database for its reliability and support for JSON and full-text search. SQLite is utilized during development for lightweight testing. The backend handles authentication, authorization, document management, request routing, AI orchestration, and database transactions.

2.3 AI/ML Components

The AI layer forms the core intelligence of the system. LangChain orchestrates interactions between multiple LLMs and tools. OpenAI's GPT models perform reasoning, answering, summarization, and text generation. LlamaIndex structures documents into vector indexes for efficient retrieval.

FAISS or HNSW is used as a high-performance vector database supporting approximate nearest neighbor (ANN) search for fast semantic similarity queries. Custom pre-processing pipelines ensure consistent text normalization and chunk-level embedding.

3 Detailed Component Breakdown

3.1 Document Processing Pipeline

The document processing pipeline transforms raw documents into structured, searchable knowledge units.

3.1.1 File Upload Handler

The file upload handler supports PDF, DOCX, and TXT formats. It performs validation on file type, size, and MIME format to prevent malicious uploads. Files are securely stored in object storage with unique identifiers.

3.1.2 Text Extraction

Text extraction converts uploaded files into raw text:

- PyPDF is used for PDF parsing.
- python-docx is used for Word document extraction.
- BeautifulSoup processes embedded HTML content.

The extracted text is cleaned to remove headers, footers, and irrelevant characters.

3.1.3 Content Processing

Extracted text undergoes normalization, tokenization, language detection, and segmentation into chunks. Chunking is essential for efficient embedding and AI context window management. Each chunk is stored for semantic indexing.

3.2 Search System

The search system supports both exact keyword-based and semantic context-based retrieval.

3.2.1 Full-Text Search

PostgreSQL full-text search provides rapid lexical matching with ranking and fuzzy matching, suitable for exact technical queries.

3.2.2 Semantic Search

Semantic search uses embedding vectors and ANN algorithms to retrieve conceptually similar text, even if the exact keywords are not matched. Hybrid ranking combines both lexical and semantic scores.

3.3 User Management System

The system implements secure identity and access control mechanisms.

3.3.1 Authentication

User authentication supports email/password login with bcrypt hashing and OAuth2 login through Google and GitHub. JWT access and refresh tokens ensure secure, stateless sessions.

3.3.2 Authorization

Role-Based Access Control (RBAC) restricts access to sensitive resources. Permission management supports personal workspaces, shared teams, and collaborative research projects.

4 Database Schema

The database schema follows a normalized relational design to ensure integrity, scalability, and performance.

4.1 Users Table

Stores authentication and profile information. UUIDs provide secure unique identification.

4.2 Documents Table

Stores metadata and raw extracted content of uploaded documents. JSONB allows flexible metadata storage.

4.3 Document Sections Table

Stores semantic chunks with vector embeddings enabling fast ANN-based retrieval.

5 API Endpoints

RESTful API design is followed for stateless client-server communication.

5.1 Authentication APIs

User registration, login, token refresh, and logout endpoints ensure secure identity management.

5.2 Document APIs

Provide CRUD operations for document upload, retrieval, deletion, and internal document search.

5.3 Search APIs

Support global search, semantic query processing, and real-time suggestion generation.

6 AI Integration

6.1 Document Analysis

AI automatically performs summarization, keyword extraction, Named Entity Recognition, and citation identification. These features enhance document understanding and organization.

6.2 Question Answering

Context-aware question answering is implemented using retrieval-augmented generation (RAG). The system fetches relevant document chunks and provides grounded answers with source attribution.

7 Deployment

7.1 Containerization

Docker is used for containerizing frontend and backend services with multi-stage optimized builds. Environment variables securely manage secrets.

7.2 Cloud Deployment

Vercel hosts the frontend. Backend services run on AWS or GCP with auto-scaling. Managed PostgreSQL ensures high availability. Object storage manages document files securely.

7.3 CI/CD Pipeline

GitHub Actions enable automated testing, linting, container builds, and staged deployments across development, staging, and production environments.

8 Security Measures

8.1 Data Security

AES-256 encryption protects data at rest. TLS 1.3 secures data in motion. Automatic backups and retention policies ensure data availability and compliance.

8.2 Application Security

Rate limiting prevents abuse. CORS policies protect cross-origin attacks. CSRF tokens, input validation, and HTTP security headers defend against common vulnerabilities.

9 Performance Optimization

9.1 Caching

Redis is used to cache frequently accessed queries, embeddings, and session data. HTTP headers optimize browser-side caching.

9.2 Database Optimization

Indexing, query optimization, connection pooling, and read replicas ensure efficient database performance under high workloads.

10 Monitoring and Logging

10.1 Monitoring

System metrics, error rates, API latency, and uptime are continuously monitored using cloud-based monitoring tools.

10.2 Logging

Structured logs are generated at multiple log levels and centralized for debugging and audit purposes.

11 Testing Strategy

11.1 Unit Testing

Unit tests ensure correctness of individual modules with more than 80% code coverage.

11.2 Integration Testing

End-to-end testing validates correctness of API communication, database operations, AI inference, and user workflows.

12 Future Enhancements

12.1 AI Features

Future upgrades include automatic research paper drafting, citation formatting, plagiarism detection, and multilingual document processing.

12.2 Collaboration

Real-time multi-user editing, threaded annotations, integrated version control, and team-based workspaces will enhance collaboration.

13 Conclusion

The AI Research Assistant demonstrates the potential of Artificial Intelligence in transforming traditional research workflows. By integrating semantic search, intelligent reasoning, and scalable cloud architecture, the system significantly improves research efficiency, accuracy, and collaboration. Its modular design allows future expansion into advanced academic automation tools.