

JARVIS: Just A Rather Very Intelligent System

Interactive AI Assistant Project Documentation



Table of Contents

- [Overview](#)
- [Technical Stack](#)
- [Architecture](#)
- [Frontend Components](#)
- [Backend Services](#)
- [RAG Implementation](#)
- [API Endpoints](#)
- [Development Process with CursorAI](#)
- [Setup and Deployment](#)
- [Future Enhancements](#)
- [Troubleshooting](#)

Overview

JARVIS (Just A Rather Very Intelligent System) is an interactive AI assistant interface inspired by the AI from Iron Man. This project showcases an immersive, futuristic UI that allows users to interact with a RAG-powered (Retrieval Augmented Generation) AI assistant that can answer questions about Timal Pathirana's professional background, skills, and projects.

The system features:

- A holographic-inspired user interface with sci-fi elements
- RAG-powered question answering using vector search
- Context-aware conversation capability
- Animated interface with bootup sequence
- Interactive panels with detailed information
- Real-time AI chat functionality

Technical Stack

Frontend

- **React 18:** Core UI framework
- **TypeScript:** For type-safe code
- **Framer Motion:** For fluid animations and transitions
- **Tailwind CSS:** For responsive styling
- **React Icons:** For iconography
- **Vite:** For fast development and bundling

Backend

- **FastAPI:** High-performance Python web framework
- **LangChain:** For connecting LLMs with vector databases
- **Pinecone:** Vector database for RAG functionality
- **OpenAI API:** For powering the LLM capabilities
- **Python 3.10+:** Backend language

DevOps & Tools

- **CursorAI:** For AI-assisted development
- **GitHub Actions:** CI/CD pipeline
- **Vercel:** Frontend deployment
- **Docker:** For containerization (optional)

Architecture

JARVIS employs a clean, separated architecture:

1. **Frontend Layer:** React application with TypeScript

- UI Components
- Animation logic
- API service interfaces

2. **API Layer:** FastAPI backend

- REST endpoints
- Request validation
- Response formatting

3. **LLM Processing Layer:**

- Query processing
- Context retrieval
- Response generation

4. **Vector Storage Layer:**

- Document storage
- Embedding calculation
- Similarity search

5. **Document Processing Layer:**

- PDF parsing
- Text chunking
- Metadata extraction

Frontend Components

The frontend is structured around reusable components:

JarvisUI (Main Component)

The core UI component that orchestrates the entire interface experience, including:

- Bootup sequence animation
- Welcome message and personalization
- Holographic panels with information
- Chat interface

Key UI Features

1. **Bootup Sequence:** Simulates system initialization with progress bars and tech-themed messages
2. **Information Panels:** Interactive panels for:
 - Profile
 - Education
 - Recommendations
 - Career
 - Projects
 - Skills
3. **JARVIS Chat Interface:**
 - Two-column layout with suggested questions
 - Real-time chat with the AI
 - Visual indicators for "thinking" state
 - Animated audio visualization
4. **Animations:**
 - Motion effects for all interactive elements
 - Particle background for depth
 - Scan lines for sci-fi aesthetic
 - Pulse effects for key interactive elements

Backend Services

The backend is organized into several specialized services:

API Service (api.py)

- Handles HTTP requests and responses
- Implements CORS for cross-origin requests
- Provides endpoints for document upload and queries
- Manages error handling and response formatting

LLM Manager (llm.py)

- Integrates with OpenAI API
- Manages conversation context
- Formats prompts with retrieved information
- Handles fallback strategies when RAG fails

Vector Store Manager (vector_store.py)

- Interfaces with Pinecone vector database
- Manages document embeddings

- Performs similarity search for RAG
- Handles index creation and maintenance

Document Processor (document_processor.py)

- Parses PDF documents
- Chunks text for effective embedding
- Extracts metadata for better retrieval
- Preprocesses text for optimal embedding

RAG Implementation

The Retrieval Augmented Generation (RAG) system works as follows:

1. Document Ingestion:

- PDF documents containing information about Timal's background are uploaded
- Documents are parsed and split into semantic chunks
- Each chunk is embedded using the OpenAI embeddings model
- Embeddings are stored in Pinecone vector database

2. Query Processing:

- User query is received through the frontend
- Query is embedded using the same embedding model
- Vector similarity search finds relevant chunks
- Retrieved chunks form the context for the LLM

3. Response Generation:

- System prompt defines JARVIS personality and scope
- Retrieved context is added to the prompt
- Conversation history is included for continuity
- OpenAI generates a contextually relevant response
- Sources are tracked for citation

4. Fallback Mechanism:

- If RAG retrieval fails, falls back to base LLM
- Error handling preserves user experience
- Backend logs capture issues for debugging

API Endpoints

Frontend API Service

The frontend communicates with the backend through a dedicated service:

```
// Key API functions
queryJarvis(query: string, conversationHistory: Array<{role: string, content: string}>)
uploadDocument(file: File)
checkApiStatus()
clearVectorStore()
```

Backend Endpoints

1. **GET /** - Health check endpoint
 - Returns API status
2. **POST /upload** - Document upload
 - Accepts PDF files
 - Returns document processing results
3. **POST /query** - RAG query endpoint
 - Accepts user query and conversation history
 - Returns AI response with source citations
4. **POST /clear** - Admin endpoint
 - Clears the vector database
 - For development use only

Development Process with CursorAI

CursorAI played a pivotal role in developing JARVIS:

Initial Setup and Scaffolding

- CursorAI generated the initial project structure
- Created component skeletons based on requirements
- Suggested appropriate libraries and dependencies

UI Development

- Generated complex Framer Motion animations
- Helped design the holographic interface elements
- Created responsive layouts with Tailwind CSS
- Implemented particle effects and animated transitions

Backend Integration

- Built FastAPI endpoints with proper validation
- Implemented RAG functionality with LangChain
- Set up Pinecone vector store integration
- Created document processing pipeline

Code Refinement

- Identified and fixed type issues in TypeScript
- Optimized React component rendering
- Improved error handling and fallback strategies
- Enhanced code documentation

Testing and Debugging

- Generated test cases for API endpoints
- Diagnosed and fixed integration issues
- Validated RAG functionality
- Ensured cross-browser compatibility

Setup and Deployment

Prerequisites

- Node.js 16+
- Python 3.10+
- OpenAI API key
- Pinecone API key

Frontend Setup

1. Install dependencies:

```
npm install
```

2. Configure environment:
Create .env file with:

```
VITE_API_BASE_URL=http://localhost:8000
```

3. Run development server:

```
npm run dev
```

Backend Setup

1. Install dependencies:

```
pip install -r backend/requirements.txt
```

2. Configure environment:
Create .env file based on env.example:

```
OPENAI_API_KEY=your_openai_key  
PINECONE_API_KEY=your_pinecone_key  
PINECONE_ENVIRONMENT=us-west4-gcp
```

3. Run development server:

```
cd backend  
uvicorn api:app --reload
```

Deployment

1. **Frontend:**
 - Deploy to Vercel or similar service
 - Configure production environment variables

2. Backend:

- Deploy to a cloud provider (AWS, GCP, Azure)
- Set up environment variables
- Configure CORS for production domain

Future Enhancements

Potential improvements for JARVIS include:

1. Voice Interface:

- Speech-to-text for voice commands
- Text-to-speech for JARVIS responses
- Voice signature matching

2. Enhanced Visualization:

- 3D holographic elements with Three.js
- Data visualization for career progression
- Interactive skill tree

3. Extended RAG Capabilities:

- Multi-modal document processing (images, videos)
- Real-time information from web APIs
- Structured data integration (SQL databases)

4. Performance Optimizations:

- Lazy loading of UI components
- Optimized vector retrieval
- Caching of common queries

Troubleshooting

Common Issues and Solutions

1. API Connection Errors:

- Ensure backend server is running
- Check CORS configuration
- Verify API base URL in frontend config

2. RAG Not Working:

- Confirm vector database has documents
- Check embedding model access
- Validate API keys

3. UI Animation Performance:

- Reduce animation complexity on lower-end devices
- Implement conditional rendering for effects
- Use React.memo for expensive components

4. PDF Processing Issues:

- Ensure PDF is text-based (not scanned)
 - Check file size limits
 - Verify PDF library dependencies
-

Conclusion

JARVIS represents a sophisticated blend of modern frontend design, AI capabilities, and software engineering best practices. By leveraging CursorAI for development, the project demonstrates how AI-assisted coding can accelerate complex application development while maintaining high code quality.

The combination of a visually stunning UI with powerful RAG functionality creates an engaging experience that showcases both technical skills and creative design thinking. The architecture emphasizes clean separation of concerns, making the system maintainable and extensible for future enhancements.

Document prepared by Timal Pathirana - Software Engineer & AI Enthusiast