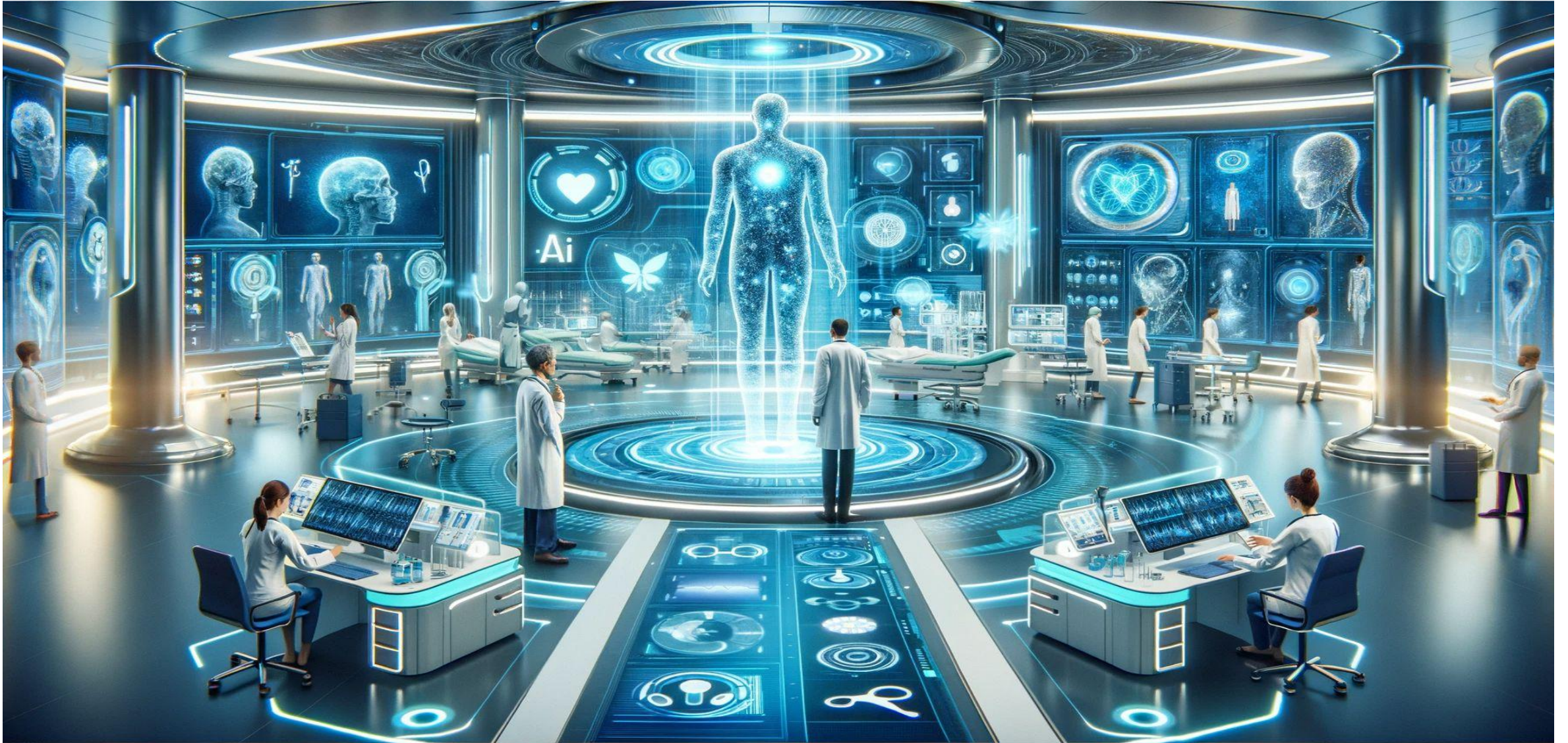


# Health AI - Intelligent Healthcare Assistant Generative AI with IBM



# Introduction

The project Health AI – Intelligent Healthcare Assistant focuses on building an advanced AI-powered system to support healthcare services. By combining intelligent automation and conversational capabilities, the assistant is designed to make healthcare delivery more efficient and user-friendly.

## **Team Members**

- Girupa Shankar K
- Bhuvanesh R
- Rohith Rekash V
- Jaya Prakash S

## 2. PROJECT OVERVIEW

The Health AI Intelligent Healthcare Assistant represents an advanced AI-powered platform engineered to enhance healthcare delivery through sophisticated intelligent support mechanisms. This comprehensive solution seamlessly integrates artificial intelligence and machine learning models with Large Language Models (LLMs) and vector databases to deliver user-centric healthcare assistance.



## Core Features

### ❑ **Conversational Interface**

The interactive conversational platform enables users to engage with the healthcare system through natural language processing, supporting both text and voice inputs. This interface facilitates health-related inquiries, policy guidance, and connection with healthcare professionals.

### ❑ **Policy Summarization**

An automated system that extracts and synthesizes complex healthcare policies into concise, accessible summaries. This feature enables patients, healthcare staff, and citizens to efficiently comprehend critical policy updates and regulatory changes.

### ❑ **Eco-Tip Generator**

A sustainability-focused module that delivers practical, actionable recommendations combining health optimization with environmental consciousness. The system promotes sustainable healthcare practices including waste minimization, energy conservation, and community health enhancement strategies.

### ❑ **Citizen Feedback Integration**

A comprehensive feedback collection system that gathers patient and public insights through multiple channels including surveys, chat interfaces, and structured forms. Advanced sentiment analysis and topic detection algorithms process this feedback to drive continuous system improvements.

### ❑ **Key Performance Indicator (KPI) Forecasting**

Predictive analytics capabilities that forecast essential healthcare metrics including bed occupancy rates, patient flow patterns, staff availability, and resource utilization. These insights support strategic planning and operational decision-making to optimize healthcare efficiency.

### ❑ **Anomaly Detection**

Intelligent monitoring systems that identify irregular patterns within healthcare data, including unexpected disease outbreaks or unusual resource consumption patterns. This capability enables proactive alerting and rapid response mechanisms to enhance patient safety.

## ❑ **Multimodal Input Processing**

Comprehensive input support that processes diverse data types including text, voice, images, and sensor data. This multimodal approach enables more intuitive, precise, and efficient interactions for both patients and medical personnel.

## ❑ **Streamlit to Gradio UI Integration**

Seamless conversion capabilities that transform Python applications into interactive web interfaces with simplified deployment processes. Gradio provides intuitive drag-and-drop components for enhanced user interaction and experience.

## **3. SYSTEM ARCHITECTURE**

### ○ **Frontend Architecture**

The frontend architecture of the Health AI system encompasses multiple interaction channels including web applications, mobile platforms, and voice-activated interfaces, ensuring accessible and user-friendly patient and citizen engagement.

### ○ **Backend Architecture**

The robust backend infrastructure manages core business logic, establishes secure connections with LLM services for predictive analytics and recommendations, and provides comprehensive APIs for secure, real-time frontend communication.

### ○ **Large Language Model Integration**

The integrated LLM processes user queries, interprets user intent, and generates contextually appropriate responses while interfacing with knowledge bases for information retrieval and personalized healthcare recommendations.

### ○ **Vector Database Implementation**

The vector database component stores embeddings of medical data, healthcare policies, and patient records, enabling rapid semantic search capabilities. This allows the AI system to retrieve contextually relevant information for accurate guidance and recommendations.

### ○ **Machine Learning Modules**

Specialized ML modules deliver advanced analytics, anomaly detection, and personalized health recommendations through comprehensive analysis of patient data, healthcare trends, and medical knowledge bases to support clinical decision-making.

## 4. SETUP AND INSTALLATION

### Prerequisites

- Python and Node.js installation with appropriate code editor
- OpenAI API access and vector database configuration
- Fundamental machine learning knowledge and required libraries
- Frontend development skills (React, Angular, or Flutter)
- Access to medical databases or pre-trained healthcare models

### Installation Procedure

- Install Python and Node.js, establish virtual environment
- Install backend dependencies and AI-specific libraries
- Configure vector database and create necessary indices
- Execute backend and frontend systems, conduct query testing and AI response validation

## 5. PROJECT STRUCTURE

**Directory Organization** • **App:** Core application codebase

- **Data:** Data storage and management files
- **Docs:** Comprehensive project documentation
- **Test:** Unit testing suites and validation scripts
- **Venv:** Virtual environment configuration
- **Txt:** Text processing utilities
- **Utils:** Utility functions and helper modules

## 6. APPLICATION DEPLOYMENT

### **Project Initialization Steps**

- Access terminal and clone repository from GitHub
- Navigate to project directory: Health AI-Intelligent Healthcare Assistant-AI
- Create and activate virtual environment
- Install dependencies using: `pip install -r requirements.txt`
- Initialize database and vector store configurations
- Configure API keys (OpenAI/LLM, database credentials, etc.)
- Launch backend server: `python backend/API/main.py`
- Execute frontend UI: `streamlit run frontend/app.py` or `gradio app.py`

### **Frontend Deployment (Streamlit)**

Navigate to frontend directory → Execute Streamlit → Access browser interface → Utilize Health AI Healthcare Assistant UI

### **Backend Deployment (FastAPI)**

Navigate to backend directory → Initialize FastAPI server → Access API documentation → Complete backend integration for Health AI system

## 7. API DOCUMENTATION

The comprehensive API documentation provides structured endpoints for patient interactions, health record management, and AI-driven recommendation systems, facilitating seamless healthcare support through conversational and analytical modules

## **8. AUTHENTICATION FRAMEWORK**

**Authentication Process Flow** : User/Client Registration: User enrollment or application registration to obtain API credentials

- API Key Generation: System generates unique API keys or JWT tokens
- Authentication Request: Client submits login credentials
- Access Token Provision: Server validates credentials and issues access tokens
- Authenticated API Calls: Include Authorization: Bearer <token> in request headers for all API interactions

## **9. USER INTERFACE DESIGN**

- Secure login portal for authenticated access
- Interactive chat interface supporting text and voice queries
- Comprehensive dashboard displaying health records and analytical reports
- AI-powered personalized recommendations and health insights

## **10. QUALITY ASSURANCE AND TESTING**

The testing framework for the Health AI Intelligent Healthcare Assistant ensures reliable and secure operation for both patients and healthcare providers. The testing protocol includes comprehensive unit testing to verify individual system modules including authentication, chat functionality, and health record access systems.

User Acceptance Testing (UAT) validates that end-users find the system functional, accurate, and user-friendly. This systematic testing approach ensures a robust and dependable healthcare assistance platform.

## **11. CURRENT LIMITATIONS**

While the Health AI Intelligent Healthcare Assistant demonstrates high efficiency, certain limitations have been identified. Occasional inaccuracies in symptom analysis may occur due to limited training data for rare medical conditions. Multimodal input processing, particularly when interpreting combined image and text data, may experience delays or misinterpretation with low-quality inputs.

Some users may encounter minor UI inconsistencies across different devices. During initial testing and deployment, several performance issues were identified that may impact system functionality or user experience. Certain features may exhibit slower response times under heavy system load, and occasional UI discrepancies have been observed across various devices. Specific data inputs may trigger unexpected errors when they don't conform to expected formats. Additionally, integration with external APIs may experience occasional connectivity issues during server downtime. Future releases will address these limitations to improve system stability and functionality.

## **12. FUTURE DEVELOPMENT ROADMAP**

Future enhancements for the Health AI Intelligent Healthcare Assistant will focus on expanding system capabilities and improving patient care quality. Planned developments include integration of advanced diagnostic tools, enhanced multimodal input support incorporating voice and image processing, and integration with additional health data sources for more personalized recommendations.

The system will implement predictive analytics for early disease detection, real-time health monitoring capabilities, and enhanced natural language understanding for more accurate and empathetic patient interactions. The user interface will be redesigned for improved accessibility and responsiveness, while automated error handling and notification systems will be implemented to ensure smoother operations. These enhancements aim to deliver a more robust, efficient, and user-friendly healthcare solution.



# Project Demonstration

The screenshot shows a Google Colab notebook interface. The browser tabs at the top include WhatsApp, healthai.py - Google Drive, colab.google, and Health AI.ipynb - Colab. The address bar shows the URL: colab.research.google.com/drive/1ifyATVYYjFQh2pGBRvvqqNaWwokZZ7W. The notebook title is 'Health AI.ipynb'. The menu bar includes File, Edit, View, Insert, Runtime, Tools, and Help. The toolbar shows a search bar, '+ Code', '+ Text', 'Run all', and a 'Connect T4' button. The code editor contains the following Python code:

```
import gradio as gr
import torch
from transformers import AutoTokenizer, AutoModelForCausalLM

# Load model and tokenizer
model_name = "ibm-granite/granite-3.2-2b-instruct"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32,
    device_map="auto" if torch.cuda.is_available() else None
)

if tokenizer.pad_token is None:
    tokenizer.pad_token = tokenizer.eos_token

def generate_response(prompt, max_length=1024):
    inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=512)

    if torch.cuda.is_available():
        inputs = {k: v.to(model.device) for k, v in inputs.items()}

    with torch.no_grad():
        outputs = model.generate(
            **inputs,
            max_length=max_length,
            temperature=0.7,
            do_sample=True,
            pad_token_id=tokenizer.eos_token_id
        )
```

The bottom of the interface shows a taskbar with various application icons, a search bar, and the system clock displaying 13:49 on 13-09-2025.

```
return generate_response(prompt, max_length=1200)

# Create Gradio interface
with gr.Blocks() as app:
    gr.Markdown("# Medical AI Assistant")
    gr.Markdown("**Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.**")

    with gr.Tabs():
        with gr.TabItem("Disease Prediction"):
            with gr.Row():
                with gr.Column():
                    symptoms_input = gr.Textbox(
                        label="Enter Symptoms",
                        placeholder="e.g., fever, headache, cough, fatigue...",
                        lines=4
                    )
                    predict_btn = gr.Button("Analyze Symptoms")

                with gr.Column():
                    prediction_output = gr.Textbox(label="Possible Conditions & Recommendations", lines=20)

            predict_btn.click(disease_prediction, inputs=symptoms_input, outputs=prediction_output)

        with gr.TabItem("Treatment Plans"):
            with gr.Row():
                with gr.Column():
                    condition_input = gr.Textbox(
                        label="Medical Condition",
                        placeholder="e.g., diabetes, hypertension, migraine...",
                        lines=2
                    )
```

Commands + Code + Text Run all

RAM  
Disk

```
[1] app.launch(share=True)

... /usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your Google Colab and restart your ses
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
  warnings.warn(
tokenizer_config.json: 8.88k/? [00:00<00:00, 767kB/s]
vocab.json: 777k/? [00:00<00:00, 26.6MB/s]
merges.txt: 442k/? [00:00<00:00, 19.5MB/s]
tokenizer.json: 3.48M/? [00:00<00:00, 96.3MB/s]
added_tokens.json: 100% 87.0/87.0 [00:00<00:00, 11.0kB/s]
special_tokens_map.json: 100% 701/701 [00:00<00:00, 73.5kB/s]
config.json: 100% 786/786 [00:00<00:00, 87.3kB/s]
`torch_dtype` is deprecated! Use `dtype` instead!
model.safetensors.index.json: 29.8k/? [00:00<00:00, 1.85MB/s]
Fetching 2 files: 0% 0/2 [00:00<?, ?it/s]
model-00001-of-00002.safetensors: 39% 1.94G/5.00G [00:59<01:20, 38.0MB/s]
model-00002-of-00002.safetensors: 100% 67.1M/67.1M [00:01<00:00, 20.8MB/s]
```

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run ``gradio deploy`` from the terminal in the working directory

# Medical AI Assistant

Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.

Disease Prediction Treatment Plans

Enter Symptoms

e.g., fever, headache, cough, fatigue...

Analyze Symptoms

Possible Conditions & Recommendations

Thank You.....