

# Medical RAG ChatBot Technical Documentation

Emre Büyükyılmaz

September 17, 2025

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Overview</b>                             | <b>1</b> |
| <b>2</b> | <b>Architecture</b>                         | <b>1</b> |
| 2.1      | High-Level . . . . .                        | 1        |
| 2.2      | Services and Ports . . . . .                | 2        |
| <b>3</b> | <b>Setup</b>                                | <b>2</b> |
| 3.1      | Prerequisites . . . . .                     | 2        |
| 3.2      | Environment Variables . . . . .             | 2        |
| 3.3      | Build & Run . . . . .                       | 2        |
| <b>4</b> | <b>Database Schema</b>                      | <b>2</b> |
| <b>5</b> | <b>REST API</b>                             | <b>2</b> |
| <b>6</b> | <b>UI (Gradio)</b>                          | <b>3</b> |
| <b>7</b> | <b>Data &amp; Evaluation Assets</b>         | <b>3</b> |
| 7.1      | Combined Evaluation Metrics (CSV) . . . . . | 3        |
| 7.2      | QA Files (Source for Evaluation) . . . . .  | 3        |
| <b>8</b> | <b>Performance &amp; Sizing</b>             | <b>3</b> |

## Overview

This project provides:

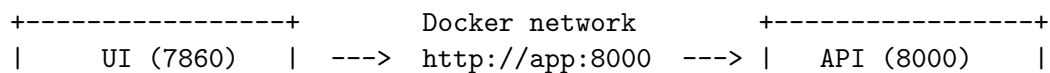
- **API** (FastAPI on Uvicorn): retrieval-augmented generation and chat/log management.
- **UI** (Gradio): a web interface to interact with the API.

Core components:

- **LLM**: meta-llama/Meta-Llama-3-8B-Instruct
- **Retriever**: FAISS index + BioBERT embeddings
- **SQLite**: API logs and chat sessions/messages
- **Docker Compose**: separate API and UI services

## Architecture

### High-Level





## Services and Ports

- **API:** 8000:8000
- **UI:** 7860:7860

## Setup

### Prerequisites

- Docker Desktop (Windows) with WSL2 backend
- NVIDIA GPU recommended (CUDA-enabled base image)
- Hugging Face account with a **read** token

### Environment Variables

Create `.env` next to `compose.yaml`:

```
HF_TOKEN=hf_XXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

### Build & Run

```
# From the project root:
docker compose up --build
# Subsequent runs:
docker compose up
```

UI: <http://localhost:7860>    API: <http://localhost:8000>

### Database Schema

SQLite at `data/logs/chat_logs.db`:

- **api\_logs:** id, timestamp, question, answer, retrieved\_sources (JSON), retrieval\_time, generation\_time, total\_time
- **chat\_sessions:** session\_id (PK), created\_at
- **chat\_messages:** session\_id (FK), timestamp, role (user—assistant), content

### REST API

Use Swagger UI for full details and interactive testing:

- Swagger UI: <http://localhost:8000/docs>
- ReDoc: <http://localhost:8000/redoc>

Endpoints:

- GET `/health`
- GET `/meta`
- POST `/query`
- GET `/logs`
- GET `/chats`

- GET /chats/{session\_id}
- DELETE /chats/{session\_id}

## UI (Gradio)

- Chat interface with example question boxes
- “New Chat” creates a fresh session and clears local history
- Auto-refresh of chat list after sends/deletes

## Data & Evaluation Assets

### Combined Evaluation Metrics (CSV)

- **Path:** data/evaluate/combined\_evaluation\_metrics.csv
- **Purpose:** main evaluation artifact combining human ratings and automated metrics.
- **Key columns** (header examples):  
 query, expected\_answer, generated\_answer, Relevance, Accuracy, Source\_Citation, Fluency, Query\_Length, Retrieval\_Time\_MS, Generation\_Time\_MS, Total\_Time\_MS, BLEU, ROUGE-1, ROUGE-2, ROUGE-L, METEOR, BERTScore-P, BERTScore-R, BERTScore-F1, Perplexity.

### QA Files (Source for Evaluation)

- A set of 100 questions (QA files) in ./data/evaluation/test\_questionsV2 directory was used to produce the rows in combined\_evaluation\_metrics.csv.

## Performance & Sizing

- **Cold start:** initial model download and load can be long on first run; subsequent runs use cache.
- **GPU memory:** 8B class model; ensure sufficient VRAM. Monitor with `nvidia-smi`.