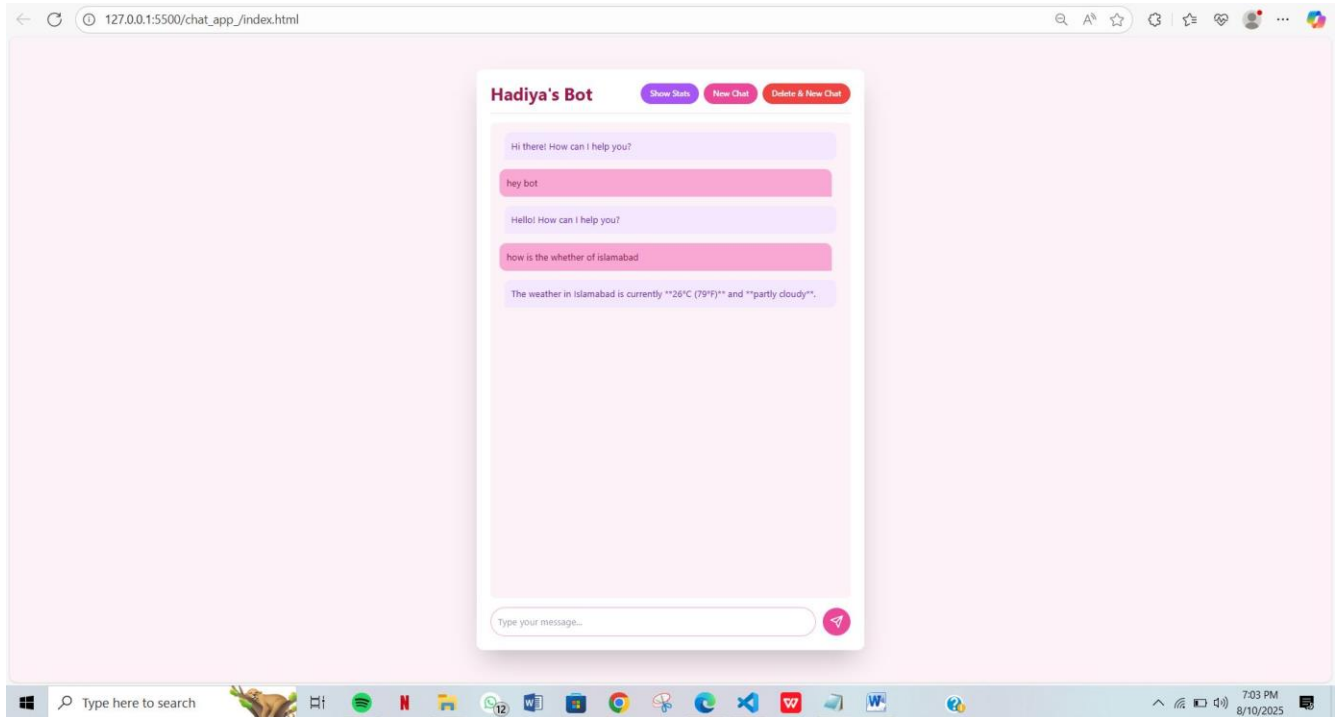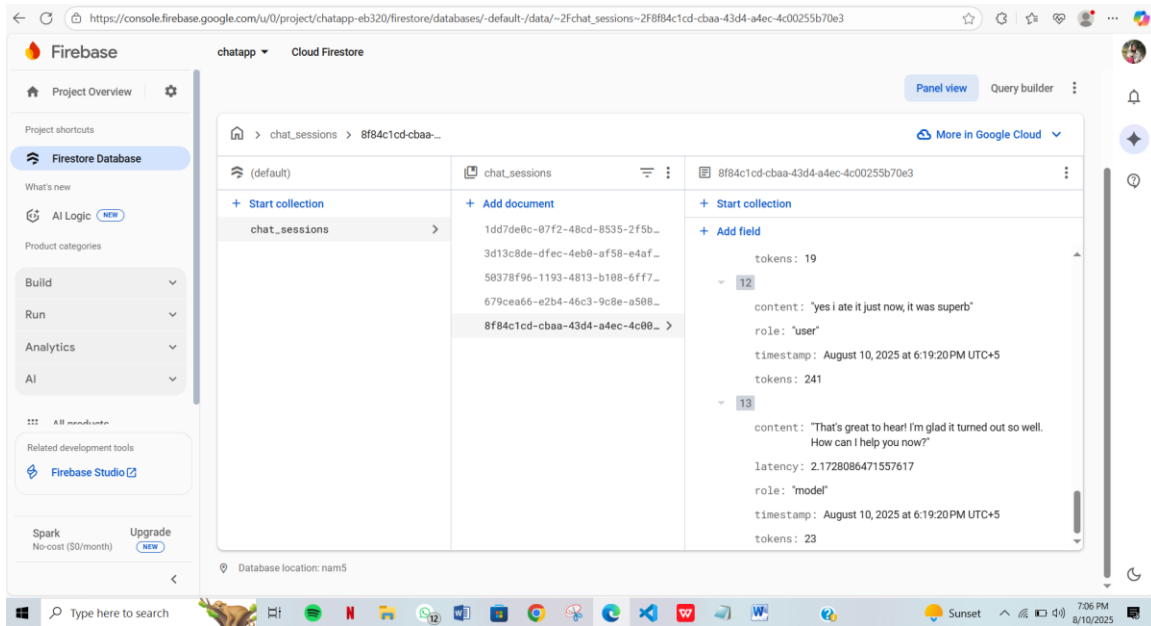# HADIYA'S BOT

## 1. Simple HTML page with streaming response

Implemented a single-page HTML client that opens a streaming connection to the backend and renders incremental chat/response tokens in real-time.
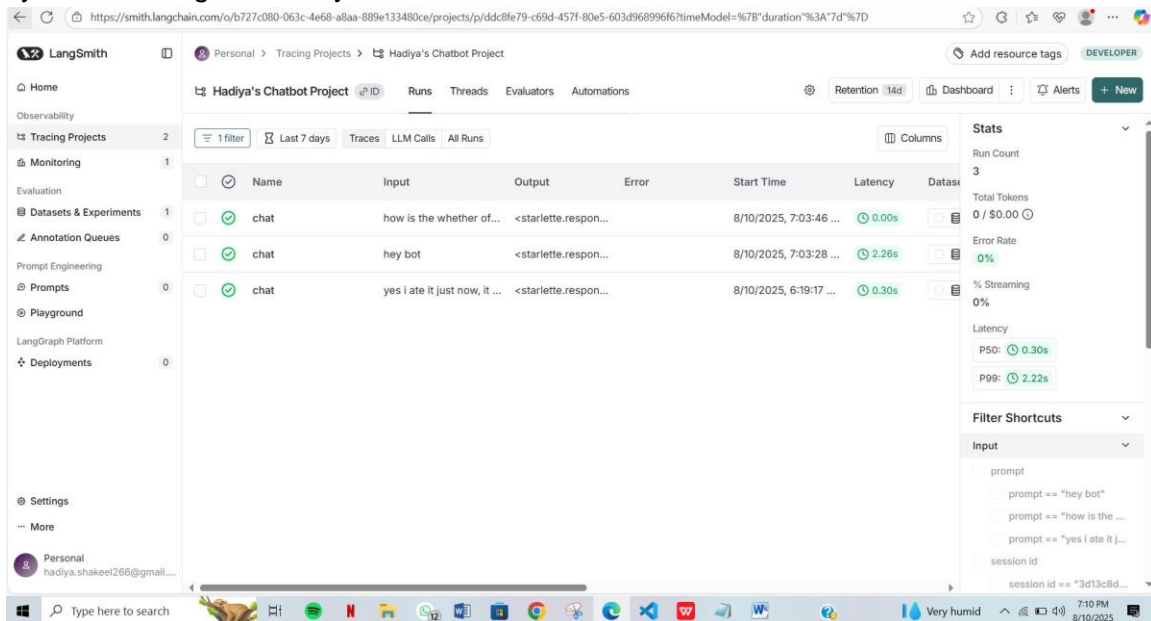


## 2. Save cost in DB

Costs (per request) are recorded to the database after each completed call: model name, prompt tokens, completion tokens, total cost (calculated), timestamp, and request id. Designed schema to support querying and aggregation.

## 3. system prompt (fetched from LangSmith)

System prompt is fetched at conversation start from LangSmith and injected as the system message for every chat session.

session.



## 4. Endpoint for statistical analysis (total cost, tokens used, latency)

Created a /stats endpoint that returns aggregated metrics: total cost, total tokens (prompt & completion), average latency per step.



.