

Enhancing Clinical Decision Support and EHR Insights through LLMs and the Model Context Protocol: An Open-Source MCP-FHIR Framework

Abul Ehtesham¹, Aditi Singh², Saket Kumar³

¹Kent State University, USA

²Department of Computer Science, Cleveland State University, USA

³Northeastern University, USA

aehtesha@kent.edu, a.singh22@csuohio.edu, kumar.sak@northeastern.edu

Abstract—Enhancing clinical decision support (CDS), reducing documentation burdens, and improving patient health literacy remain persistent challenges in digital health. This paper presents an open-source, agent-based framework that integrates Large Language Models (LLMs) with HL7 FHIR data via the Model Context Protocol (MCP) for dynamic extraction and reasoning over electronic health records (EHRs). Built on the established MCP-FHIR implementation, the framework enables declarative access to diverse FHIR resources through JSON-based configurations, supporting real-time summarization, interpretation, and personalized communication across multiple user personas, including clinicians, caregivers, and patients. To ensure privacy and reproducibility, the framework is evaluated using synthetic EHR data from the SMART Health IT sandbox (<https://r4.smarthealthit.org/>), which conforms to the FHIR R4 standard. Unlike traditional approaches that rely on hardcoded retrieval and static workflows, the proposed method delivers scalable, explainable, and interoperable AI-powered EHR applications. The agentic architecture further supports multiple FHIR formats, laying a robust foundation for advancing personalized digital health solutions.

Index Terms—Clinical Decision Support, Electronic Health Records, Model Context Protocol, FHIR, Large Language Models, Agentic Workflow, Health Literacy, Explainable AI, Interoperability,

I. INTRODUCTION

Despite the widespread adoption of electronic health records (EHRs) [1], significant challenges persist in enhancing clinical decision support (CDS), reducing physician documentation burdens, and improving patient comprehension of health information. While the integration of HL7 Fast Healthcare Interoperability Resources (FHIR) [2] and mandates such as the 21st Century Cures Act [3] have expanded access to structured medical records, a substantial gap remains between data availability and its meaningful interpretation by end users particularly clinicians, caregivers, and patients [4].

Recent advances in large language models (LLMs) [5], including OpenAI’s GPT-4 [6] and emerging open-source alternatives [7], offer promising capabilities in summarizing, interpreting, and simplifying complex medical content. However, integrating LLMs into clinical workflows [8] continues to face obstacles related to effective context injection, consistent access to structured EHR data, and ongoing concerns

about replicability, safety, and explainability in medical AI systems [9], [10].

Earlier efforts, such as the LLM on FHIR system [11], demonstrated that mobile applications could retrieve and interpret patient records by leveraging LLMs and function-calling mechanisms built on Stanford’s Spezi ecosystem [12]. While these solutions showed the feasibility of converting structured medical data into user-friendly narratives, their broader clinical adoption has been limited due to platform-specific constraints, static data pipelines, and inconsistent LLM outputs. To address these limitations, this paper presents an open-source, agent-based framework that integrates Large Language Models with the Model Context Protocol (MCP) [13] and a dynamic FHIR server, using the publicly available MCP-FHIR implementation [14]. Through declarative JSON configurations and RESTful interactions, the system enables LLMs to retrieve and summarize diverse FHIR-based patient records without the need for custom hardcoding. This modular design facilitates transparent, reproducible reasoning workflows across heterogeneous EHR systems and varying FHIR formats. Figure 1 illustrates the architectural contrast between traditional, manually wired integrations and the MCP-based approach. While the former relies on hardcoded API calls and static pipelines, MCP enables dynamic, declarative tool invocation and composable agent workflows.

The framework also incorporates the MCP-Agent module [15], which abstracts server interactions and enables dynamic orchestration of AI agents. By supporting composable design patterns, MCP-Agent simplifies the development of robust, model-agnostic agents and enhances the flexibility and interoperability of the system.

The remainder of this paper is structured as follows. Section II reviews related work on LLM integration in clinical applications. Section III details the system architecture, including the roles of the MCP server, LLM engine, and patient-facing interface. Section IV describes the implementation workflow. Section V presents a practical use case demonstration. Finally, Section VI concludes the paper.

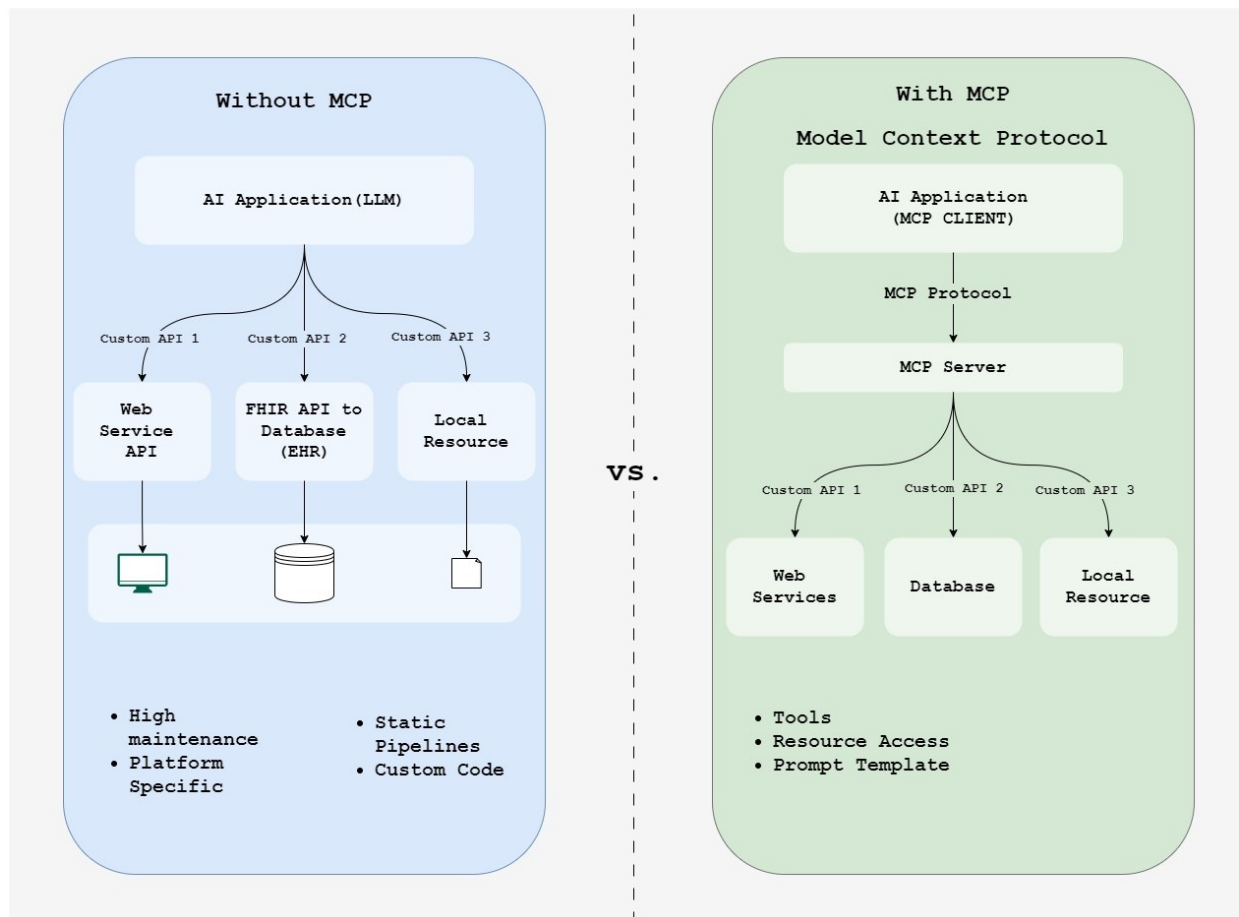


Fig. 1. Comparison of traditional AI application integration (left), where the LLM directly calls multiple custom APIs leading to high maintenance and platform-specific pipelines—versus the Model Context Protocol (MCP) approach (right), which centralizes access to tools, resources, and prompt templates through a unified MCP server.

II. RELATED WORK

The rapid advancement of large language models (LLMs) in recent years has catalyzed a wide range of healthcare applications [16], spanning clinical decision support, documentation, and patient education. Early research demonstrated the potential of LLMs to assist clinicians with documentation and summarization tasks [9], [17]. These studies showed that LLMs could extract meaningful insights from unstructured text; however, reliably incorporating structured electronic health record (EHR) data into clinical workflows remained a significant challenge.

Initiatives such as the LLM on FHIR system [18] served as early proof-of-concept efforts, employing mobile applications and function-calling mechanisms built on Stanford’s Spezi ecosystem. While these implementations successfully transformed complex medical data into accessible narratives, they were constrained by platform-specific dependencies, static configurations, and the need for manual prompt engineering. These limitations led to inconsistent outputs and hindered scalability and broader clinical adoption [10].

Recent work has shifted toward modular, open-source frameworks that prioritize interoperability and reproducibil-

ity [19]. The integration of the Model Context Protocol (MCP) with FHIR resources represents a significant step in this direction, offering a standardized, declarative interface for accessing diverse healthcare data. Unlike prior approaches that relied on hardcoded API calls, the MCP-FHIR framework enables flexible querying of FHIR resource types, supporting dynamic context injection, transparent reasoning workflows, and more consistent performance in LLM-assisted clinical applications [6], [14], [20].

These developments reflect a broader movement away from fragmented, custom-built integrations toward unified solutions that emphasize scalability, explainability, and cross-platform interoperability in healthcare AI. By addressing long-standing limitations in context management and structured data access, the proposed framework builds on this progress to advance clinical decision support and improve patient health literacy.

III. SYSTEM OVERVIEW

The proposed system addresses persistent challenges in clinical decision support (CDS), documentation overload, and patient health literacy by integrating Large Language Models (LLMs) with FHIR resources using the Model Context Proto-

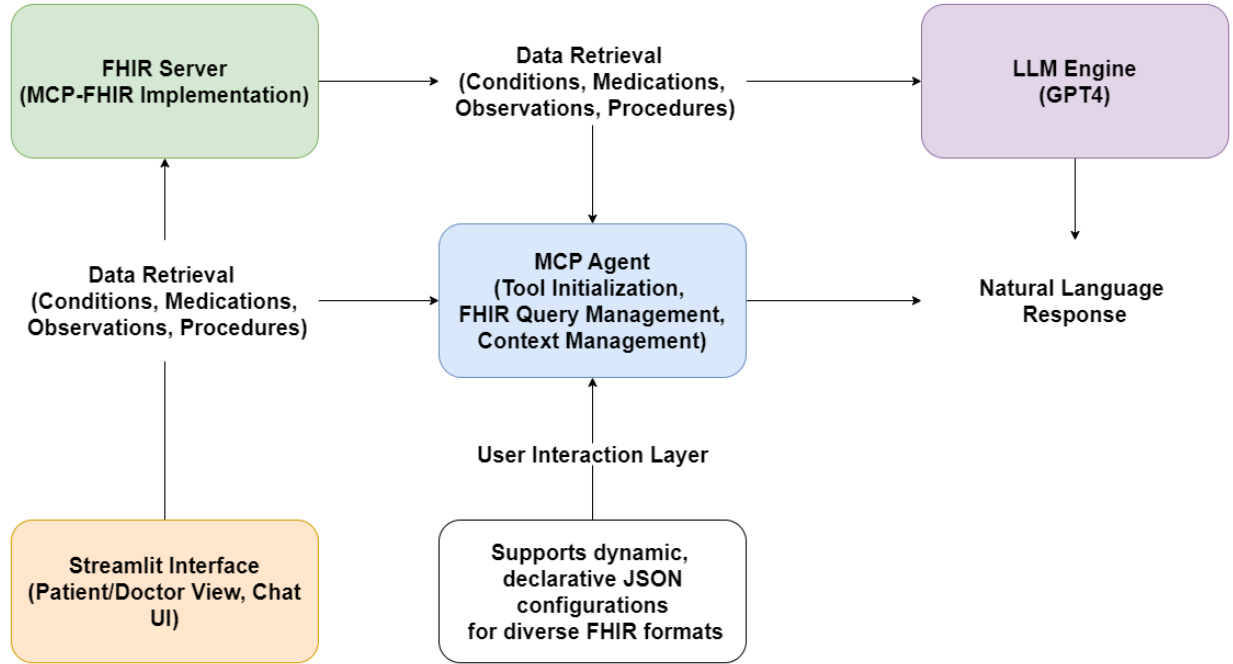


Fig. 2. System Architecture using MCP, FHIR, and LLM workflows.

col (MCP). This agent-based framework leverages a dynamic, MCP-enabled approach to query patient-specific FHIR data in real time and deliver natural language explanations tailored to various personas, including clinicians, caregivers, and patients.

At the core of the system is an LLM agent that orchestrates secure interactions between structured FHIR resources and the LLM reasoning engine. The integration is facilitated through a lightweight, open-source MCP-FHIR server that supports declarative access to diverse FHIR resource types via JSON configurations. A user-friendly interface, built using the Streamlit framework, allows end users to select patients, review summaries, and engage in conversational interactions with the system.

The operational workflow begins with the initialization of the MCP agent, which establishes connections to the FHIR server and dynamically retrieves patient records, such as conditions, medications, observations, and procedures. Based on the selected patient and persona, the agent composes a context-aware prompt that is transmitted to the LLM. The LLM then produces a synthesized, natural language explanation of the patient data, thereby supporting clinical reasoning and enhancing patient understanding.

This modular and extensible design promotes interoperability, transparency, and reproducibility in clinical workflows. In addition, the architecture is adaptable to various FHIR formats, as exemplified by the MCP-FHIR implementation, and can be extended with additional modules, such as imaging viewers,

lab analyzers, and multilingual translation tools.

IV. SYSTEM ARCHITECTURE

The system architecture is designed to support dynamic interactions between LLMs and structured FHIR data via the Model Context Protocol (MCP). This section details the primary components, workflow, and extensibility features of the framework.

A. Core Components

The architecture is composed of the following key components:

- **MCP Agent:** Serves as the central orchestrator of clinical workflows by initializing tool connections, selecting appropriate FHIR queries, and managing the context for LLM interactions.
- **FHIR Server:** Based on an open-source MCP-FHIR implementation, this server provides access to FHIR resources. It supports operations such as reading and searching for patient data via declarative JSON configurations.
- **LLM Engine:** Utilizes advanced LLMs (e.g., OpenAI's GPT-4) to synthesize information from FHIR data, delivering concise summaries and interpretations that enhance clinical reasoning and patient understanding.
- **Streamlit Interface:** A lightweight frontend that enables users to select personas (clinician, caregiver, patient),

review patient summaries, and interact with the system through conversational queries.

B. Workflow Overview

The system workflow proceeds as follows:

- 1) **Initialization:** The MCP agent establishes connections with the FHIR server and attaches the LLM engine.
- 2) **Data Retrieval:** Upon selecting a patient and a user persona, the agent queries the FHIR server for relevant EHR data, such as conditions, medications, observations, and procedures.
- 3) **Contextual Prompt Composition:** Based on the user-selected persona and retrieved patient data, the agent composes a tailored, context-aware prompt.
- 4) **LLM Inference:** The composed prompt is transmitted to the LLM engine, which returns a synthesized natural language response that informs clinical decision-making or enhances patient comprehension.

C. Extensibility

The architecture is inherently modular and extensible:

- The MCP framework enables dynamic integration of additional tools and resources, such as imaging viewers or lab analyzers.
- The use of declarative JSON configurations allows for seamless expansion to support additional FHIR resource types or alternative data formats.
- The composable design encourages reuse and interoperability across various EHR systems and clinical workflows.

Figure 2 illustrates the overall architecture, emphasizing the flow between the MCP agent, FHIR server, LLM engine, and user interface.

V. SYSTEM WORKFLOW AND IMPLEMENTATION

In this section, the implementation details of the open-source MCP-FHIR clinical assistant system are described. The architecture integrates a configurable Model Context Protocol (MCP) agent, a publicly accessible FHIR server, and an OpenAI LLM through a persona-based Streamlit user interface. This solution supports dynamic querying of FHIR resources and generates contextual clinical summaries in real time, ensuring adaptability, transparency, and applicability in both clinical and patient-facing scenarios.

A. Architecture Overview

Figure 2 illustrates the overall system architecture. The user interacts with a Streamlit-based front end, where they select a persona (e.g., Clinician, Caregiver, Patient) and choose a patient record for review. The MCP agent functions as an orchestration layer by invoking tools to fetch FHIR resources via standard APIs and then forwarding the structured results to an LLM (e.g., GPT-4o). The LLM generates personalized and comprehensible summaries or responses based on the persona-specific instructions.

B. Agent Configuration with FHIR Server

The system utilizes the open-source `mcp-fhir` server (available from Flexpa) which is integrated into the MCP agent through a YAML configuration. This configuration specifies the FHIR base URL and access credentials, thereby enabling the agent to interact with any HL7-compliant resource. The modular configuration approach supports dynamic discovery and retrieval of diverse FHIR resource types in accordance with HL7 standards.

C. Agentic Orchestration and Dynamic Prompt Generation

The MCP agent is initialized with an instruction set tailored to each persona (for instance, “You are a helpful assistant for a clinician working with EHRs”). The agent leverages standardized tool calls to retrieve clinical data such as Conditions, Medications, Observations, and Procedures from the FHIR server. Each retrieved resource is summarized and subsequently incorporated into a context-aware prompt. For example, a typical prompt may be structured as follows:

```
Persona: Clinician.  
Patient Name: John Doe.  
Conditions: Asthma; Hypertension.  
Medications: Metoprolol; Albuterol.
```

This prompt, which integrates both demographic and clinical details, is submitted to the LLM through an abstraction layer that manages context history and enforces consistency across interactions. The resulting LLM output is a concise clinical summary or explanatory response intended to support clinical decision-making and enhance patient understanding.

D. Interactive Multimodal User Interface

The user interface is implemented using Streamlit and is designed to be both intuitive and interactive. Key features include:

- A preview of patient lists with demographic summaries.
- Persona-specific, predefined question templates to facilitate initial queries.
- A free-form chat interface that supports continuous session-aware interactions.

Responses are progressively streamed into the UI, enabling real-time conversational interaction with the clinical assistant.

E. Agentic Benefits

The overall agentic workflow provides several advantages:

- **Modularity:** The system’s architecture permits independent development and reuse of individual tools.
- **Interoperability:** It can connect with any standard FHIR server via the MCP, ensuring broad compatibility.
- **Explainability:** Detailed prompt history and source data are available, supporting transparent clinical reasoning.
- **Scalability:** The architecture is designed for seamless integration with additional modules (e.g., on-device LLMs, clinical validation tools) in future extensions.

In summary, the system enables real-time clinical decision support and improved health literacy through a scalable, standards-based architecture. It dynamically supports diverse FHIR data, leverages agentic orchestration for context-aware interactions, and provides interpretable outputs for multiple user personas.

VI. USE CASE AND WORKFLOW EXAMPLE

This section demonstrates a typical end-to-end scenario that highlights the practical value of the MCP-FHIR framework for clinical decision support and patient education. The use case illustrates how dynamic FHIR querying, persona-based prompt generation, and LLM-based summarization are integrated into an interactive workflow.

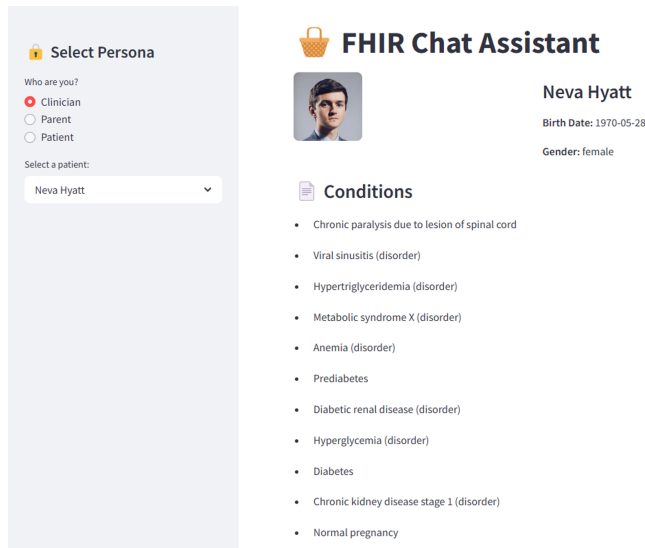


Fig. 3. An illustration of Clinician Persona.

Observations

- Estimated Glomerular Filtration Rate
- Blood Pressure
- Hemoglobin A1c/Hemoglobin.total in Blood
- Creatinine
- Body Mass Index
- Pain severity - 0-10 verbal numeric rating [Score] - Reported
- Platelet mean volume [Entitic volume] in Blood by Automated count
- Sodium
- Total Cholesterol

Fig. 4. An illustration of Observation by Clinician.

Procedures

- Intramuscular injection
- Peripheral blood smear interpretation
- Auscultation of the fetal heart
- Hepatitis B Surface Antigen Measurement
- Insertion of subcutaneous contraceptive
- Screening for chromosomal aneuploidy in prenatal amniotic fluid
- Auscultation of the fetal heart
- Chlamydia antigen test

Fig. 5. A Snippet of Procedures listed by Clinician.

A. Persona-Based Clinical Decision Support Interaction

In the clinical scenario, a user selects the “Clinician” persona from the Streamlit-based interface. The system then retrieves a list of recent patients, displaying key demographic details such as name, birthdate, and gender. Upon selecting a patient, the MCP agent dynamically fetches multiple FHIR resource types (e.g., Condition, MedicationRequest, Observation, and Procedure). These resources are organized into structured segments, allowing the clinician to quickly gain an overview of the patient’s clinical history.

B. Predefined and Freeform Query Options

To facilitate efficient clinical decision-making, the interface offers predefined query templates, such as:

- What treatment options are available for this patient’s conditions?
- Summarize the medications and conditions for this patient.
- List recent procedures and relevant laboratory insights.

When a predefined query is selected, the system composes a context-rich prompt that integrates persona information, patient demographics, and the retrieved FHIR data. The generated prompt is then submitted to the LLM, which returns a concise summary. For example, a typical clinician prompt may yield:

Patient John Doe, diagnosed with Type 2 Diabetes and Hypertension, is currently prescribed Metformin and Lisinopril. Recent laboratory results indicate elevated HbA1c levels. Consider revising the treatment regimen and advising lifestyle modifications.

C. Adaptive Workflow for Multiple Personas

The framework’s modular design also supports workflows for other personas, such as caregivers and patients. Under the “Caregiver” persona, the system alters the tone and complexity of the response to ensure the information is accessible. For instance, a caregiver’s query might result in:

◆ Pre-defined Questions

Select a question to ask:

List recent procedures and lab insights.

Ask Selected Question

Persona: Clinician. Patient Name: Denny Bins. Conditions: Laceration of thigh; Acute bronchitis (disorder) Medications: Chlorpheniramine Maleate 2 MG/ML Oral Solution; Ibuprofen 200 MG Oral Tablet; Acetaminophen 325 MG Oral Tablet Lab Results: Total Cholesterol; Tobacco smoking status NHIS; Blood Pressure; Body Weight; Blood Pressure; Pain severity - 0-10 verbal numeric rating [Score] - Reported; Blood Pressure; Low Density Lipoprotein Cholesterol; Triglycerides; Tobacco smoking status NHIS; Low Density Lipoprotein Cholesterol; Total Cholesterol; High Density Lipoprotein Cholesterol; Total Cholesterol; High Density Lipoprotein Cholesterol; Pain severity - 0-10 verbal numeric rating [Score] - Reported; Body Height; Platelet mean volume [Entitic volume] in Blood by Automated count; Platelet distribution width [Entitic volume] in Blood by Automated count; Platelets [#/volume] in Blood by Automated count; MCHC [Mass/volume] by Automated count; MCH [Entitic mass] by Automated count; Hematocrit [Volume Fraction] of Blood by Automated count; Hemoglobin [Mass/volume] in Blood; Leukocytes [#/volume] in Blood by Automated count; Low Density Lipoprotein Cholesterol; Body Height; Body Mass Index; Triglycerides; High Density Lipoprotein Cholesterol; MCV [Entitic volume] by Automated count; Erythrocytes [#/volume] in Blood by Automated count; Triglycerides; Erythrocyte distribution width [Entitic volume] by Automated count; Body Mass Index; Body Mass Index; Body Weight; Body Height; Pain severity - 0-10 verbal numeric rating [Score] - Reported; Tobacco smoking status NHIS; Body Weight List recent procedures and lab insights.

Here are the recent procedures and lab insights for Denny Bins:

Recent Procedures:

1. Suture Open Wound
 - Status: Completed
 - Code: Suture open wound (SNOMED Code: 288086009)
 - Performed: 2020-02-27
 - Reason: Laceration of forearm.

Lab Insights:

1. Blood Pressure:
 - Recorded: 2021-03-16
 - Systolic: 111.22 mmHg
 - Diastolic: 86.70 mmHg
2. Triglycerides:
 - Recorded: 2013-03-08
 - Value: 119.64 mg/dL
3. Total Cholesterol:

Type your question here...

Fig. 6. An illustration of Pre-defined Questions in a clinician scenario.

John Doe has diabetes and high blood pressure. He is taking medications to control these conditions, and recent tests suggest that his blood sugar levels could be better managed. It is recommended to review his treatment plan with his physician.

This adaptability ensures that clinical insights are communicated effectively according to the user's role and expertise.

D. Explainability and Session Continuity

To maintain transparency, each LLM output includes traceable references to the corresponding FHIR resources. This linkage ensures that clinical recommendations can be verified against the original data (e.g., medication names from `medicationCodeableConcept.text` and laboratory results from `Observation.code.text`). Moreover, the system retains a session history that preserves context across multiple turns in the conversation, thereby supporting cohesive multi-turn interactions and progressive decision support.

E. Impact on Clinical Workflows

The described workflow exemplifies several key advantages of the MCP-FHIR solution:

- **Real-Time Data Access:** Dynamic querying eliminates the need for custom API integrations, enabling immediate access to up-to-date FHIR resources.

- **Personalization:** The system adapts the complexity and tone of output to match the selected persona, thereby enhancing usability across clinical and patient settings.
- **Traceability and Explainability:** Each response is directly linked to the underlying FHIR data, ensuring that recommendations are both verifiable and interpretable.
- **Scalability:** The modular architecture supports future extensions, including on-device LLM processing and additional clinical data integrations.

In summary, the presented use case demonstrates how the MCP-FHIR framework provides a scalable, standards-based solution for real-time clinical decision support and patient education. The system's ability to dynamically integrate diverse FHIR data, generate context-aware prompts, and produce interpretable outputs underscores its potential to improve clinical workflows and health literacy.

VII. CONCLUSION

This paper presented a scalable, standards-based framework that integrates Large Language Models (LLMs) with FHIR resources through the Model Context Protocol (MCP) to enhance clinical decision support and EHR reasoning. The proposed system leverages a modular, agent-based architecture that supports dynamic querying, context-aware prompt generation, and real-time LLM-based summarization. By abstracting the complexities of FHIR data retrieval and integrating persona-specific workflows via a Streamlit interface, the framework

offers transparent and reproducible interactions that enhance both clinical decision-making and patient health literacy.

The implementation demonstrates the potential of combining LLMs with standardized FHIR access to generate concise and interpretable clinical insights. Key benefits include the modularity and scalability of the system, the ability to dynamically adjust language output to different user roles, and the provision of traceable, explainable outputs. These advantages are particularly relevant in addressing the growing challenges of documentation overload and the need for personalized health information access.

In summary, the integration of LLMs with MCP-enabled FHIR data access represents a significant step toward more intelligent, interactive, and patient-centered clinical decision support systems.

REFERENCES

- [1] Centers for Medicare & Medicaid Services, "Electronic health records," <https://www.cms.gov/priorities/key-initiatives/e-health/records>, accessed: Apr. 11, 2025.
- [2] Health Level Seven International, "FHIR Overview," <https://www.hl7.org/fhir/overview.html>, accessed: Apr. 11, 2025.
- [3] U.S. Food and Drug Administration, "21st Century Cures Act," <https://www.fda.gov/regulatory-information/selected-amendments-fdc-act/21st-century-cures-act>, 2025, accessed: Apr. 11, 2025.
- [4] S. Graham and J. Brookey, "Do patients understand?" *The Permanente Journal*, vol. 12, no. 3, pp. 67–69, 2008, published Online: September 1, 2008.
- [5] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, and J.-R. Wen, "A survey of large language models," 2025. [Online]. Available: <https://arxiv.org/abs/2303.18223>
- [6] OpenAI, "Gpt-4 technical report," 2023. [Online]. Available: <https://openai.com/research/gpt-4>
- [7] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao, "Large language models: A survey," 2025. [Online]. Available: <https://arxiv.org/abs/2402.06196>
- [8] M. M. Davis, R. Gunn, M. Cifuentes, P. Khatri, J. Hall, E. Gilchrist, C. J. Peek, M. Klowden, J. A. Lazarus, B. F. Miller, and D. J. Cohen, "Clinical workflows and the associated tasks and behaviors to support delivery of integrated behavioral health and primary care," *Journal of Ambulatory Care Management*, vol. 42, no. 1, pp. 51–65, 2019, pMCID: PMC6278604.
- [9] H. Nori, N. King, S. M. McKinney, D. Carignan, and E. Horvitz, "Capabilities of gpt-4 on medical challenge problems," 2023. [Online]. Available: <https://arxiv.org/abs/2303.13375>
- [10] J. A. Yeung, Z. Kraljevic, A. Luintel, A. Balston, E. Idowu, R. J. Dobson, and J. T. Teo, "AI chatbots not yet ready for clinical use," *Frontiers in Digital Health*, vol. 5, p. 1161098, Apr 2023.
- [11] P. Schmiedmayer, A. Rao, P. Zagar, V. Ravi, A. Zahedivash, A. Fereydooni, and O. Aalami, "Llm on fhir – demystifying health records," 2024. [Online]. Available: <https://arxiv.org/abs/2402.01711>
- [12] Stanford Spezi, "Stanford Spezi: Open-source framework for digital health applications," <https://github.com/StanfordSpezi>, 2025, accessed: Apr. 11, 2025.
- [13] A. Singh, A. Ehtesham, S. Kumar, and T. T. Khoei, "A survey of the model context protocol (mcp): Standardizing context to enhance large language models (llms)," *Preprints*, April 2025. [Online]. Available: <https://doi.org/10.20944/preprints202504.0245.v1>
- [14] Flexpa, "Flexpa mcp-fhir server," 2024. [Online]. Available: <https://github.com/flexpa/mcp-fhir>
- [15] L. AI, "Mcp agent: Build effective agents using model context protocol and simple workflow patterns," <https://github.com/lastmile-ai/mcp-agent>, 2025, accessed: 2025-04-11.
- [16] K. He, R. Mao, Q. Lin, Y. Ruan, X. Lan, M. Feng, and E. Cambria, "A survey of large language models for healthcare: from data, technology, and applications to accountability and ethics," 2025. [Online]. Available: <https://arxiv.org/abs/2310.05694>
- [17] N. H. Shah, D. Entwistle, and M. A. Pfeffer, "Creation and adoption of large language models in medicine," *JAMA*, vol. 330, no. 9, pp. 866–869, 2023.
- [18] P. Schmiedmayer, A. Rao, P. Zagar, V. Ravi, A. Zahedivash, A. Fereydooni, and O. Aalami, "Llm on fhir: Demystifying health records," *arXiv preprint arXiv:2402.01711*, 2024.
- [19] A. Ehtesham, A. Singh, G. K. Gupta, and S. Kumar, "A survey of agent interoperability protocols: Model context protocol (mcp), agent communication protocol (acp), agent-to-agent protocol (a2a), and agent network protocol (anp)," 2025. [Online]. Available: <https://arxiv.org/abs/2505.02279>
- [20] H. Touvron *et al.*, "Llama 2: Open foundation and fine-tuned chat models," 2023, meta AI. [Online]. Available: <https://ai.meta.com/llama/>