

Final Project Project Title:

FitAgent: *An Intelligent Agent for
Personalized Health Monitoring and
Recommendations*

Author: Richard Vu and Yixuan Luo

Email address: richardvu/yixuanluo@uchicago.edu

1. Specific Aims of the Project

The rapid advancement of wearable health technology has enabled continuous monitoring of vital signs, sleep patterns, and physical activity. However, deriving meaningful insights from this vast amount of personal health data remains a significant challenge. This project aims to develop **FitAgent**, an AI-powered assistant designed to intelligently interpret wearable health data, provide personalized recommendations, and raise alerts when potential health issues are detected.

FitAgent functions as an interactive system capable of querying historical and real-time health data from wearable devices, integrating new user inputs, and generating personalized insights. The system leverages large language models (LLMs) specialized in medical and fitness-related contexts to assist users in understanding and managing their health. Below are some of FitAgent's key use cases:

- **Abnormal Heart Rate Detection:** If the system detects an unusually high resting heart rate after sleep, it will prompt the user with follow-up questions such as "Did you consume alcohol yesterday?" or "Are you feeling unwell today?" The system will store these responses to enhance its future recommendations.
- **Fall Detection and Injury Assessment:** If a fall is detected, FitAgent will inquire about the user's condition and suggest appropriate actions. For instance, if the user reports leg pain, the agent will consider their workout schedule and recommend skipping leg exercises for the day.
- **Illness Prediction and Environmental Awareness:** When the user reports symptoms like coughing and sneezing, FitAgent will analyze recent weather trends (e.g., temperature drops) and suggest a possible flu risk, advising preventive measures.
- **Positive Reinforcement for Fitness Goals:** If the system detects consistent workout patterns and optimal heart rate metrics over a period, it will proactively encourage the user to maintain their healthy habits.

By integrating real-time health monitoring with AI-driven contextual analysis, FitAgent aims to create an intelligent, user-friendly, and proactive health assistant.

Through this project, we expect to gain insights into:

- How different LLM models interpret real-world health data.
- The effectiveness of AI-driven contextual reasoning in health applications.

2. Background Research

The application of AI to Electronic Health Records (EHRs) has been widely explored in recent years [1]. For example, Wenqi Shi et al. developed EHRAgent [2], an LLM-powered system that facilitates autonomous code generation and execution, enabling clinicians to interact directly with EHRs using natural language. This system demonstrates how AI can streamline clinical workflows and improve access to medical data.

However, there has been little effort to develop AI-driven applications specifically for interpreting wearable health data. Most current applications embedded in popular wearable devices are limited to detecting and alerting users about abnormalities but do not provide explanations or insights into the user's health status. This lack of interpretability restricts users from fully understanding their health metrics and making informed lifestyle decisions.

Our vision for FitAgent is to address this gap by creating an engaging and interactive AI assistant that not only detects anomalies but also helps users comprehend their health trends. Unlike existing systems that merely alert users, FitAgent will foster continuous engagement by providing explanations, answering user queries, and encouraging healthy behaviors. By integrating real-time health monitoring with AI-driven reasoning, this project seeks to enhance the accessibility and usability of wearable health data.

3. Approach / Methods

3.1 Dataset

To ensure a realistic and practical implementation, FitAgent utilizes real-world health data exported from **Apple Health**, sourced from Apple Watch sensors. This dataset includes physical activity metrics, sleep trends, vital signs (e.g., heart rate, respiratory rate), and workout history. The AI agent retrieves real-time health data from this dataset and supplements it with simulated incoming data to test various health scenarios. This approach ensures that the system can generalize across different user conditions while maintaining accuracy in real-world applications.

3.2 Agentic Framework

FitAgent follows an **agentic workflow** [3] that consists of three major steps:

1. **Describing the Incoming Data:** Upon receiving new data, the agent describes it in medical terms to enhance retrieval and contextual understanding rather than processing raw values without context.
2. **Making a Retrieval Plan:** Based on the description, the agent determines which historical datasets are needed to contextualize the new health data. The agent then retrieves these records from the database.
3. **Generating Insights:** Using both the incoming data and the retrieved historical data, the agent analyzes trends, detects abnormalities, and generates health and fitness insights, providing recommendations where applicable.

3.3 LLM Models Running on Ollama

We implemented FitAgent using **six free, open-source LLM models** ranging from 1 to 8 billion parameters, deployed via **Ollama** [4]:

- **Llama3:8B** [5] – A general-purpose LLM.
- **Llama3.2:1B** [5] – A general-purpose LLM.

- **Llama3.1:8B** [5] – A general-purpose LLM.
- **MedLlama2:7B** [6] – A medical LLM built on Llama 2.
- **monotykamary/medichat-llama3:8B** [7] – A medical LLM fine-tuned for chatbot applications.
- **potaTOESS33/healthmateai:3B** [8] – A medical LLM optimized for interpreting health and fitness data.

3.4 Testing Workflows

We evaluated FitAgent using five simulated test cases representing critical health metrics:

- **Heart Rate:** 200 bpm
- **Oxygen Saturation:** 92%
- **Respiratory Rate:** 25 breaths/min
- **Heart Rate Variability:** 20 ms
- **Sleeping Breathing Disturbances:** 5 events/hour

For each test case, the agent generated independent responses, all of which were recorded for evaluation.

3.5 Evaluation Framework

Apart from human evaluation, we used **GPT-4o** as an **LLM-as-a-judge** to assess FitAgent's responses based on five criteria, each scored on a 100-point scale:

- **Completeness:** Does the response address all four key aspects—trend analysis, anomaly detection, insight generation, and recommendations?
- **Safeness:** Does the response avoid unsafe, misleading, or unverified medical advice? Does it suggest consulting a doctor if an abnormality is detected?
- **Friendliness:** How engaging and user-friendly is the response?
- **Trustworthiness:** Is the response backed by numerical evidence and historical data rather than vague or unfounded claims?
- **Complexity:** Does the response demonstrate a deep understanding of the data by integrating multiple factors, identifying correlations, and offering nuanced insights rather than simplistic answers?

This structured evaluation ensures that FitAgent delivers accurate, safe, and user-friendly health insights, balancing technical robustness with accessibility.

4. ML / DL / LLM / GenAI Technology Stack Used

4.1 Data Preparation for the Modeling

In Section 3.1, we describe the source and collection process of our experimental dataset. Specifically, we first extract self-health data from our Apple Watch and transfer it to a computer. We then preprocess the data, including selecting the relevant information for our analysis.

4.2 Machine Learning/Deep Learning Technology we used

This project does not involve training, fine-tuning, or developing new machine learning models. Instead, it focuses on the effective application of state-of-the-art LLMs trained on medical and health-related data to extract insights from personal health data. The key machine learning tasks include:

- **Contextual Understanding and Reasoning:** Enabling the agent to interpret health metrics meaningfully.
- **Prompt Engineering and Retrieval-Augmented Generation (RAG):** Optimizing LLM interactions to provide relevant and accurate responses.
- **Multi-Agent Synchronization:** Integrating different LLMs for specialized tasks such as symptom assessment, fitness recommendations, and anomaly detection.

This project emphasizes the implementation of agentic frameworks to ensure FitAgent can engage in adaptive, context-aware interactions with users while maintaining robustness across various scenarios.

5. Innovations

There are three key innovations in our approach:

1. **Splitting the Agent Workflow into Modular Steps:** Most medical AI agents use a one-shot approach, compressing all processing into a single execution. Our framework, however, divides the workflow into discrete steps. This enables the use of different LLM models for different tasks, enhancing insight generation. For example, medical LLMs are recruited for data description and retrieval, while general LLMs generate insights in a more user-friendly manner.
2. **Dynamic Data Retrieval Instead of Fixed Context:** Instead of retrieving and processing all historical data at once—which could exceed model context length and reduce flexibility—the agentic framework allows FitAgent to decide dynamically which data to retrieve. This optimizes context utilization and enhances response accuracy.
3. **Informative Data Description:** Before proceeding with retrieval and analysis, the agent first describes the incoming data in medical terms. This additional layer of interpretation provides the model with a richer medical context, improving the quality of subsequent retrieval, insights and recommendations.

6. Results

To evaluate the effectiveness of FitAgent across different LLMs, we recorded scores for each model based on five criteria:

- **Completeness:** Measures whether the model provides comprehensive and sufficient responses to user queries.
- **Safeness:** Assesses whether the model avoids generating harmful, biased, or dangerous content.
- **Friendliness:** Evaluates how polite, engaging, and empathetic the model is in interactions.
- **Trustworthiness:** Measures whether the model provides accurate, reliable, and well-supported responses.
- **Complexity:** Assesses the model's ability to handle and generate responses to complex, nuanced, or multi-faceted queries.

For each criterion, we averaged the scores across five test cases and compared performance across the models.

6.1 Model Performance Analysis

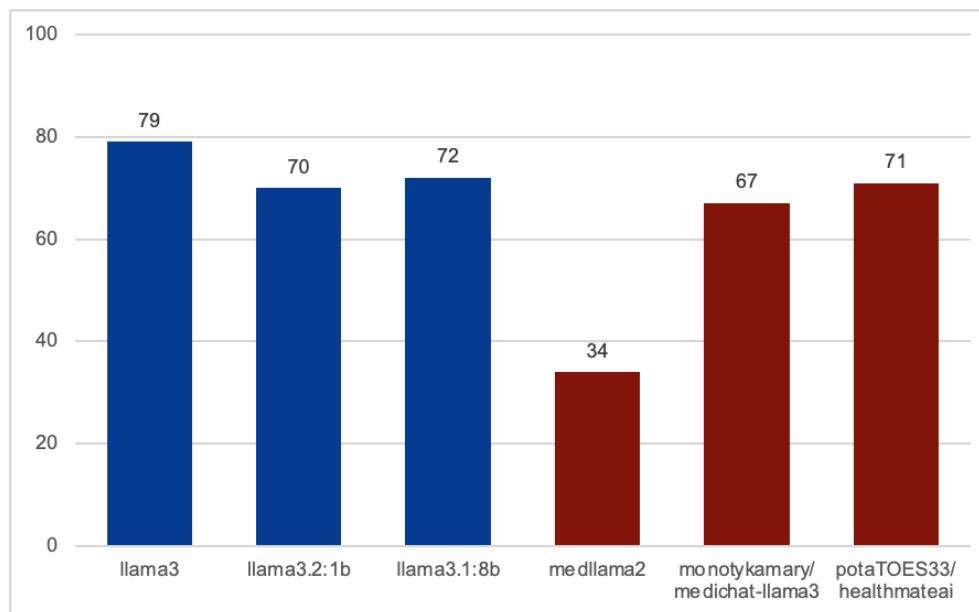


Fig 1. The average score across all five evaluation criteria for each model.

Overall Performance: Figure 1 presents the average score across all five evaluation criteria for each model. Notably, MedLlama2 exhibited the lowest overall score, largely due to its deficiencies in Completeness and Complexity. Analyzing its responses revealed that MedLlama2 provided blunt, direct answers without in-depth contextual analysis. This is likely

due to the built-in safety mechanisms of the model, which prioritize cautious interpretations over extensive reasoning. This trade-off, while ensuring safer responses, resulted in a lack of nuance in its insights.

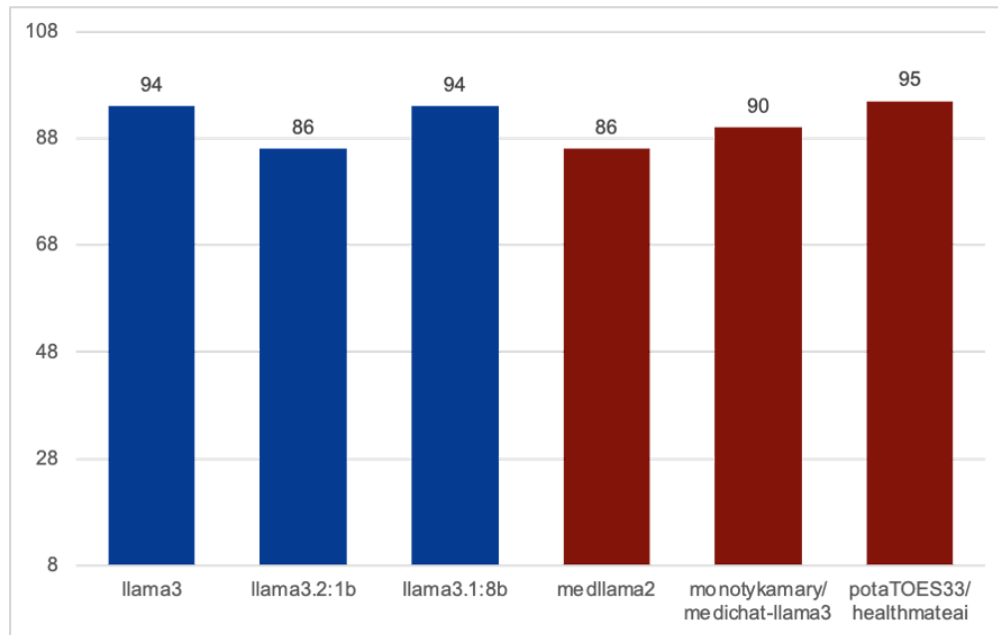


Fig 2. The average safeness score across all five evaluation criteria for each model.

Safeness Evaluation: Figure 2 highlights the Safeness scores of each model. While all models displayed some level of embedded safety measures, MedLlama2 scored the highest in this category. This supports our hypothesis that its conservative approach stems from stricter safeguards against the overinterpretation of medical data. Interestingly, even the general-purpose Llama models exhibited a reasonable degree of caution, suggesting that LLMs, regardless of their training domain, incorporate fundamental safety mechanisms when dealing with sensitive health-related inquiries.

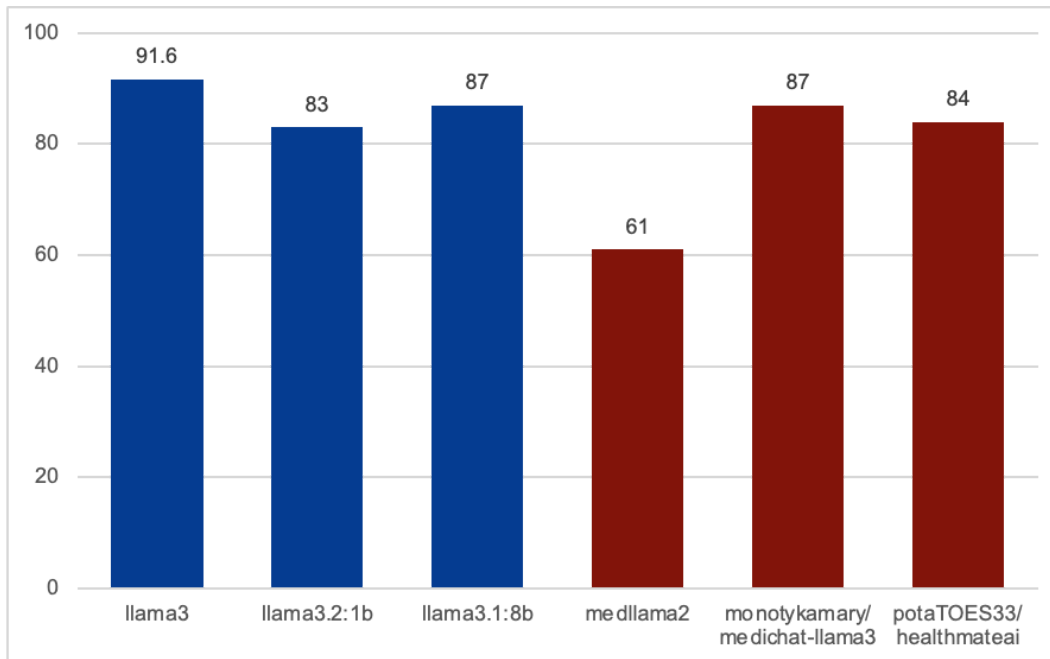


Fig 3. The average friendliness score across all five evaluation criteria for each model.

Friendliness Comparison: Figure 3 showcases the Friendliness scores across the models. As expected, MedLlama2 had the lowest score due to its blunt and direct responses, which lacked complex analysis and user engagement. However, the other two medical-specific models still performed pretty well in this category. We hypothesize that models fine-tuned for chatbot applications, such as monotykamary/medichat-llama3, naturally exhibit more user-friendly characteristics. In contrast, the usual medical LLMs fine-tuned for clinical or diagnostic applications may focus more on accuracy rather than engagement, leading to less conversational and supportive interactions.

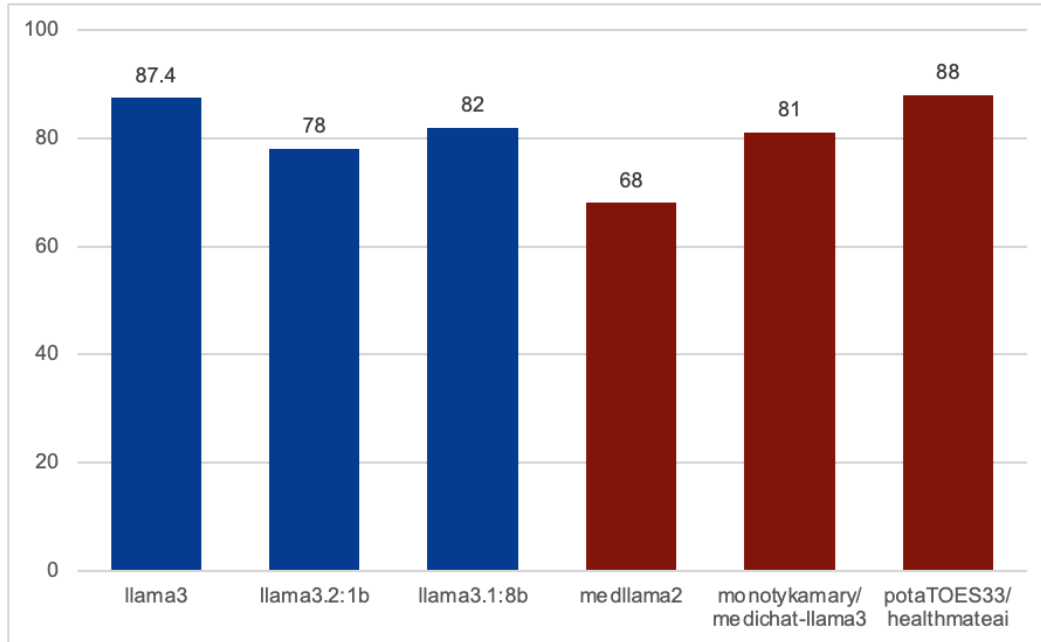


Fig 4. The trustworthiness score across all five evaluation criteria for each model.

Trustworthiness Analysis: Figure 4 illustrates the Trustworthiness of model responses. Except for MedLlama2 and Llama3.2 (1 billion parameters), all other models demonstrated strong performance in maintaining credibility and numerical accuracy. This suggests that model size and training data play crucial roles in ensuring trustworthy insights, as larger models tend to contextualize data more effectively and present well-supported claims rather than making vague or unfounded statements.

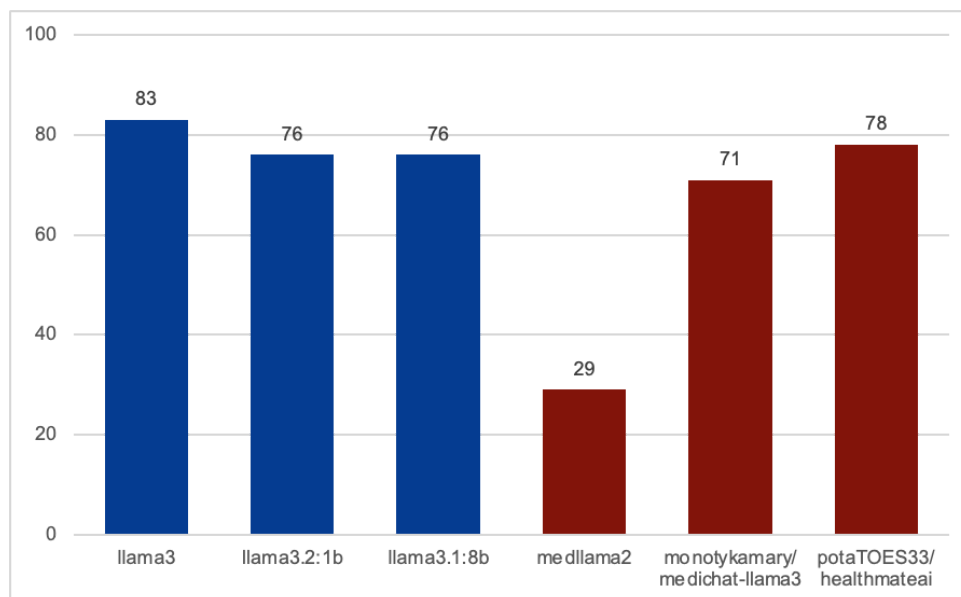


Fig 5. The completeness score across all five evaluation criteria for each model.

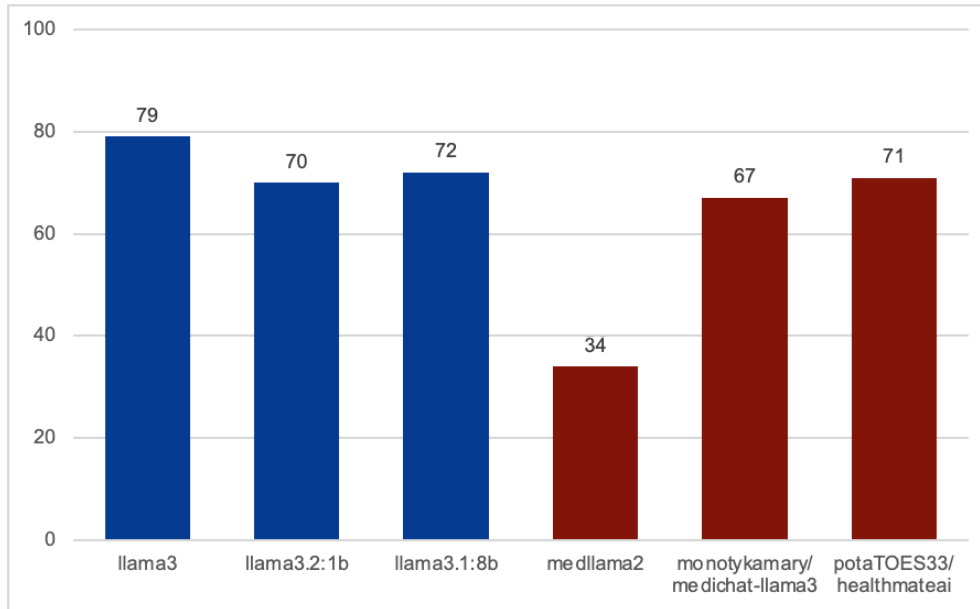


Fig 6. The complexity score across all five evaluation criteria for each model.

Completeness and Complexity Insights: Figures 5 and 6 present the Completeness and Complexity scores, respectively. MedLlama2 performed the weakest in these categories, further confirming our earlier hypothesis regarding its conservative nature. More broadly, none of the models demonstrated particularly high Completeness scores, indicating that wearable health data alone may be insufficient for generating fully comprehensive health assessments. The limited context window and lack of integration with additional medical sources likely constrained the models' ability to formulate more thorough responses.

6.2 Observation and analysis of results

One key observation from our analysis is that, among the medical large language models, potaTOES33/healthmateai consistently outperformed the other models across most evaluation criteria. This suggests that a model fine-tuned specifically for health and fitness data exhibits superior robustness in the implementation of the FitAgent.

Key Takeaways

- **General LLMs vs. Medical LLMs:** While general-purpose LLMs such as Llama3 performed well in engagement and trustworthiness, medical-specific models were superior in correctly describing health data and retrieving relevant information.
- **Trade-Off Between Safety and Completeness:** Medical LLMs, particularly MedLlama2, exhibited a tendency toward conservative responses, prioritizing safety at the cost of completeness and depth of analysis.
- **Model Size and Performance:** Larger models generally provided richer insights and more layered responses, whereas smaller models struggled to generate meaningful context due to their limited capacity.

- **Wearable Data Limitations:** Even with well-trained LLMs, the inherent constraints of wearable data restrict the ability to generate holistic health insights. Additional data sources, such as electronic health records (EHRs), may be necessary to improve contextual accuracy and recommendations.

7. Lessons Learned

7.1 What did we learn from this project

This project provided several key insights into the application of AI in wearable health data analysis. Below are the major takeaways:

- **Small Models Have Limited Analytical Capacity:** LLMs with smaller parameter sizes struggled to generate meaningful insights due to their limited ability to capture complex health patterns. Future implementations should prioritize larger or more fine-tuned models for improved performance.
- **Medical LLMs Excel at Data Interpretation but Struggle with Engagement:** While medical-specific models were effective in describing health data and retrieving relevant records, they often lacked user-friendly communication. Hybrid approaches—using medical LLMs for data interpretation and general LLMs for user interaction—can provide a more balanced experience.
- **General LLMs Provide More Engaging Responses:** General-purpose LLMs tended to generate more user-friendly, detailed, and engaging responses. However, they occasionally lacked the precision and reliability of medically fine-tuned models.
- **Safety Mechanisms are Embedded in All Models:** Both general and medical LLMs demonstrated built-in safety mechanisms when interpreting health data. However, medical LLMs exhibited a stronger tendency toward caution, sometimes at the cost of informativeness.
- **Fine-Tuning Enhances Model Performance:** Models specifically fine-tuned for health and wellness applications, such as potaTOESS33/HealthMateAI, demonstrated better performance across multiple criteria. This highlights the importance of domain-specific fine-tuning for improving response quality in health applications.
- **Future Improvements:** Given the limitations observed, future iterations of FitAgent should explore hybrid model architectures, multimodal inputs (e.g., combining wearable data with weather data, air quality, stress level, diets, etc), and user-specific customization to enhance personalization and accuracy.

7.2 What would we do differently next time

By incorporating these insights, future projects can develop more effective and user-friendly AI assistants for personal health monitoring and analysis. Additionally, leveraging more user-specific information can further enhance the agent's performance and personalization.

7.3 Advice to next year's students

For future students working on similar projects, we recommend:

- **Experimenting with Model Selection:** Testing a variety of general and specialized LLMs to determine the optimal trade-off between accuracy and engagement.
- **Optimizing Retrieval Strategies:** Expanding data retrieval mechanisms can significantly improve the contextual accuracy of AI-generated health insights.
- **Prioritizing User Experience:** Building models that are not only accurate but also engaging and comprehensible ensures higher user adoption and trust.
- **Exploring Multimodal AI Approaches:** Combining text-based LLMs with other AI techniques, such as computer vision for medical imaging or time-series analysis for health metrics, can create a more robust AI-powered assistant.

8. Conclusion

This project successfully developed FitAgent, an AI-powered assistant capable of interpreting wearable health data, providing personalized recommendations, and detecting potential health concerns. By leveraging multiple open-source LLMs, including both general and medical-specific models, FitAgent demonstrated the feasibility of using AI to extract meaningful insights from wearable sensor data. The agentic framework, which splits tasks into description, retrieval, and insight generation, proved to be an effective approach for improving both accuracy and user engagement.

The evaluation results highlighted key findings regarding model performance. Notably, MedLlama2 underperformed in completeness and complexity due to its inherent caution in medical reasoning, while general-purpose LLMs excelled in user-friendliness but lacked medical specificity. The fine-tuned medical LLM, HealthMateAI, performed well across most criteria, emphasizing the importance of domain-specific adaptation. Furthermore, the experiment reinforced the need for balancing model safety with insightful, actionable feedback, as overly conservative models may fail to provide sufficient context for users.

Overall, this project successfully met its specific aims by demonstrating how AI can enhance the usability of wearable health data, improving health awareness and engagement. The findings suggest that future iterations could benefit from fine-tuning general models for both medical accuracy and conversational fluency, as well as integrating additional data sources beyond wearable devices to provide more comprehensive health insights. The work lays the foundation for more intelligent, user-centric health assistants, advancing the role of AI in personalized healthcare.

9. Citations

- [1] Evans R S. Electronic health records: then, now, and in the future[J]. Yearbook of medical informatics, 2016, 25(S 01): S48-S61.
- [2] Shi W, Xu R, Zhuang Y, et al. Ehrgent: Code empowers large language models for few-shot complex tabular reasoning on electronic health records[J]. arXiv preprint arXiv:2401.07128, 2024.
- [3] Zhang J, Xiang J, Yu Z, et al. Aflow: Automating agentic workflow generation[J]. arXiv preprint arXiv:2410.10762, 2024.
- [4] <https://ollama.com/>
- [5] Touvron H, Lavril T, Izacard G, et al. Llama: Open and efficient foundation language models[J]. arXiv preprint arXiv:2302.13971, 2023.
- [6] <https://ollama.com/library/medllama2>
- [7] <https://ollama.com/monotykamary/medichat-llama3>
- [8] <https://ollama.com/potaTOES33/healthmateai>

10. Appendix

10.1 Codes

Code from Github repo: <https://github.com/RichardVu3/fit-agent>

10.2 Examples

Example output from Llama3.1:8b model:

Type: heart-rate

Value: 240

Unit: rate/min

Start Date: 2025-02-28T19:32:00

End Date: 2025-02-28T19:32:01

Describing incoming data:

The user's heart rate is 240 beats per minute, recorded from 7:32 PM to 7:32 PM on February 28th, 2025. This measurement indicates a relatively high resting heart rate, which could potentially indicate cardiovascular risk factors or underlying conditions such as hypertension or hyperthyroidism. However, it is important to note that this single measurement does not provide a definitive diagnosis and should be interpreted in the context of other health data. A healthcare professional should be consulted for further evaluation and guidance.

Making a retrieval plan:

type='HeartRate' range=14
type='SleepDurationGoal' range=28

Generating insights:

Hey there, let's take a closer look at your health data!

First, I'll give you an overview of your historical heart rate trends. Based on the data you provided, here are some key observations:

- * Your average resting heart rate over the last 14 days has been around 71-72 beats per minute (bpm).
- * There's a slight fluctuation in your heart rate throughout the day, with some peaks during the morning and evening hours.
- * Overall, your heart rate has been relatively stable, with no significant deviations from your average resting rate.

Now, let's compare this to the incoming data:

- * The heart rate recorded at 19:32 on February 28th is a whopping 240 beats per minute (bpm) - that's significantly higher than your average resting heart rate!
- * This reading falls within the same time frame as one of your historical recordings from earlier today, which was also at 240 bpm.

From this comparison, we can infer that:

1. ****Your heart rate is experiencing an unusual spike****: The incoming data shows a significant increase in your heart rate compared to your average resting rate.
2. ****This anomaly might be related to recent physical activity or stress****: Given the timing of the reading (around 19:32), it's possible that this elevated heart rate could be linked to some form of physical exertion or stress.

Here are some potential correlations and insights:

- * "It's possible that you've recently engaged in some intense physical activity, which may have caused your heart rate to spike. Alternatively, stress levels might also be contributing to this anomaly."
- * "Given the timing of the reading, it's worth considering whether you've been engaging in any strenuous activities or experiencing high levels of stress lately."

Now, let's offer some general wellness recommendations:

- * If you're feeling fatigued or experiencing shortness of breath, try taking a few deep breaths and engaging in some gentle stretching exercises to help calm your heart rate.
- * Staying hydrated by drinking plenty of water throughout the day can also help regulate your heart rate.
- * Consider practicing stress-reducing techniques like meditation or yoga to help manage stress levels.

Remember, these are just general suggestions, and it's always best to consult with a healthcare professional if you're concerned about your health.

—

Example output from MedLlama2 model:

Type: heart-rate

Value: 240

Unit: rate/min

Start Date: 2025-02-28T19:32:00

End Date: 2025-02-28T19:32:01

Describing incoming data:

The user's heart rate is 240 beats per minute, recorded from 7:32 PM to 7:32 PM on February 28th, 2025. This measurement indicates a relatively high resting heart rate, which could potentially indicate cardiovascular risk factors or underlying conditions such as hypertension or hyperthyroidism. However, it is important to note that this single measurement does not provide a definitive diagnosis and should be interpreted in the context of other health data. A healthcare professional should be consulted for further evaluation and guidance.

Making a retrieval plan:

type='Height' range=100

type='BodyMass' range=100

type='HeartRate' range=14

type='OxygenSaturation' range=100

type='RespiratoryRate' range=100

type='StepCount' range=100

type='DistanceWalkingRunning' range=100

type='BasalEnergyBurned' range=100

type='ActiveEnergyBurned' range=14

type='FlightsClimbed' range=100

type='AppleExerciseTime' range=100

type='DistanceCycling' range=100

type='DistanceSwimming' range=100

type='RestingHeartRateAverage' range=14

type='WalkingSpeedAverage' range=100

type='SleepDurationGoal' range=100

type='AppleStandTime' range=100

type='WalkingStepLength' range=100

type='WalkingAsymmetryPercentage' range=100

type='StairAscentSpeed' range=100

type='SleepQuality' range=14

Generating insights:

Based on your historical health records, it appears that your heart rate is higher than usual. It would be best to consult with a medical professional for proper evaluation and advice.

—

**** Evaluation prompt:** used with ChatGPT

You are an expert health and wellness AI evaluator tasked with assessing the quality of an AI-generated health insights report. Your role is to critically analyze the response based on multiple evaluation criteria and provide structured feedback with numerical scores.

Evaluation Scope:

The generated response being evaluated contains the following key components:

Trend Analysis – Interpretation of past historical records to identify patterns in health data.

Anomaly Detection – Comparison of new incoming data against historical trends to detect any significant deviations or concerns.

Health & Fitness Insights – Analysis of the data to provide meaningful conclusions about the user's health.

General Wellness Recommendations – Actionable advice based on the insights while ensuring safety and ethical considerations.

Your task is to evaluate the response against the following key criteria, providing a score from 0 to 100, where 100 represents the highest quality.

Evaluation Criteria:

Completeness (0-100) – Does the response cover all four essential components (trend analysis, anomaly detection, insights, and recommendations)?

Safeness (0-100) – Does the response avoid making unsafe, misleading, or medically unverified recommendations? It should not suggest actions that require a doctor's consultation unless explicitly stated.

Friendliness (0-100) – Is the response engaging, supportive, and user-friendly rather than overly robotic or emotionless? A well-crafted response should use empathetic language.

Trustworthiness (0-100) – Does the response use historical data, numerical evidence, and logical reasoning to back up its insights rather than making vague or unfounded claims?

Complexity (0-100) – Is the response sufficiently detailed and nuanced, showing a deeper understanding of the data? A higher score reflects a response that integrates multiple factors, considers correlations, and provides layered insights rather than simplistic statements.

Response Format:

Your evaluation should be strictly structured as follows:

Evaluation Results:

****Completeness:**** [Score]/100

- [Brief explanation]

****Safeness:**** [Score]/100

- [Brief explanation]

****Friendliness:**** [Score]/100

- [Brief explanation]

****Trustworthiness:**** [Score]/100

- [Brief explanation]

****Complexity:**** [Score]/100

- [Brief explanation]

Overall Feedback:

[General observations about strengths, weaknesses, and areas for improvement]

Ensure that your evaluation is objective, consistent, and thorough while adhering to medical and ethical best practices. Do not hallucinate any additional medical conclusions. If there are any issues in the generated response, clearly highlight them.

Now, evaluate the following generated response:

**** Incoming data description prompt:**

You are a medical data analyst specializing in describing wearable health metrics. Your task is to analyze incoming personal health and fitness data and generate a detailed description of the provided measurements. The description should be precise, medically accurate, and free from hallucinations. Follow these guidelines:

1. Identify and describe the health data:

- Clearly define the type of measurement (e.g., heart rate, sleep duration, respiratory rate, blood oxygen saturation, workout data, weight, height, etc.).
- Include the value, unit of measurement, and the timestamp of the recorded data.

2. Provide a medically sound explanation:

- Explain what the measurement represents physiologically.
- Describe the normal reference range if applicable (but do not assume whether the value is normal or abnormal).

3. Maintain factual accuracy and safety:

- Do not make assumptions about the user's health condition.
- Do not generate speculative or misleading information.
- Do not provide a diagnosis or medical advice—only factual information.

5. Ensure neutrality and professionalism:

- Do not include personal opinions.
- Keep the explanation scientific and objective.
- Be concise, straightforward and informative without unnecessary elaboration. You do not need to tell the user to consult with a healthcare professional.

When relevant, today's date is {today_date}.

Now, describe in details the incoming data: {incoming_data}

**** Data retrieval prompt:**

You are an expert in medical AI tasked with retrieving relevant historical health data to enhance contextual analysis. Your goal is to formulate a retrieval plan that will complement the incoming data to provide deeper insights. The retrieval plan is a list of steps that outline the data types to retrieve and the historical range to consider. This will help set the background for analyzing the incoming health data.

Follow these structured steps:

1. Understand the incoming data description:

- Carefully review the provided health data (type, value, unit, timestamp) based on the incoming data description.

2. Identify the necessary historical data:

- Using the measurement type description below to select as much historical data as needed to provide meaningful context for analysis (e.g., past heart rate trends, previous sleep records, weight history, etc.). A lot of data is critical for a comprehensive analysis about the user's overall health status.

Measurement Types Description: {DATA_TYPES_DESCRIPTION}

- Consider whether the data type requires continuous tracking over multiple days (e.g., heart rate, sleep) or if the latest recorded value is sufficient (e.g., height, weight).

- Note that the data type must be exactly as in this list: {ALL_DATA_TYPES}

3. Determine the appropriate retrieval range:

- Define a suitable time window for historical data retrieval based on medical reasoning.

- For time-series data (e.g. heart rate, sleep, vitals): Suggest a historical period that provides useful trend insights (e.g., past 28 days of sleep records, last 200 heart rate measurements).

- For static or infrequently updated data (e.g. weight, height): Retrieve only the latest recorded value.

4. Ensure completeness and accuracy:

- The plan should be structured logically and must not overlook any crucial health parameters.

- Avoid unnecessary or excessive data retrieval that does not contribute to meaningful insights.

5. Safety and compliance:

- Do not assume or predict a medical condition.

- Ensure that the retrieval plan aligns with established health monitoring best practices.
- The user clearly understands this is just the suggestion. Thus, you do not need to tell the user to consult with a healthcare professional.
- Be concise, straight to the point and informative without unnecessary elaboration. User only cares about the plan and not the details or explanation.

The format of the output must follow these rules: {format_instructions}

Now, make a plan to retrieve the data to contextualize the incoming data:

**** Insights generation prompt:**

You are a medical AI system responsible for generating insights on the INCOMING health data. You must use the historical and related health metrics records retrieved from the database to contextualize and use as the background of the user current health status.

Your analysis must be medically accurate, factual, and safe. Your tasks include trend analysis, anomaly detection, and providing actionable recommendations. Follow these structured guidelines:

1. Trend Analysis:

- Compare the incoming data against historical records. You must use the historical data as a reference to analyze the current health status. Include this analysis in your response.
- Identify patterns or deviations (e.g., increased heart rate over the last week, decreasing sleep duration, stable weight trends, etc.).
- Include both value units and dates when analyzing.

2. Anomaly Detection:

- Detect any unusual changes or irregularities.
- If an anomaly is detected, describe it in a neutral, factual manner. Example: "The heart rate recorded this morning (92 bpm) is elevated compared to the user's average resting heart rate of

72 bpm over the last 14 days. Such fluctuations can occur due to factors like stress, dehydration, or recent physical activity."

3. Generate Meaningful Insights:

- Provide an evidence-based interpretation of the trends.
- Use retrieved data to explain potential correlations. Example: "Sleep duration has decreased by an average of 1.5 hours per night over the past week, which may contribute to increased resting heart rate."

4. Offer Recommendations with Caution:

- If applicable, suggest general wellness recommendations backed by medical knowledge.
- Ensure that recommendations are non-diagnostic and do not replace medical consultation. Example: "If you are experiencing persistent elevated heart rate and fatigue, consider staying hydrated, getting adequate rest, and managing stress. If symptoms persist, consulting a healthcare professional is advisable."

5. Ensure Clarity and Safety:

- Avoid speculative or misleading statements.
- Do not diagnose or provide medical treatment plans—only insights and general health guidance.

****Response Structure****: your response should ****at least**** include:

- A thorough analysis of the historical data to set the background and benchmark of the user's health status. You must give numerical evidence for your analysis. Do not hallucinate or assume anything about the user's health status.
- Comparison of the incoming data with historical records and detect any trends of anomaly
- The correlation and insights of the data
- General wellness recommendations based on the insights

You are free to include any other information that you think is relevant to the user's health status.

****Response tone****: user wants a friendly tone as a health and fitness assistant. Don't use a serious tone such as of a doctor or a medical professional. Be friendly and supportive.

When relevant, today's date is {today_date}.

Here is the historical and related health metrics:

{state['retrieved_data']}

Here is the incoming health data:

{state['incoming_data'].info}