
FHIR-Hopper: A Neuro-Symbolic Agent for Reasoning over EHRs

Anonymous Authors¹

Abstract

Electronic health records (EHRs), increasingly exchanged via the Fast Healthcare Interoperability Resources (FHIR) standard, encode patient histories as deeply nested, multi-relational bundles that are difficult for large language models (LLMs) to reason over directly. We introduce *FHIR-Hopper*, a neuro-symbolic agentic framework that projects a FHIR bundle into an episodic chronological graph, materializes a budgeted, saliency-ranked linearization as the LLM’s initial context, and exposes deterministic graph-traversal tools so the agent can recover details on demand. Across three realistic FHIR clinical benchmarks, FHIR-Hopper attains the highest accuracy across multiple base LLMs *while keeping average input-token usage stable at $\sim 20k$ tokens, a typical $\sim 10\times$ reduction over the strongest retrieval baseline on long records*. These results suggest that decoupling structure-aware retrieval from neural reasoning is an effective design for clinical question answering over structured EHRs.

1 Introduction

While the widespread adoption of the HL7 Fast Healthcare Interoperability Resources (FHIR) standard (Bender & Sartipi, 2013) has improved syntactic interoperability for electronic health records (EHRs), leveraging large language models (LLMs) to reason over these complex records remains inefficient (Makhni et al., 2025; Li et al., 2024; Idrissi-Yaghir et al., 2025). Despite the momentum of frontier LLMs in clinical applications (Saab et al., 2024; Singhal et al., 2023; Shmatko et al., 2025; Heydari et al., 2025), the prevailing heuristic of flattening and concatenating massive JSON- or XML-formatted FHIR bundles into expanded context windows is fundamentally inefficient and suboptimal for two reasons. First, FHIR’s system-to-system serialization injects immense structural boilerplate and metadata,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

exacerbating context overflow and the “lost in the middle” phenomenon (Liu et al., 2024). Second, sequential text processing destroys the inherently multi-relational, directed graph structure of clinical data (e.g., an `Observation` referencing an `Encounter` anchoring a `Condition`). Forcing LLMs to track scattered universally unique identifier (UUID) cross-references in working memory severs the edges between temporally distributed events. Details about the FHIR standard are provided in Section A.1.

Existing mitigation strategies fall short: format conversions (e.g., JSON-to-YAML) fail to restore graph topology, standard retrieval-augmented generation (RAG) struggles with rigid identifiers and temporal dependencies (Cao et al., 2026; Tang & Yang, 2024; Amugongo et al., 2025), and prior agentic systems that directly query FHIR APIs suffer from hallucinated parameters and compounding multi-hop errors (Lee et al., 2025a; Jiang et al., 2025). Due to space constraints, a comprehensive discussion of related work—including tool-using LLM agents, graph-augmented retrieval, and prior EHR systems—is provided in Section A.2.

To address these limitations, we introduce *FHIR-Hopper*, a neuro-symbolic agentic framework for clinical question answering (Figure 1). FHIR-Hopper projects a patient’s bundle into an *episodic chronological graph*, resolving UUID references into explicit edges and grouping events into time-ordered episode nodes. From this graph, FHIR-Hopper materializes a token-budgeted, query-conditioned linearization as the LLM’s initial context. The agent subsequently invokes deterministic tools (traversing references, filtering temporally, and searching the hidden graph) to recover details on demand. This delegates deterministic lookups to a symbolic engine, freeing the LLM to focus strictly on higher-order reasoning. More specifically, in this work, our contributions are:

- An *episodic chronological graph* representation for FHIR bundles that preserves multi-relational and temporal structure while minimizing token overhead.
- *FHIR-Hopper*, a neuro-symbolic agent combining budgeted, saliency-ranked graph linearization with deterministic traversal tools for multi-hop clinical reasoning.
- Empirical evaluations across three FHIR benchmarks and three frontier LLMs, demonstrating that FHIR-Hopper achieves state-of-the-art accuracy while reducing aver-

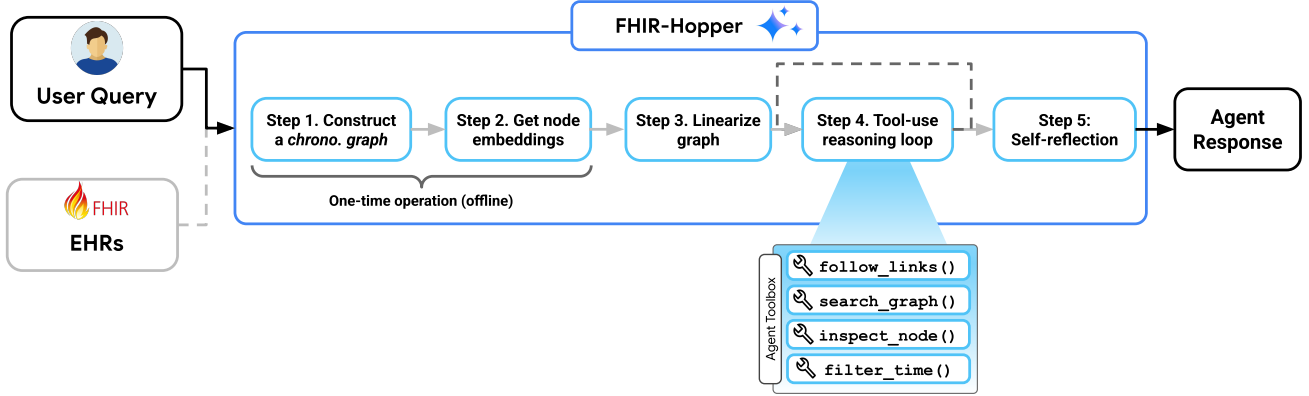


Figure 1. **Overview of the FHIR-Hopper architecture.** The framework projects a raw FHIR bundle into an episodic chronological graph and applies Saliency-Budgeted Linearization (SBL) to produce a token-efficient skeleton view within a specified budget. The LLM agent uses this skeleton together with deterministic tools (`follow_links`, `search_graph`, `inspect_node`, `filter_time`) to interactively traverse the graph and answer complex clinical questions.

age input-token usage by approximately $10\times$ over the strongest retrieval baseline on long records.

2 Methods

2.1 The FHIR-Hopper Framework

FHIR-Hopper consists of three coupled stages: (1) *episodic graph projection*, (2) *saliency-budgeted linearization* (SBL), and (3) *graph-constrained agentic reasoning*.

Episodic Graph Projection. While FHIR records form a directed graph $G = (V, E)$ of clinical resources V and references E (Tomaszuk et al., 2025a; ALMutairi et al., 2024; Tomaszuk et al., 2025b), this flat view ignores that clinical events cluster into encounters, and reference edges are often sparse. We project the bundle into an *episodic chronological graph* $H = (\mathcal{E}, \mathcal{R})$, where $\mathcal{E} = \{E_1, \dots, E_T\}$ is a time-ordered set of episode nodes (clinical encounters), and \mathcal{R} contains intra- and inter-episode references (Section A.4). To resolve sparsity, we assign orphan resources (lacking explicit `Encounter` references) via *latent episode inference* (Alg. 1): orphans are attached to the nearest preceding `Encounter` within a temporal window δ , or grouped into a synthetic δ -snapped episode.

Saliency-Budgeted Linearization (SBL). Because H typically exceeds an LLM’s token budget β , we construct a token-budgeted linearization L . We first score each resource $v \in V$ against the query q using cosine similarity $\phi(v, q) \in [-1, 1]$ in a shared text-embedding space (Gemini Embedding 1, $d = 3072$ (Lee et al., 2025b)).

Budgeted Selection. Let $x_v \in \{0, 1\}$ indicate if v is in L , and $S = \{v : x_v = 1\}$ be the active set. Let $A(S)$ be the set of *active episodes* containing at least one resource in S . Emitting an active episode incurs a header cost C_{head} , while

skipped episodes are compressed into a `[GAP]` marker at cost C_{gap} . The selection objective is:

$$\begin{aligned} \max_{S \subseteq V} \quad & \sum_{v \in S} \phi(v, q) \\ \text{s.t.} \quad & |A(S)|C_{\text{head}} + \sum_{v \in S} c(v) + N_{\text{gap}}(S)C_{\text{gap}} \leq \beta, \end{aligned} \quad (1)$$

where $c(v)$ is the token cost of v , and $N_{\text{gap}}(S)$ is the number of contiguous inactive episode runs. To showcase performance under extreme constraints, we set $\beta = 4000$, $< 0.1\%$ of average tokens in real-world EHRs (see Section A.3).

Greedy Density Solution. Since $N_{\text{gap}}(S)$ is non-linear, we solve (1) greedily. For a candidate $v \in E_t$, we define its marginal saliency-per-token as:

$$\rho(v | S) = \frac{\phi(v, q)}{c(v) + \mathbf{1}[E_t \notin \text{active}(S)]C_{\text{head}}} \quad (2)$$

We iteratively add the resource maximizing ρ , maintaining chronological order and emitting gaps between non-adjacent episodes. This concentrates the budget on high-saliency nodes while compressing inactive periods.

Graph-Constrained Agentic Reasoning. The LLM agent receives the SBL view as its *initial* context. To navigate beyond this skeleton and to increase its context, the agent uses four deterministic tools via standard function-calling APIs: **1) `inspect_node(id)`:** Returns full JSON of a resource for precise values (e.g., dosages, labs). **2) `search_graph(keywords)`:** Semantically searches the hidden graph for pruned nodes. **3) `follow_links(id)`:** Traverses explicit relational edges (e.g., `MedicationRequest` \rightarrow `Condition`) for multi-hop reasoning. **4) `filter_time(...)`:** Scopes nodes by date range and type. To prevent runaway loops, we enforce a strict 15-step limit, triggering a graceful termination callback that forces final answer synthesis.

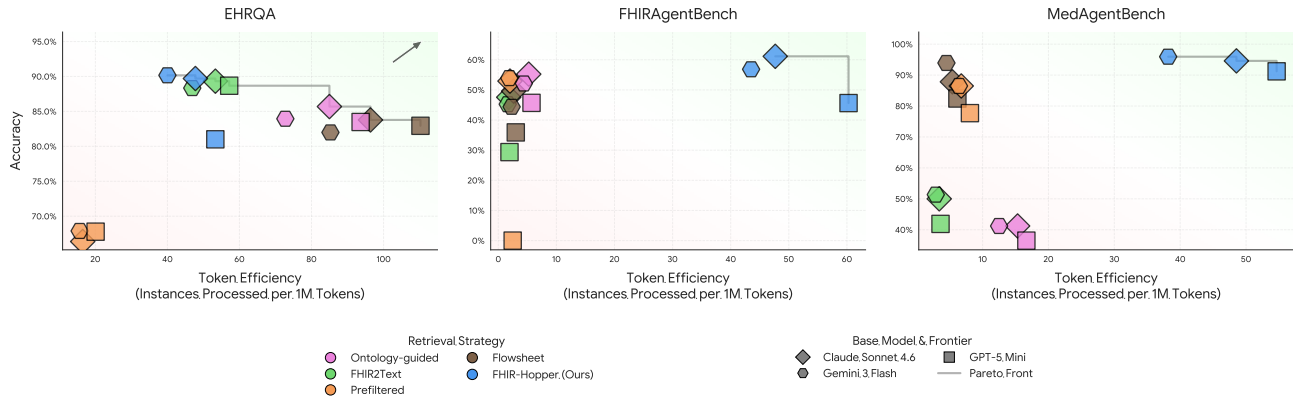


Figure 2. Pareto efficiency of retrieval strategies. Axes: accuracy vs. instances processed per 1M tokens (up/right is better, as pointed by the gray arrow in the leftmost plot). FHIR-Hopper strictly dominates the Pareto front on complex records (FHIRAgentBench, MedAgentBench). On short EHRQA records, FHIR-Hopper’s lower efficiency stems from its $\mathcal{O}(1)$ fixed overhead (agent prompts/tools). While this overhead outweighs pruning savings on trivial bundles, it avoids the $\mathcal{O}(N)$ context explosion of baselines, enabling scalable reasoning on long, real-world histories.

2.2 Datasets

We evaluate on three benchmarks of increasing difficulty (details are provided in Section A.3): **EHRQA**: 5,133 LLM-generated QA pairs over 20 synthetic Synthea bundles (Walonoski et al., 2018), testing basic single/two-hop retrieval. **MedAgentBench** (Jiang et al., 2025): 148 clinician-authored retrieval tasks over 100 synthetic profiles containing 700K+ clinical data elements. **FHIR-AgentBench** (Lee et al., 2025a): 2,931 clinician-sourced questions on de-identified MIMIC-IV-FHIR records (avg. 9K resources/patient), requiring complex multi-hop, temporal, and large-bundle reasoning.

2.3 Baselines

We compare FHIR-Hopper against four diverse retrieval strategies: **FHIR2Text**: Heuristically flattens key FHIR resources into natural language sentences and removing all non-text characters. **Ontology-guided** (Kabak et al., 2026): Expands query tokens via a static medical ontology, returning positively scored resources in reverse chronological order alongside fixed anchors (e.g., Patient). **Pre-filtered** (Schmiedmayer et al., 2025): Filters by clinical heuristics (e.g., keeping only active meds, deduplicating repeat observations) and exposes a condensed identifier triplet (type, name, date) to save tokens. **Flowsheet (SQL on FHIR)** (Grimes et al., 2025): Our implementation of the SOTA tabularization approach, which flattens hierarchical bundles into chronologically ordered Markdown tables (Date, Concept, Value) prepended with patient demographics to support identity-based reasoning.

2.4 Evaluation

We evaluate clinical retrieval performance along two axes: **accuracy** (fraction of correct answers evaluated by an LLM-

as-a-judge against deterministic ground truth rules, limiting subjective bias) and **token usage** (total context plus all multi-turn tool-call tokens).

3 Results

Accuracy. FHIR-Hopper consistently attains the highest accuracy across all benchmarks (Table 1, Figure S6). On the simpler benchmark, EHRQA, most baselines exhibit ceiling effects, but FHIR-Hopper alone clears 0.9 (0.902 ± 0.008 , Gemini 3 Flash). On the substantially more difficult benchmark, FHIR-AgentBench, FHIR-Hopper (Claude Sonnet 4.6) reaches 0.611 ± 0.018 , significantly outperforming all baselines (two-sided McNemar’s test, $p < 0.001$). On MedAgentBench, FHIR-Hopper (Gemini 3 Flash) dominates at 0.959 ± 0.030 ($p < 0.01$), while FHIR2Text and Ontology-guided collapse to the 0.4–0.5 range. Crucially, because Ontology-guided uses the exact same Gemini embedding for retrieval, FHIR-Hopper’s superiority stems directly from its structural representation and graph-constrained reasoning, rather than embedding quality (ablation detailed in Section A.7).

Context Token Usage. Average input-token usage is reported in Table 1. FHIR2Text and Prefiltered scale poorly on complex records, with FHIR2Text exploding up to 614,787 tokens on FHIR-AgentBench (frequently triggering context overflows). In contrast, FHIR-Hopper maintains a stable footprint of $\sim 20,000$ tokens across all benchmarks, an approximate $10\times$ reduction over the strongest retrieval baseline (Ontology-guided) on FHIR-AgentBench, and up to $30\times$ relative to FHIR2Text.

Ultimately, FHIR-Hopper is strictly pareto-dominant on the longer, real-world benchmarks (FHIR-AgentBench and MedAgentBench), achieving maximum accuracy at the lowest token cost (Figure 2). We explicitly note a deliberate

Table 1. Accuracy and token usage of FHIR retrieval strategies across base LLMs and benchmarks. We compare FHIR-Hopper against four baselines. Accuracy is reported as mean \pm 95% confidence interval (Acc \uparrow). Token usage per query is reported as the mean (Tok \downarrow), where lower is more desirable. The best configuration (highest accuracy, smallest context footprint) on each benchmark is in bold. Note that average token usage could exceed the available context window, since it includes context tokens for all queries. A dash (–) denotes where the model failed to generate responses due to exceeding the context window limit.

Strategy	LLM Backbone	EHRQA		FHIR-AgentBench		MedAgentBench	
		Acc \uparrow	Tok \downarrow	Acc \uparrow	Tok \downarrow	Acc \uparrow	Tok \downarrow
ONTOLOGY-GUIDED	GEMINI 3 FLASH	0.839 \pm 0.010	13,748	0.521 \pm 0.019	227,060	0.412 \pm 0.078	79,949
	CLAUDE SONNET 4.6	0.857 \pm 0.010	11,769	0.552 \pm 0.018	190,705	0.412 \pm 0.078	65,295
	GPT-5 MINI	0.835 \pm 0.010	10,680	0.457 \pm 0.018	174,843	0.365 \pm 0.081	60,077
FHIR2TEXT	GEMINI 3 FLASH	0.883 \pm 0.009	21,325	0.453 \pm 0.018	614,787	0.514 \pm 0.084	345,232
	CLAUDE SONNET 4.6	0.893 \pm 0.008	18,757	0.476 \pm 0.018	556,485	0.500 \pm 0.084	292,471
	GPT-5 MINI	0.887 \pm 0.009	17,498	0.294 \pm 0.016	518,598	0.419 \pm 0.078	277,005
PREFILTERED	GEMINI 3 FLASH	0.679 \pm 0.013	63,978	0.539 \pm 0.018	529,171	0.865 \pm 0.054	155,893
	CLAUDE SONNET 4.6	0.664 \pm 0.013	60,698	0.529 \pm 0.019	496,584	0.865 \pm 0.054	147,948
	GPT-5 MINI	0.678 \pm 0.013	49,849	–	–	0.777 \pm 0.071	123,823
FLOWSHEET	GEMINI 3 FLASH	0.819 \pm 0.010	11,804	0.446 \pm 0.039	418,423	0.939 \pm 0.037	221,666
	CLAUDE SONNET 4.6	0.838 \pm 0.010	10,438	0.496 \pm 0.017	378,380	0.878 \pm 0.051	184,698
	GPT-5 MINI	0.829 \pm 0.009	9,129	0.361 \pm 0.016	331,548	0.824 \pm 0.061	161,991
FHIR-HOPPER (OURS)	GEMINI 3 FLASH	0.902 \pm 0.008	24,998	0.568 \pm 0.018	22,984	0.959 \pm 0.030	26,196
	CLAUDE SONNET 4.6	0.897 \pm 0.008	20,935	0.611 \pm 0.018	20,978	0.946 \pm 0.037	20,603
	GPT-5 MINI	0.810 \pm 0.011	18,761	0.456 \pm 0.018	16,589	0.912 \pm 0.044	18,301

scaling trade-off on the short, synthetic EHRQA bundles, where FHIR-Hopper consumes more tokens than the Flow-sheet and Ontology-guided baselines. This occurs because our agentic framework incurs a fixed initialization overhead (tool definitions, system prompts, and episode headers). On trivial records, this fixed overhead outweighs the savings from graph pruning. However, as patient history grows $\mathcal{O}(N)$, baseline context sizes explode linearly, whereas FHIR-Hopper’s budgeted linearization bounds token usage to a stable $\mathcal{O}(1)$ footprint, initialized by the SBL budget, allowing it to scale better to massive clinical histories.

4 Discussion

In this work, we proposed the FHIR-Hopper, an agentic approach that leverages graph representations of FHIR bundles to achieve state-of-the-art performance on clinical EHR question-answering tasks. Our results demonstrate that the FHIR-Hopper not only outperforms other harnesses (FHIR2Text, Ontology-Guided Retrieval, Prefiltered) and serialization baselines in accuracy, but also does so with a significantly reduced context token footprint. The recent momentum in enterprise health AI, e.g. the launch of *GPT Health*, *Perplexity Health*, and *Verily Me*, highlights an unprecedented opportunity to interface LLMs directly with EHRs. However, sustaining this interaction securely and affordably requires moving beyond naive data serialization. FHIR-Hopper provides a scalable blueprint for such integration. By projecting FHIR bundles into navigable graphs, we demonstrate that LLMs can achieve state-of-the-art clinical

reasoning while shrinking the context footprint to a stable $\sim 20,000$ tokens, even for complex queries (e.g. in FHIR AgentBench).

This token efficiency yields critical operational advantages. **Privacy and Cost:** Transmitting only a budgeted graph skeleton and dynamically retrieved nodes to cloud inference APIs minimizes the exposure of sensitive Protected Health Information (PHI) and drastically cuts token costs. **Auditability:** Unlike monolithic black-box LLMs, FHIR-Hopper produces a deterministic trace of traversed nodes and invoked tools, establishing a transparent reasoning path essential for trust.

While optimized for question answering, this neuro-symbolic architecture naturally extends to **administrative workflows** (automating billing and prior authorizations), **patient-facing tools** (synthesizing longitudinal care plans), **data quality and maintenance** (identifying data inconsistencies and data gaps), and **biomedical research** (structuring EHR cohorts).

Limitations & Trade-offs. First, multi-turn tool calling introduces sequential inference latency, presenting a trade-off against single-shot context stuffing for real-time decision support. Second, the framework relies heavily on FHIR reference integrity; sparse databases will degrade graph traversal. Finally, our current implementation focuses on structured tabular data, currently treating other resources such as unstructured clinical notes (DocumentReference) and medical imaging as opaque nodes.

References

- ALMutairi, M., AlKulaib, L., Wang, S., Chen, Z., ALMutairi, Y., Alenazi, T. M., Luther, K., and Lu, C.-T. Fhirviz: Multi-agent platform for fhir visualization to advance healthcare analytics. In *Proceedings of the 15th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, BCB '24*, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400713026. doi: 10.1145/3698587.3701392. URL <https://doi.org/10.1145/3698587.3701392>.
- Amugongo, L. M., Mascheroni, P., Brooks, S., Doering, S., and Seidel, J. Retrieval augmented generation for large language models in healthcare: A systematic review. *PLOS Digital Health*, 4(6):e0000877, 2025.
- Anthropic. Claude sonnet 4.6 system card, 2026. URL <https://anthropic.com/claude-sonnet-4-6-system-card>.
- Bender, D. and Sartipi, K. HI7 fhir: An agile and restful approach to healthcare information exchange. In *Proceedings of the 26th IEEE International Symposium on Computer-Based Medical Systems*, pp. 326–331, 2013. doi: 10.1109/CBMS.2013.6627810.
- Cao, L., Chen, Q., and Guo, Y. EHR-RAG: Bridging long-horizon structured electronic health records and large language models via enhanced retrieval-augmented generation, 2026. URL <https://arxiv.org/abs/2601.21340>.
- DeepMind, G. Gemini 3.0 flash preview system card, 2025. URL <https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Flash-Model-Card.pdf>.
- Edge, D., Trinh, H., Cheng, N., Bradley, J., Chao, A., Mody, A., Truitt, S., Metropolitansky, D., Ness, R. O., and Larson, J. From local to global: A graph rag approach to query-focused summarization, 2025. URL <https://arxiv.org/abs/2404.16130>.
- Grimes, J., Brush, R., Rhyzhikov, N., Szul, P., Mandel, J., Gottlieb, D., Grieve, G., Sadjad, B., and Sanyal, A. SQL on FHIR - tabular views of FHIR data using FHIRPath. *npj Digital Medicine*, 8(1):342, 2025. doi: 10.1038/s41746-025-01708-w. URL <https://doi.org/10.1038/s41746-025-01708-w>.
- Heydari, A. A., Gu, K., Srinivas, V., Yu, H., Zhang, Z., Zhang, Y., Paruchuri, A., He, Q., Palangi, H., Hammerquist, N., Metwally, A. A., Winslow, B., Kim, Y., Ayush, K., Yang, Y., Narayanswamy, G., Xu, M. A., Garrison, J., Lee, A. A., Vafeiadou, J., Graef, B., Galatzer-Levy, I. R., Schenck, E., Barakat, A., Perez, J., Shreibati, J., Hernandez, J., Faranesh, A. Z., Prieto, J. L., Heneghan, C., Liu, Y., Zhan, J., Malhotra, M., Patel, S., Althoff, T., Liu, X., McDuff, D., and Xu, X. O. The anatomy of a personal health agent, 2025. URL <https://arxiv.org/abs/2508.20148>.
- Idrissi-Yaghir, A., Arzideh, K., Schäfer, H., Eryilmaz, B., Bahn, M., Wen, Y., Borys, K., Hartmann, E., Schmidt, C., Pelka, O., et al. Using a diverse test suite to assess large language models on fast health care interoperability resources knowledge: Comparative analysis. *Journal of Medical Internet Research*, 27:e73540, 2025.
- Jiang, J., Zhou, K., Zhao, W. X., Song, Y., Zhu, C., Zhu, H., and Wen, J.-R. Kg-agent: An efficient autonomous agent framework for complex reasoning over knowledge graph, 2024. URL <https://arxiv.org/abs/2402.11163>.
- Jiang, Y., Black, K. C., Geng, G., Park, D., Zou, J., Ng, A. Y., and Chen, J. H. Medagentbench: A virtual ehr environment to benchmark medical llm agents. *NEJM AI*, 2(9):AIdbp2500144, 2025. doi: 10.1056/AIdbp2500144. URL <https://ai.nejm.org/doi/full/10.1056/AIdbp2500144>.
- Johnson, A. E. W., Bulgarelli, L., Shen, L., Gayles, A., Shammout, A., Horng, S., Pollard, T. J., Hao, S., Moody, B., Gow, B., Lehman, L.-w. H., Celi, L. A., and Mark, R. G. Mimic-iv, a freely accessible electronic health record dataset. *Scientific Data*, 10(1):1, 2023. doi: 10.1038/s41597-022-01899-x. URL <https://doi.org/10.1038/s41597-022-01899-x>.
- Kabak, Y., Erturkmen, G. B. L., Gencturk, M., Namli, T., Sinaci, A. A., Corcoles, R. A., Ballesteros, C. G., Abizanda, P., and Dogac, A. FHIR-RAG-MEDS: Integrating HL7 FHIR with retrieval-augmented large language models for enhanced medical decision support, 2026. URL <https://arxiv.org/abs/2509.07706>.
- Lee, G., Hwang, H., Bae, S., Kwon, Y., Shin, W., Yang, S., Seo, M., Kim, J.-Y., and Choi, E. EHRSQL: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35: 15589–15601, 2022.
- Lee, G., Bach, E., Yang, E., Pollard, T., Johnson, A., Choi, E., Lee, J. H., et al. Fhir-agentbench: Benchmarking llm agents for realistic interoperable ehr question answering. *arXiv preprint arXiv:2509.19319*, 2025a.
- Lee, J., Chen, F., Dua, S., Cer, D., Shanbhogue, M., Naim, I., Ábrego, G. H., Li, Z., Chen, K., Vera, H. S., Ren, X., Zhang, S., Salz, D., Boratko, M., Han, J., Chen, B., Huang, S., Rao, V., Suganthan, P., Han, F., Doumanoglou, A., Gupta, N., Moiseev, F., Yip, C., Jain, A., Baumgartner,

- 275 S., Shahi, S., Gomez, F. P., Mariserla, S., Choi, M., Shah,
276 P., Goenka, S., Chen, K., Xia, Y., Chen, K., Duddu, S.
277 M. K., Chen, Y., Walker, T., Zhou, W., Ghiya, R., Gle-
278 icher, Z., Gill, K., Dong, Z., Seyedhosseini, M., Sung,
279 Y., Hoffmann, R., and Duerig, T. Gemini embedding:
280 Generalizable embeddings from gemini, 2025b. URL
281 <https://arxiv.org/abs/2503.07891>.
- 282 Li, Y., Wang, H., Yerebakan, H. Z., Shinagawa, Y., and Luo,
283 Y. FHIR-GPT enhances health interoperability with large
284 language models. *NEJM AI*, 1(8):A1cs2300301, 2024.
285 doi: 10.1056/A1cs2300301. URL <https://ai.nejm.org/doi/full/10.1056/A1cs2300301>.
- 286 Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua,
287 M., Petroni, F., and Liang, P. Lost in the middle: How
288 language models use long contexts. *Transactions of*
289 *the Association for Computational Linguistics*, 12:157–
290 173, 2024. doi: 10.1162/tacl.a.00638. URL <https://aclanthology.org/2024.tacl-1.9/>.
- 291 Makhni, S., Cerrato, P., Rico, J., Niazi, S., O’Horo, J.,
292 Peters, S., Shah, V., and Halamka, J. Meeting the
293 challenges of electronic health record (ehr) optimiza-
294 tion. *NPJ Digit Med*, 9(1):8, Dec 2025. ISSN 2398-
295 6352 (Electronic); 2398-6352 (Linking). doi: 10.1038/
296 s41746-025-02178-w.
- 297 OpenAI. Openai gpt-5 system card, 2025. URL <https://arxiv.org/abs/2601.03267>.
- 298 Saab, K., Tu, T., Weng, W.-H., Tanno, R., Stutz, D., Wul-
299 czyn, E., Zhang, F., Strother, T., Park, C., Vedadi, E.,
300 Chaves, J. Z., Hu, S.-Y., Schaekermann, M., Kamath,
301 A., Cheng, Y., Barrett, D. G. T., Cheung, C., Mustafa,
302 B., Palepu, A., McDuff, D., Hou, L., Golany, T., Liu,
303 L., baptiste Alayrac, J., Hounsby, N., Tomasev, N., Frey-
304 berg, J., Lau, C., Kemp, J., Lai, J., Azizi, S., Kanada, K.,
305 Man, S., Kulkarni, K., Sun, R., Shakeri, S., He, L., Caine,
306 B., Webson, A., Latysheva, N., Johnson, M., Mansfield,
307 P., Lu, J., Rivlin, E., Anderson, J., Green, B., Wong,
308 R., Krause, J., Shlens, J., Dominowska, E., Eslami, S.
309 M. A., Chou, K., Cui, C., Vinyals, O., Kavukcuoglu, K.,
310 Manyika, J., Dean, J., Hassabis, D., Matias, Y., Web-
311 ster, D., Barral, J., Corrado, G., Sementurs, C., Mahdavi,
312 S. S., Gottweis, J., Karthikesalingam, A., and Natarajan,
313 V. Capabilities of gemini models in medicine, 2024. URL
314 <https://arxiv.org/abs/2404.18416>.
- 315 Schick, T., Dwivedi-Yu, J., Dessi, R., Raileanu, R., Lomeli,
316 M., Zettlemoyer, L., Cancedda, N., and Scialom, T.
317 Toolformer: Language models can teach themselves to
318 use tools, 2023. URL <https://arxiv.org/abs/2302.04761>.
- 319 Schmiedmayer, P., Rao, A., Zagar, P., Aalami, L., Ravi,
320 V., Zahedivash, A., han Yao, D., Fereydooni, A.,
321 and Aalami, O. Llmomfhir. *JACC: Advances*, 4
322 (6_Part.1):101780, 2025. doi: 10.1016/j.jacadv.2025.
323 101780. URL <https://www.jacc.org/doi/abs/10.1016/j.jacadv.2025.101780>.
- 324 Shmatko, A., Jung, A. W., Gaurav, K., Brunak, S.,
325 Mortensen, L. H., Birney, E., Fitzgerald, T., and Gerstung,
326 M. Learning the natural history of human disease with
327 generative transformers. *Nature*, 647(8088):248–256,
328 2025. doi: 10.1038/s41586-025-09529-3. URL <https://doi.org/10.1038/s41586-025-09529-3>.
- 329 Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung,
H. W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S.,
Payne, P., Seneviratne, M., Gamble, P., Kelly, C., Babiker,
A., Schärli, N., Chowdhery, A., Mansfield, P., Demner-
Fushman, D., Agüera y Arcas, B., Webster, D., Corrado,
G. S., Matias, Y., Chou, K., Gottweis, J., Tomasev, N.,
Liu, Y., Rajkomar, A., Barral, J., Sementurs, C., Karthike-
salingam, A., and Natarajan, V. Large language models
encode clinical knowledge. *Nature*, 620(7972):172–180,
2023. doi: 10.1038/s41586-023-06291-2. URL <https://doi.org/10.1038/s41586-023-06291-2>.
- Tang, Y. and Yang, Y. Multihop-RAG: Benchmarking
retrieval-augmented generation for multi-hop queries.
arXiv preprint arXiv:2401.15391, 2024.
- Tomaszuk, D., Smajevic, A., Sagi, T., and Hose, K. Fhir
lens: A graph-based approach to semantic ehr explo-
ration. In *2025 IEEE 38th International Symposium
on Computer-Based Medical Systems (CBMS)*, pp. 1–6,
2025a. doi: 10.1109/CBMS65348.2025.00133.
- Tomaszuk, D., Smajevic, A., Sagi, T., and Hose, K. Fhir
lens: A graph-based approach to semantic ehr explo-
ration. In *2025 IEEE 38th International Symposium
on Computer-Based Medical Systems (CBMS)*, pp. 1–6,
2025b. doi: 10.1109/CBMS65348.2025.00133.
- Walonoski, J., Kramer, M., Nichols, J., Quina, A., Moe-
sel, C., Hall, D., Duffett, C., Dube, K., Gallagher, T.,
and McLachlan, S. Synthea: An approach, method,
and software mechanism for generating synthetic pa-
tients and the synthetic electronic health care record.
*Journal of the American Medical Informatics Associa-
tion*, 25(3):230–238, 03 2018. ISSN 1527-974X. doi:
10.1093/jamia/ocx079. URL <https://doi.org/10.1093/jamia/ocx079>.
- Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang,
L., Zhang, X., Zhang, S., Liu, J., Awadallah, A. H., White,
R. W., Burger, D., and Wang, C. Autogen: Enabling next-
gen llm applications via multi-agent conversation, 2023.
URL <https://arxiv.org/abs/2308.08155>.

330 Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan,
331 K., and Cao, Y. React: Synergizing reasoning and act-
332 ing in language models, 2023. URL [https://arxiv.](https://arxiv.org/abs/2210.03629)
333 [org/abs/2210.03629](https://arxiv.org/abs/2210.03629).

334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384

A Appendix

A.1 Fast Healthcare Interoperability Resources (FHIR)

Established by HL7 International in 2011, the Fast Healthcare Interoperability Resources (FHIR) framework serves as the preeminent standard for healthcare data storage and exchange. The schema is predicated on several core components:

- **Resources:** Modular primitives representing distinct clinical entities (e.g., `Patient`, `Condition`, `Encounter`, `Medication`, `Observation`).
- **References:** Directional edges that map relationships from a source resource to a target resource, establishing the overarching graph topology.
- **Patient:** The root node for most clinical trajectories, encapsulating administrative and demographic metadata (e.g., identifiers, gender, birth date).
- **Encounter:** Represents a specific patient-provider interaction, establishing the spatiotemporal context (e.g., inpatient admission, emergency visit) that anchors subsequent clinical events.
- **Condition:** Documents diagnoses, health concerns, or problems, tracking the longitudinal severity and clinical status (e.g., active, relapse, remission) of a patient's ailments.
- **Observation:** The primary mechanism for recording measurements and assertions, spanning vital signs, laboratory assays, and social history.
- **Medication:** Represents a specific pharmacological substance or mixture, typically instantiated contextually via linked `MedicationRequest` or `MedicationAdministration` resources.

Simplified Patient FHIR bundle

```
[
  {
    "resourceType": "Patient",
    "id": "0a8eebfd-...",
    "gender": "female",
    "birthDate": "2128-05-06",
    "extension": [
      // ... Nested custom demographic attributes (e.g., Race, Ethnicity)
      { "url": ".../us-core-race", "valueString": "White" }
    ]
  },
  {
    "resourceType": "Condition",
    "id": "b9ea8d38-...",
    "subject": { "reference": "Patient/0a8eebfd-..." }, // Relational link to
      Patient
    "encounter": { "reference": "Encounter/6c68c032-..." },
    "code": {
      "coding": [{
        "system": ".../mimic-diagnosis-icd9", // Standardized medical ontology
        "code": "78959",
        "display": "Other ascites"
      }]
    }
  },
  {
    "resourceType": "Encounter",
    "id": "18f74cab-...",
    "subject": { "reference": "Patient/0a8eebfd-..." },
    "period": {
      "start": "2180-07-23T14:00:00-04:00",
      "end": "2180-07-23T23:50:47-04:00"
    }
  }
] // ... {3,607 resources}
```

Figure S1. Simplified Patient FHIR bundle.

Unlike flat tabular architectures, FHIR structures a patient’s longitudinal history as a directed graph of JSON objects interconnected via References (Figures S1 and S2). Consequently, Resources are deeply nested; reconstructing a coherent clinical sequence typically requires relational joins and temporal grouping by Encounter. Further specifications are detailed in the official FHIR documentation (<https://www.hl7.org/fhir/>).

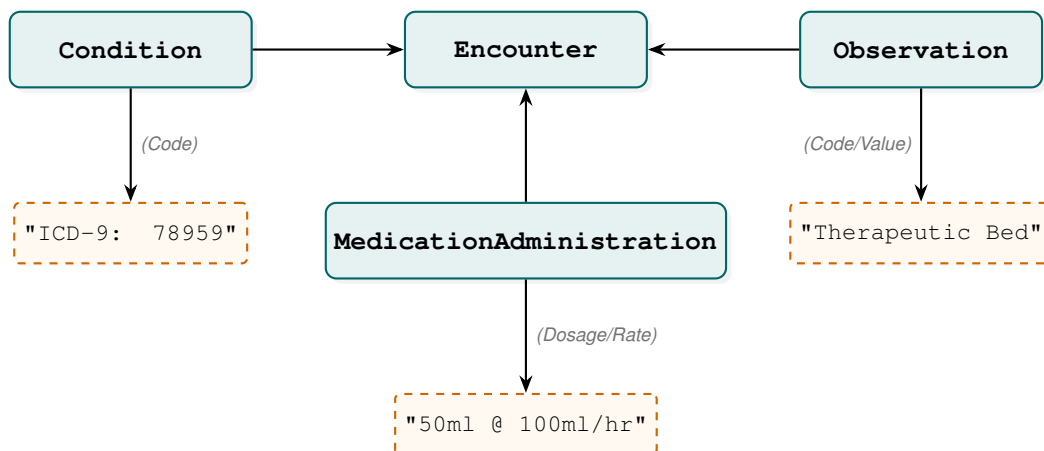


Figure S2. **Graph structure of a simplified FHIR bundle.** Resources such as Condition, Observation, and MedicationAdministration reference their parent Encounter, forming a directed graph.

A.2 Related Work

Tool-using LLM agents. A line of recent work equips LLMs with external tools and lets them interleave reasoning with tool calls, including ReAct (Yao et al., 2023), Toolformer (Schick et al., 2023), AutoGen (Wu et al., 2023). Similar to these systems, FHIR-Hopper exposes deterministic tools and uses standard function-calling APIs, but the tools are specialized to graph traversal over a structured clinical representation rather than general web or code execution.

Graph and structure-augmented retrieval. GraphRAG (Edge et al., 2025) and related approaches (Jiang et al., 2024) build a knowledge graph from a corpus and let an LLM traverse or condition on it during retrieval. Our approach is closest in spirit: we treat the FHIR bundle itself as a typed, time-stamped graph and provide both a budgeted linearization and traversal tools. The differences are domain-specific: edges and node types come from the FHIR schema rather than open-domain extraction, and the graph is anchored on an episodic spine rather than a flat similarity graph.

LLMs over FHIR and EHR. A number of systems consume FHIR data with LLMs. LLMs-On-FHIR (Schmiedmayer et al., 2025) relies on heuristic prefiltering and de-duplication; FHIR-RAG-Meds (Kabak et al., 2026) performs ontology-guided retrieval; FHIR-AgentBench (Lee et al., 2025a) and MedAgentBench (Jiang et al., 2025) construct agentic benchmarks in which LLMs call FHIR APIs directly. These systems serve as our baselines or evaluation targets. Other clinical LLMs such as Med-PaLM (Singhal et al., 2023) and Med-Gemini (Saab et al., 2024) are trained or instruction-tuned on clinical text, but do not specifically address structured-resource reasoning over FHIR bundles.

A.3 Datasets

A.3.1 MIMIC-IV FHIR DEMO (FROM FHIR AGENTBENCH)

This benchmark contains 2,931 clinician-sourced questions and answers, originally from the EHRSQL (Lee et al., 2022) dataset grounded in de-identified MIMIC-IV (Johnson et al., 2023) FHIR patient records. Each patient record contains roughly 9K FHIR resources on average, and some questions require reasoning over more than 2,000 relevant resources. MIMIC-IV is a publicly available database from the Beth Israel Deaconess Medical Center covering emergency department and ICU admissions between 2008-2019. The MIMIC-IV Clinical Database Demo on FHIR is a subset of these data that was converted into FHIR R4 specifications for 100 randomly selected patients (94 of which are present in FHIR AgentBench). FHIR resources included patient, organization, observations from specimen samples, medications, charted observations, and billing. Free-text clinical notes were excluded.

The queries in FHIR AgentBench are focused on single patient natural language questions. The benchmark tests agents’

ability to perform multi-step retrieval and reasoning over nested FHIR data, including temporal logic, cross-resource reference traversal, case-sensitive terminology handling, and queries with no relevant matches. We show the token count distribution of the FHIR bundles present in this dataset in Figure S3.

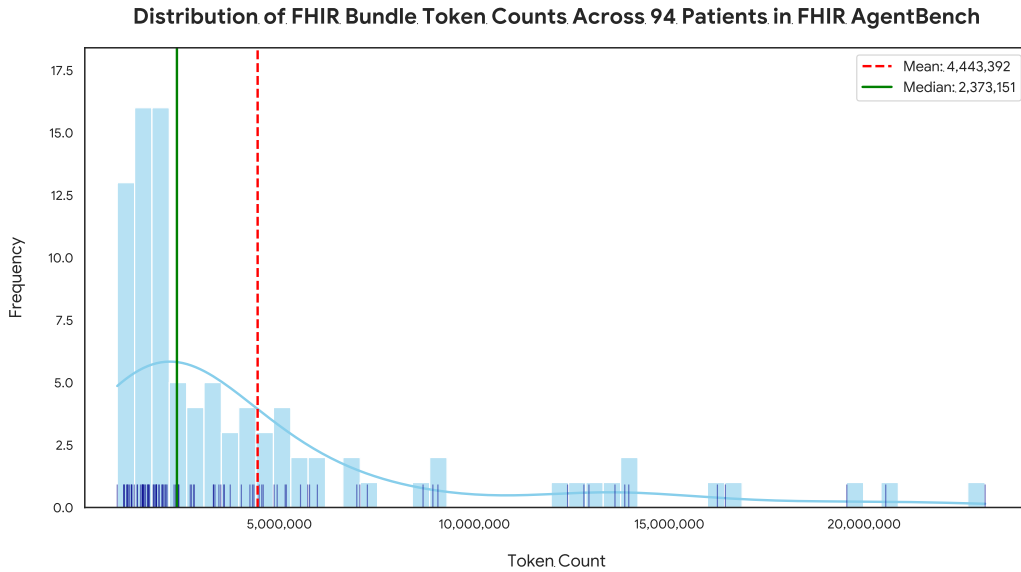


Figure S3. Distribution of Tokens (computed through Gemini’s tokenizer) for bundles in the FHIR AgentBench.

A.3.2 MEDAGENTBENCH

MedAgentBench (Jiang et al., 2025) consists of 300 clinically-relevant and verifiable tasks from 10 categories written by licensed human clinicians. The profiles represent 100 patients with over 700,000 individual clinical data elements, including laboratory results, vital signs, procedures, diagnoses, and medication orders. Given the scope of our work, we subset the dataset to include *all retrieval* queries (i.e. only patient information retrieval, laboratory result retrieval, and patient data aggregation; “POST” tasks such as *Test ordering*, *Medication ordering*, *Referral ordering*, and *Recording patient data* were excluded), resulting in 148 tasks. We present the token count distribution of the FHIR bundles, for the retrieval queries, in Figure S4.

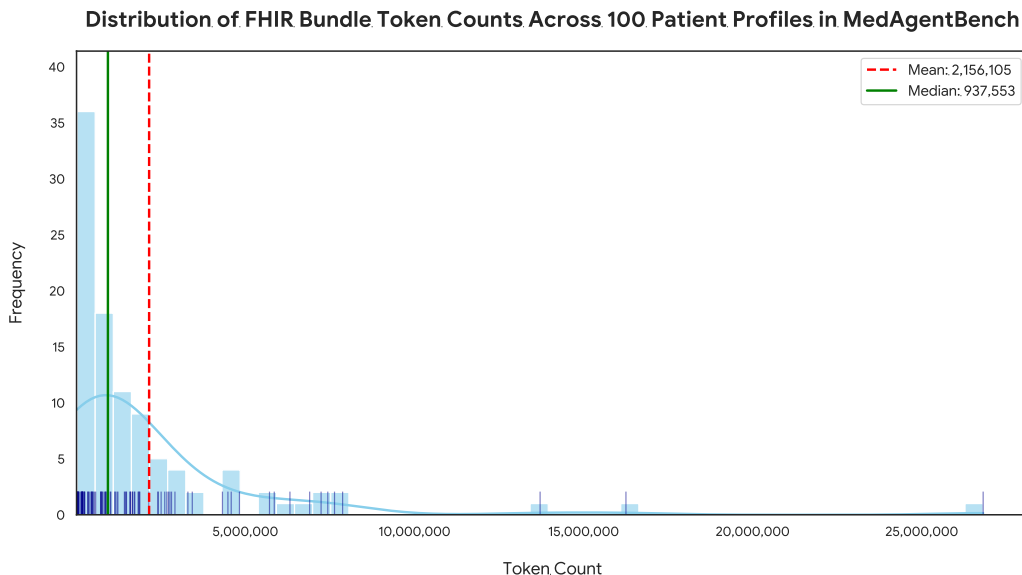


Figure S4. Distribution of Tokens (computed through Gemini’s tokenizer) for bundles in the MedAgentBench dataset.

A.3.3 EHRQA

We curated this dataset for granular testing and validation on a simple benchmark. This dataset contains 5,133 LLM-generated multiple-choice question-answer pairs for 20 synthetic FHIR bundles produced by Synthea (Walonoski et al., 2018). All synthetic bundles fit within a 32k context window, and clinical notes were not included in EHRQA. We provide an example of an entry in this dataset in below.

EHRQA Example

```

QAEExample (
  patient_id='78a7ab19-bdd0-2238-e4e5-de757ee074d1', question='What past nasal
  sinus condition does the patient have?\n',
  answer_choices=['(A) acute empyema of nasal sinus', '(B) fungal sinusitis', '(C)
  chronic sinusitis', '(D) viral sinusitis'],
  correct_answer='D', ehr_fact=ConditionExistenceInHistoryFact (fact_type=<
  ConditionFactType.EXISTENCE_IN_HISTORY: 'existence_in_history'>,
  fact_variation_type=<ConditionFactVariationType.FINDING_SITE: 'finding_site'>,
  evidence=[],
  evidence_refs=['Condition/6ba0a993-8f90-f024-dc46-89243cb67d17'],
  code='2095001',
  code_type=<MedicalCodeset.SNOMED: 'http://snomed.info/sct'>, is_transitive=False
  ,
  intermediate_codes=())
  
```

Because both the bundles and the questions are synthetic, EHRQA exercises only basic single- or two-hop retrieval and exhibits ceiling effects on competent baselines; we report it primarily to verify that strategies do not regress on simpler tasks. Moreover, we wanted to show that while some approaches may achieve promising results on such synthetic benchmarks, they may exhibit significant performance drop once they are put to test in real world setting, as evident by performance results on FHIR AgentBench and MedAgentBench.

The complete processing recipe and the dataset, including the generation pipeline, synthetic personas, and the query set, are available to download at (GITHUB LINK WILL BE PROVIDED UPON DE-ANONYMIZATION / ACCEPTANCE).

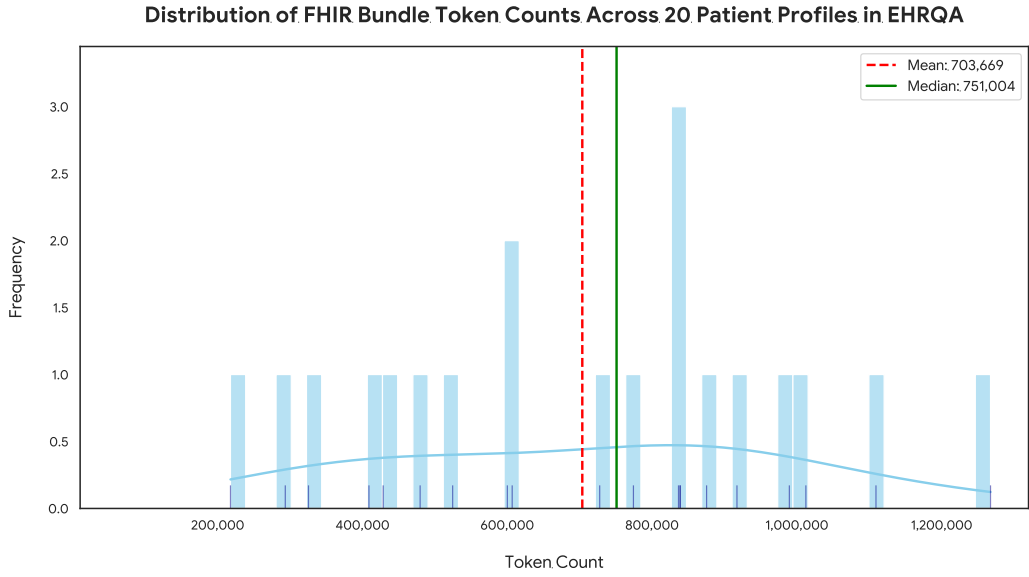


Figure S5. Distribution of Tokens (computed through Gemini's tokenizer) for bundles in the EHRQA benchmark.

A.4 FHIR-Hopper Implementation Details

A.4.1 PSEUDOCODE FOR EPISODIC GRAPH PROJECTION

Algorithm 1 Episodic Graph Projection & Latent Episode Inference

Require: Resource bundle \mathcal{B} , timestamp function $\tau(v)$, encounter-reference function $\text{ref}(v)$, temporal window δ

Ensure: Chronological episodic spine $\mathcal{E}_{\text{spine}}$

```

1:  $\mathcal{H}_{\text{explicit}} \leftarrow \emptyset$  {Map of Encounter ID  $\rightarrow$  explicit episode nodes}
2:  $\mathcal{H}_{\text{latent}} \leftarrow \emptyset$  {Map of  $\delta$ -snapped timestamp  $\rightarrow$  latent episode nodes}
3:  $\mathcal{O} \leftarrow \emptyset$  {Queue of non-Encounter resources}
4: // Phase 1: initialize backbone from explicit Encounters
5: for  $r \in \mathcal{B}$  do
6:   if  $r.\text{type} == \text{'Encounter'}$  then
7:     Initialize episode node  $E$  with spine resource  $r$ 
8:      $\mathcal{H}_{\text{explicit}}[r.\text{id}] \leftarrow E$ 
9:   else
10:     $\mathcal{O}.\text{enqueue}(r)$ 
11:   end if
12: end for
13: // Phase 2: assign orphans, falling back to latent episodes
14: for  $v \in \mathcal{O}$  do
15:    $id_{\text{ref}} \leftarrow \text{ref}(v)$  {Encounter reference if present}
16:   if  $id_{\text{ref}} \neq \text{null} \wedge id_{\text{ref}} \in \mathcal{H}_{\text{explicit}}$  then
17:      $\mathcal{H}_{\text{explicit}}[id_{\text{ref}}].\text{add}(v)$ 
18:   else
19:      $t \leftarrow \tau(v)$ 
20:     if  $t == \text{null}$  then continue
21:      $E^* \leftarrow$  nearest preceding Encounter with start time in  $[t - \delta, t]$ 
22:     if  $E^* \neq \text{null}$  then
23:        $\mathcal{H}_{\text{explicit}}[E^*.\text{id}].\text{add}(v)$ 
24:     else
25:        $t' \leftarrow$  snap  $t$  to a grid of width  $\delta$ 
26:       if  $t' \notin \mathcal{H}_{\text{latent}}$  then
27:         Create new latent episode  $E_{\text{latent}}$  anchored at  $t'$ 
28:          $\mathcal{H}_{\text{latent}}[t'] \leftarrow E_{\text{latent}}$ 
29:       end if
30:        $\mathcal{H}_{\text{latent}}[t'].\text{add}(v)$ 
31:     end if
32:   end if
33: end for
34: // Phase 3: chronological merge
35:  $\mathcal{E}_{\text{union}} \leftarrow \text{Values}(\mathcal{H}_{\text{explicit}}) \cup \text{Values}(\mathcal{H}_{\text{latent}})$ 
36:  $\mathcal{E}_{\text{spine}} \leftarrow \text{SortByTime}(\mathcal{E}_{\text{union}})$  {Chronologically ordered episodic spine}
37: return  $\mathcal{E}_{\text{spine}}$ 

```

A.4.2 SYSTEM PROMPT

System Prompt

Role

You are a Clinical Graph Navigator operating on a chronological hypergraph of patient data. Providing medical advice or a diagnosis is outside the scope of your role.

Objective:

Your goal is to provide your best answer to the user's query by navigating this graph and utilizing the tools at your disposal. Note that you are a query system, not a chat model, and should not converse with the user by asking follow-up questions.

1. THE CONTEXTUAL CONSTRAINT

You are viewing a **Budget-Constrained Saliency Map** (the Skeleton).
 * **Around 99% of the data is HIDDEN** to fit in your context window.
 * **Relevance Scores** (`[0.92]`) show semantic similarity to your query.
 * **Gaps** (`[... GAP ...]`) indicate skipped timeframes with low relevance.

You will be provided with a budget-constrained Saliency Map and a user query. Make sure to consider the user query and the **"BUDGETED SKELETON (TOP HITS)"** that will show you the top hits in the graph that are most relevant to the query. However, it is likely that this representation may be missing some data that is relevant to the query. Therefore, it is crucial for you to carefully assess whether the information you see in the SKELETON is SUFFICIENT for you to answer the query, and also check the HIDDEN graph to see if there are any additional clues or data that may be relevant to best answer the user's query.

2. BUDGET WARNING (CRITICAL)

You have a strict maximum budget of `{{ max_llm_calls }}` steps. Each step can consist of multiple tool calls. You must gather all necessary information and provide your final answer before exceeding this limit. Plan your tool calls efficiently.

3. TEMPORAL AWARENESS (CRITICAL)

You should take into account the current time when answering questions including relative time specifiers, like "last month". For example, "last month" means the 30 days prior to current time. You can use the `'filter_time'` tool to query this specific date range.

4. CONTEXTUAL INFORMATION

You are provided with the following contextual information for your query:

```
{{ context_str }}
```

5. YOUR TOOLS

You must navigate this graph dynamically to find the truth.

- `'inspect_node(resource_id)'`: Fetches the full JSON content of a specific FHIR resource.
 - Usage**: "Read the file." Fetches the full JSON content.
 - Mandatory**: You CANNOT verify values (BP, Dates, Dosage) without this.
 - Args**:
 - `'resource_id'`: A specific ID from the skeleton (e.g., "Observation/123").

2. ``search_graph(keywords)``: Searches the FULL (hidden) graph for resources containing keywords.
- * **Usage:** "Search the dark." Scans the HIDDEN portion of the graph.
 - * **When to use:** The Skeleton is missing a key node.
 - * **Args:**
 - * ``keywords``: Clinical keywords (e.g., "Troponin", "Discharge Summary"). Also use this to find atemporal data like Patient demographics, Location details, etc.
3. ``follow_links(resource_id)``: Return the IDs and Types of resources directly referenced by the node.
- * **Usage:** "Traverse the graph based on references in resources." Returns the IDs of resources *referenced by* the target node.
 - * **When to use:** You have a node (e.g., MedicationRequest) and need to find its reason (reasonReference -> Condition) or its author (requester -> Practitioner), but those linked nodes are not visible in the current skeleton.
 - * **Args:**
 - * ``resource_id``: The ID of the node you want to trace FROM.
4. ``filter_time(start_date, end_date, resource_type, keywords)``: Filters graph nodes by time range, resource type, and keywords.
- * **Usage:** "Find all observations between 2023-01-01 and 2023-12-31"
 - * **When to use:** The skeleton does not have the chronological details you need, and you want to specifically query a time period. Crucial when the user specifies a timeline ("since August 2180" or "last two weeks"). Remember to frame the boundaries relative to simulation's end date: `{{ max_timestamp }}`.
 - * **Args:**
 - * ``start_date``: ISO format YYYY-MM-DD
 - * ``end_date``: ISO format YYYY-MM-DD
 - * ``resource_type``: The FHIR resource type to filter by (e.g., "Observation")
 - * ``keywords``: Keywords to search for in the node data

Note that THERE IS NO ``finish`` tool. You should NOT hallucinate on calling a ``finish`` tool.

Make sure to also refer to the docstring that will be provided to you for each tool to double check usage and arguments that you need to pass to each tool in case things have changed.

6. NAVIGATION STRATEGY

Follow this priority queue to solve the query:

- * **PHASE 1: SCAN (The Skeleton)**
 - * Look for High-Relevance (``[0.9+]``) nodes in the Skeleton.
 - * **Constraint:** If the question is about Medications, it may not be good to trust the skeleton's absence. Medications often have low semantic similarity to condition queries. Move to Phase 2.
- * **PHASE 2: EXPAND (The Neighborhood)**
 - * If you found a relevant node (e.g., a "Diabetes" Diagnosis) but lack details, use ``follow_links`` on it.
- * **PHASE 3: HUNT (The Void)**
 - * If the Skeleton and Links are empty, use ``search_graph`` or ``filter_time``.
- * **PHASE 4: VERIFY (The Grounding)**
 - * **CRITICAL:** Never guess a value. You must ``inspect_node`` to see the actual JSON numbers before answering.

A.5 Judge for Comparing Model Response against Ground Truth

Because models generate free-form text, exact-match grading is overly rigid. Since the tasks are verifiable with the provided ground truth, we use an LLM-as-a-judge for a simple correctness assessment: given the query, the ground truth, and the model response, the judge returns a correctness verdict and a short rationale. Because the QA tasks resolve to deterministic ground truth (specific values, dates, codes), the judge’s room for subjective interpretation is bounded. As a small validation study, we used humans to verify our auto-rater on a randomly selected subset of the FHIR-AgentBench ($n = 291$, $\sim 10\%$ of all queries), and found the LLM judge to not make any mistakes in assessing answer correctness.

Judge Prompt

You are a clinical auditor. Evaluate the Model Response against the Ground Truth.

Return a JSON object with two fields:

1. "is_correct": boolean (true if the model answer accurately answer the question that aligns with the ground truth)
2. "reason": string (brief explanation)

```

**Patient Query**: {user_query}
**Ground Truth Answer**: {ground_truth}
**Model Response**: {model_response}

```

JSON Output:

A.6 Model Hyperparameters

To assess generalizability, we evaluate every method against three frontier LLMs spanning different providers and cost/latency profiles:

- **Gemini 3 Flash** (DeepMind, 2025): a highly efficient model developed by Google DeepMind, optimized for speed and cost, with a large context window suitable for long clinical records.
- **Claude Sonnet 4.6** (Anthropic, 2026): a model from Anthropic offering strong reasoning capability with moderate inference cost.
- **GPT-5 Mini** (OpenAI, 2025): a lightweight OpenAI model with a smaller context window, used to evaluate robustness under tighter resource constraints.

All models use default top- p and a max output of 8,192 tokens. We set temperature to 0.6 for Gemini 3 Flash and Claude Sonnet 4.6; GPT-5 Mini uses its API-default 1.0. FHIR-Hopper uses the same temperature as its backbone LLM. Models are accessed via Google Vertex AI (Gemini 3 Flash, Claude Sonnet 4.6) and the OpenAI API (GPT-5 Mini). Saliency scoring uses the Gemini Embedding 1 model ($d = 3,072$) consistently across configurations to isolate the effect of the backbone LLM from the embedding. The judge is Claude Sonnet 4.6 with temperature 0.0 and max output 8,192 tokens.

A.7 Ablation Studies

Table S1. Ablation study. We ablate FHIR-Hopper by removing tool access (w/o TOOLS), forcing the agent to rely solely on the SBL view of the graph.

STRATEGY	BASE MODEL	EHRQA	FHIR-AGENTBENCH	MEDAGENTBENCH
FHIR-HOPPER	GEMINI 3 FLASH	0.902 ± 0.008	0.568 ± 0.018	0.959 ± 0.030
	CLAUDE SONNET 4.6	0.897 ± 0.009	0.611 ± 0.018	0.946 ± 0.037
	GPT-5 MINI	0.810 ± 0.011	0.456 ± 0.018	0.912 ± 0.044
W/O TOOLS	GEMINI 3 FLASH	0.844 ± 0.010	0.595 ± 0.017	0.784 ± 0.064
	CLAUDE SONNET 4.6	0.859 ± 0.010	0.591 ± 0.018	0.784 ± 0.064
	GPT-5 MINI	0.838 ± 0.010	0.233 ± 0.015	0.473 ± 0.084

We ablate FHIR-Hopper by removing tool access, forcing the agent to rely on the SBL view alone (Table S1). Tools are most valuable on EHRQA and MedAgentBench (e.g., MedAgentBench with GPT-5 Mini drops from 0.912 to 0.473) and for the weaker GPT-5 Mini backbone, where tools partially compensate for weaker base reasoning. On FHIR-AgentBench the effect is small and not in a single direction, suggesting the SBL view alone already surfaces most of the evidence needed there.

A.8 Additional Results

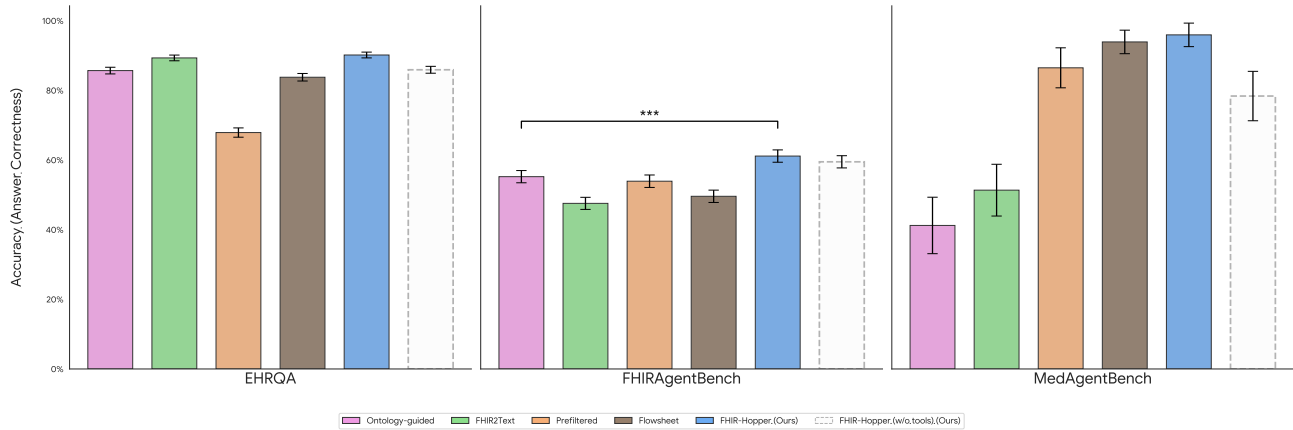


Figure S6. Accuracy of FHIR retrieval strategies across benchmarks. We compare FHIR-Hopper against the baseline strategies on EHRQA, FHIR-AgentBench, and MedAgentBench; for each strategy we report the configuration with the highest-performing backbone LLM. FHIR-Hopper attains the highest accuracy on every benchmark. Asterisks denote statistical significance ($p < 0.001$, McNemar’s test) of the improvement over the second-best strategy that is not ours.

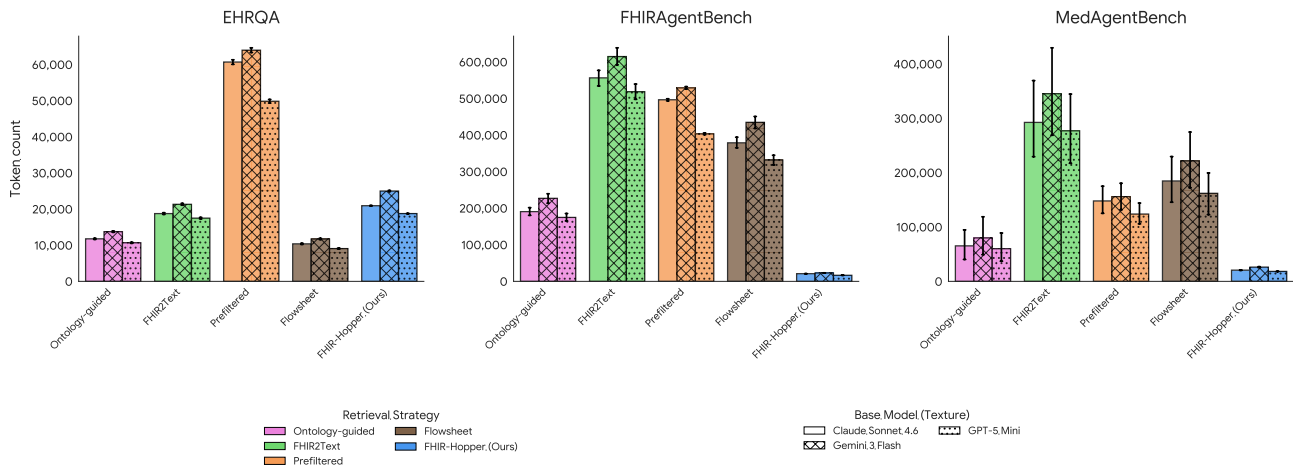


Figure S7. Token usage per strategy and LLM (lower is better). We report the average number of tokens consumed per query with the three LLMs used. FHIR-Hopper maintains a stable token footprint across all benchmarks, particularly on FHIR-AgentBench where baselines exceed context-window limits. Error bars indicate 95% confidence intervals.

A.9 Broader Impact and Ethics

Audibility: Unlike black-box LLMs, our agent produces a trace of nodes visited. This is crucial for clinical validation and trust. Additionally, by sending only relevant graph nodes to the LLM (rather than the full record), we minimize the data footprint exposed to cloud inference API.

990 The deployment of autonomous agents in clinical settings carries significant ethical, safety, and transparency implications.
991 Foremost among these is the risk of hallucinations, which could lead to severe patient harm if left unchecked. Therefore,
992 systems like the FHIR-Hopper must be deployed strictly as assistive tools rather than autonomous decision-makers.
993 Rigorous human-in-the-loop verification remains essential to ensure that AI-generated insights are clinically sound, safe,
994 and transparent before influencing patient care.
995
996
997
998
999

1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044