FROM CONVERSATION TO QUERY EXECUTION: BENCHMARKING USER AND TOOL INTERACTIONS FOR EHR DATABASE AGENTS

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031 032 033

034

037

038

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Despite the impressive performance of LLM-powered agents, their adoption for Electronic Health Record (EHR) data access remains limited by the absence of benchmarks that adequately capture real-world clinical data access flows. In practice, two core challenges hinder deployment: query ambiguity from vague user questions and value mismatch between user terminology and database entries. To address this, we introduce EHR-ChatQA, an interactive database question answering benchmark that evaluates the end-to-end workflow of database agents: clarifying user questions, using tools to resolve value mismatches, and generating correct SQL to deliver accurate answers. To cover diverse patterns of query ambiguity and value mismatch, EHR-ChatQA assesses agents in a simulated environment with an LLM-based user across two interaction flows: Incremental Query Refinement (IncreQA), where users add constraints to existing queries, and Adaptive Query Refinement (AdaptQA), where users adjust their search goals mid-conversation. Experiments with state-of-the-art LLMs (e.g., o4-mini and Gemini-2.5-Flash) over five i.i.d. trials show that while agents achieve high Pass@5 of 90–95% (at least one of five trials) on IncreQA and 60-80% on AdaptQA, their Pass^5 (consistent success across all five trials) is substantially lower by 35–60%. These results underscore the need to build agents that are not only performant but also robust for the safety-critical EHR domain. Finally, we provide diagnostic insights into common failure modes to guide future agent development.

1 Introduction

Large Language Models (LLMs) are increasingly operating as autonomous agents, interacting with external environments to solve complex tasks. One key application is interfacing with structured databases, which can substantially enhance data accessibility for non-technical users. This capability is particularly impactful in high-stakes domains such as Electronic Health Records (EHRs), where enabling natural language queries over vast patient data repositories has the potential to fundamentally transform both clinical research and patient care (Ohno-Machado, 2011; Yang et al., 2022). To assess such capabilities, the prevailing evaluation paradigm has relied on text-to-SQL benchmarks (e.g., Spider (Yu et al., 2018), MIMICSQL (Wang et al., 2020), EHRSQL (Lee et al., 2022), BIRD (Li et al., 2023)), which measure a model's ability to translate natural language questions into SQL queries. However, these benchmarks primarily emphasize the isolated task of mapping a single, well-formed question to SQL, a setting that fails to capture the complexities of how clinicians interact with EHRs.

To further investigate the gap between this evaluation paradigm and real-world clinical needs, we collaborated with a major academic medical center and reviewed internal text-to-SQL query logs. This analysis revealed that existing benchmarks are not suitable for capturing two core real-world challenges, which significantly hinders the adoption of these models: (1) *Query Ambiguity*¹: Users often pose vague questions that do not fully capture their true intent (e.g., a request such as "Show me recent labs" lacks specifics such as test types or time ranges) (Wang et al., 2023; Saparina & Lapata,

¹While linguistics distinguishes it from vagueness (a lack of specificity), we adopt the common NLP convention of using "ambiguity" to refer to any query requiring clarification.

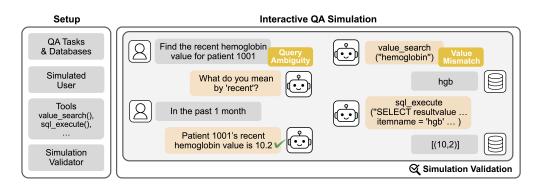


Figure 1: Overview of EHR-ChatQA. Our benchmark places an agent in a simulated environment with an LLM-powered user and tools. The agent must manage the full workflow: clarifying vague user questions, resolving terminology mismatches, generating and executing accurate SQL, and returning the final answer. Each interaction trace is validated for accurate evaluation.

2024; Dong et al., 2025). (2) *Value Mismatch*: Clinical terms in EHRs vary due to specialized database naming conventions and often differ from everyday language (e.g., "Lopressor" may not match the database entry "metoprolol tartrate"), creating a terminological gap (Holmes et al., 2021; De Mello et al., 2022) that must be bridged for accurate SQL generation.

Addressing these challenges requires moving beyond static, single-turn SQL generation to an environment where an agent can clarify a user's intent, invoke necessary tools to navigate complex EHR schemas and clinical values, and synthesize all relevant information to generate an accurate SQL query. To bridge this gap, we introduce *EHR-ChatQA*, an interactive database question answering (QA) benchmark designed to assess this end-to-end agentic workflow, from conversation to query execution, in the EHR domain. By placing agents in a simulated environment with both an LLM-based user and a suite of tools, EHR-ChatQA provides a holistic evaluation of agent capabilities in interactive clinical data access flows, with each interaction trace verified by a dedicated validator. Grounded in real-world clinical QA scenarios and two publicly available EHR databases (MIMIC-IV (Johnson et al., 2023) and eICU (Pollard et al., 2018)), our benchmark consists of tasks categorized into two different interaction flows designed for various query ambiguities and value mismatches: Incremental Query Refinement (IncreQA), which evaluates scenarios where users add new constraints to a query, and Adaptive Query Refinement (AdaptQA), which assesses an agent's ability to reformulate its plan when users' goals are modified mid-conversation.

Our evaluation of various state-of-the-art LLMs on EHR-ChatQA reveals a critical lack of robustness under diverse conversation paths. While agents often succeed in at least one of five attempts on a task (Pass@5), their ability to succeed consistently in all five attempts (Pass^5) is substantially lower. This performance gap exceeds 30% for IncreQA and 50% for AdaptQA. This inconsistency highlights a crucial lack of reliability in current agents, raising significant concerns for their deployment in safety-critical domains such as EHRs and pointing to key areas for future research.

The main contributions of our work are summarized as follows:

- We propose EHR-ChatQA, the first interactive benchmark for EHR QA that holistically evaluates agents' interactive, end-to-end workflows using simulated users and a set of customizable tools for schema exploration, value exploration, and web search.
- Grounded in real-world clinical QA scenarios and two publicly available EHR databases, the benchmark contains two interaction flows to reflect various query ambiguity and value mismatch patterns.
- Our evaluation of various LLMs reveals a critical performance gap between an agent's optimistic success (Pass@5) and its consistent success (Pass^5), providing diagnostic insights for developing more performant and reliable agents in interactive EHR QA.

2 RELATED WORK

Text-to-SQL Benchmarks Text-to-SQL research has largely focused on translating a single, well-defined question into an SQL query. Benchmarks such as Spider (Yu et al., 2018) and BIRD (Li

Table 1: Comparison of recent benchmarks categorized by core agent capabilities. EHR-ChatQA is the first benchmark to comprehensively evaluate database agents in the aspects ranging from conversational ability to effective tool use in the EHR domain. "Value Explor." indicates the mapping of user terminology to database entries (e.g., "WBC" \rightarrow "white blood cell count"). \triangle denotes resolution of ambiguity through SQL suggestions instead of user clarification.

Benchmark	Conversat	ional Ability	Tool-U	Domain	
Denemiai k	User Multi-turn	Query Ambiguity	Tool Use	Value Explor.	EHR
Spider (Yu et al., 2018)	×	×	×	×	×
SParC (Yu et al., 2019b)	\checkmark	×	×	×	×
CoSQL (Yu et al., 2019a)	\checkmark	\checkmark	×	×	×
EHRSQL (Lee et al., 2022)	×	×	×	×	\checkmark
BIRD (Li et al., 2023)	×	×	×	×	×
AgentBench (Liu et al.)	×	×	\checkmark	×	×
EHR-SeqSQL (Ryu et al., 2024)	\checkmark	×	×	×	\checkmark
PRACTIQ (Dong et al., 2025)	\checkmark	\triangle	×	×	×
Tau-Bench (Yao et al., 2025)	\checkmark	\checkmark	\checkmark	×	×
ToolDial (Shim et al., 2025)	\checkmark	\checkmark	\checkmark	×	×
MedAgentBench (Jiang et al., 2025)	×	×	\checkmark	×	\checkmark
MedAgentGym (Xu et al., 2025)	×	×	\checkmark	×	\checkmark
EHR-ChatQA (Ours)	√	√	√	√	√

et al., 2023) are prominent examples that have shaped this paradigm. While other benchmarks such as SParC (Yu et al., 2019b), CoSQL (Yu et al., 2019a), and PRACTIQ (Dong et al., 2025) introduced conversational context, their evaluation scope is often limited by predefined interaction patterns, such as requiring the model to generate an SQL query each turn or to expect a type of responses based on the script. This evaluation setting prevents them from capturing the open-ended and exploratory nature of realistic user interactions, thereby falling short of testing agents' ability to resolve query ambiguity and value mismatch. EHR-ChatQA is a benchmark dedicated to bridging this gap by explicitly evaluating agents' exploratory and interactive capabilities through user simulations.

EHR Question Answering The unique challenges of the medical domain have inspired several QA benchmarks for structured EHR data. Text-to-SQL benchmarks, from single-turn settings such as MIMICSQL (Wang et al., 2020) and EHRSQL (Lee et al., 2022) to the multi-turn EHR-SeqSQL (Ryu et al., 2024), have advanced querying on complex EHR schemas but lack support for resolving query ambiguity and value mismatch. More recently, agent-based benchmarks for EHRs such as MedAgentBench (Jiang et al., 2025) and MedAgentGym (Xu et al., 2025) tackle a broad range of clinical and biomedical tasks. However, their reliance on initial non-ambiguous task instructions still bypasses the need for dynamic, interactive resolution of query ambiguity and value mismatch, which are essential for real-world, interactive clinical QA. EHR-ChatQA aims to bridge this gap by explicitly requiring the simulated user to start from a vague question and allowing the agent to autonomously resolve the task.

Conversational and Tool-Using Agent Evaluation Evaluating LLM agents in dynamic environments has spurred progress in two key areas: tool-using agent benchmarks (Yao et al., 2022; Liu et al.), which assess instrumental competence, and frameworks for evaluating conversational skills in task-oriented dialogues, from large-scale curated datasets (Budzianowski et al., 2018; Rastogi et al., 2020) to dynamic user simulations (Sekulić et al., 2024). Recent works such as Tau-Bench (Yao et al., 2025) and ToolDial (Shim et al., 2025) combine these paradigms to evaluate the critical dual interaction loop between agents, users, and tools. EHR-ChatQA extends this framework to question answering over EHR databases. As shown in Table 1, although prior works focus on resolving query ambiguity and leveraging general tools for complex tasks, they do not specifically capture the challenges of interactive clinical data access, particularly value mismatch (Holmes et al., 2021; De Mello et al., 2022). EHR-ChatQA is designed to holistically address these challenges, grounded in real-world QA scenarios and EHR databases.

3 THE EHR-CHATQA BENCHMARK

3.1 TASK FORMULATION

The task instances in EHR-ChatQA can be formulated in the POMDP framework as (S, A, T, R, Ω) . This formulation captures the inherent uncertainty and the sequential nature of

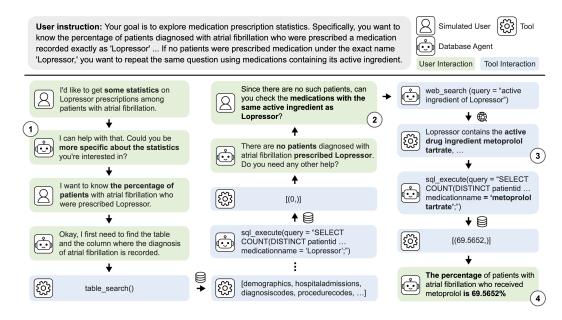


Figure 2: Example interaction trace in EHR-ChatQA (AdaptQA). Given a vague initial user question, the agent clarifies the request (1) by refining "some statistics" into the specific goal of "the percentage." When the initial query yields no matching records due to a value mismatch, the user guides the agent to search for the active ingredient instead (2). The agent then incorporates clinical knowledge retrieved via web_search() to reformulate the query (3), and finally executes the revised query to produce the correct answer (4).

translating an ambiguous user request into an executable query. The state $s \in \mathcal{S}$ is latent and includes all information relevant to solving the task: the user's true but unobserved intent (i.e., user instruction), the full conversation history, the contents of the EHR database, and any accessible external knowledge. The agent's action space \mathcal{A} models the dual interaction loop: (1) user interaction, including asking clarifying questions or providing answers, and (2) tool interaction, including invoking calls to explore the database schema, search relevant values, or access external web knowledge when necessary. After taking an action $a \in \mathcal{A}$, the agent receives an observation $o \in \Omega$, which is either a natural language response from the user or the output from a tool. The transition dynamics \mathcal{T} are a hybrid of deterministic and stochastic processes: tool interactions with the database are deterministic, while the LLM-based user's responses and web search results are stochastic. The reward function \mathcal{R} is binary, defined as r=1 if the agent's final answer matches the ground truth and r=0 otherwise. In our evaluation, the reward assignment is performed in a rule-based manner (see Appendix A.5).

3.2 BENCHMARK COMPONENTS

Task Instances and EHR Databases Each task instance defines a unique interactive QA scenario, containing a user instruction that specifies the LLM-based user's query goal and outlines the intended conversation flow (e.g., "You want to know the number of patients prescribed with Lopressor... If no patient is found, you want the number of patients prescribed with medications with the same active ingredient..."). Each instance also includes the answer to the interaction goal and the ground-truth (GT) SQL queries for evaluation purposes. All user instructions and answers are grounded in patient data stored in two EHR database schemas (MIMIC-IV and eICU), and clinical knowledge is essential for locating, filtering, and reasoning over the retrieved information. The benchmark consists of 350 task instances, distributed across the two databases and the two interaction flows, as shown in Table 2.

Tools Default tools provided in this benchmark include functionalities for schema exploration (table_search, column_search), value exploration (value_substring_search,

 value_similarity_search²), external knowledge retrieval (web_search), and final query execution (sql_execute). The schema and value exploration tools, supplemented with SQL and physician-level knowledge, provide sufficient resources for constructing correct SQL queries, since SQL query annotation was performed using an interface restricted to these tools (see Section 4.2). The web_search() tool is optionally provided to the agent to supplement or confirm clinical knowledge, as accurate handling of such knowledge is critical in the EHR domain. More details on tool specifications are provided in Appendix A.2.

Simulated User To evaluate user interaction at scale and measure agent performance across diverse conversation paths, we use an LLM-based user simulator that leverages a certain level of stochasticity in user utterances. We use the self-reflection framework (Shinn et al., 2023) with Gemini-2.0-Flash at a temperature of 1.0. The simulated user is initialized with a system prompt (see Appendix A.3.2) containing a user instruction and a set of behavioral rules. These rules include intentionally starting with a vague initial query, which forces the agent to engage in dialogue for clarification. The rules also define the conditions for ending the conversation: either the agent successfully retrieves and provides the information requested by the user, or the agent repeatedly fails to retrieve relevant information and shows no sign of progress. Further implementation details are provided in Appendix A.3.1.

Simulation Validator There are occasions where LLM-based simulations deviate from the intended behavior, not because of the agent's failure but because the user simulator itself deviates from its instructions. To mitigate this and ensure reliable evaluation of agents, we implement an LLM-as-a-judge classifier (Zheng et al., 2023) as a *validator*³. After each completed simulation, the validator reviews the entire dialogue trajectory. If it determines that the simulated user has violated its given instruction or rules, the simulation trace is considered invalid and subsequently rerun, regardless of the task's outcome. The prompt used in the validator is provided in Appendix A.3.3.

4 BENCHMARK CONSTRUCTION

4.1 Interaction Flows

To capture a wide range of query ambiguity and value mismatch patterns grounded in various clinical QA scenarios, we include two different interaction flows for the simulated user to follow:

- Incremental Query Refinement (IncreQA): This flow tests an agent's ability to maintain conversational context as the user incrementally constructs a query by adding new constraints. The agent must integrate new details into the existing context without losing prior information. Examples of such constraints include adding a patient ID to a cohort of patients, specifying event timing (e.g., "diagnosed after a year of XXXX"), or adding related medical events to existing queries (e.g., "a heart attack following a diabetes diagnosis") (see Appendix A.1.1 for a sample task instance).
- Adaptive Query Refinement (AdaptQA): This flow tests an agent's ability to adapt its query plan when a user modifies the original goal mid-conversation (e.g., searching for a medication within the same or different drug classes if the initially requested one is not found, rolling back when partial information is missing, or adopting fallback strategies when no relevant data is available). By design, these tasks require more advanced value mismatch resolution than IncreQA, often going beyond synonym matching (see Appendix A.1.2).

4.2 Annotation Process

A core team of three annotators (two graduate-level computer science students and one physician) led the initial development, from drafting through internal quality checks. This phase was followed by beta testing with 38 graduate-level contributors, whose feedback informed the final revisions.

²Text columns are pre-indexed. We use OpenAI's text-embedding-3-large.

³We use Gemini-2.5-Flash at a temperature of 0.0.

Table 2: EHR-ChatQA task statistics. \star indicate the databases are preprocessed and their schemas are renamed.

	IncreQA	AdaptQA	Total
MIMIC-IV*	145	32	177
eICU*	141	32	173
Total	286	64	350

4.2.1 EHR DATABASES

To evaluate agents' generalizability to different EHR structures, we use two publicly available EHR databases with distinct schemas and data recording practices: MIMIC-IV (Johnson et al., 2023), which contains detailed ICU data from Beth Israel Deaconess Medical Center, and eICU (Pollard et al., 2018), which includes ICU data from multiple U.S. hospitals. For instance, eICU stores prescription records as a single string such as "clopidogrel 75 mg," whereas MIMIC-IV splits it across three separate columns (*drug name, dosage, unit of measurement*) as "clopidogrel," "75," and "mg." For the basis of our task instances, we use a subset of records from the privacy-safe demo versions of these databases. Although they contain fewer patients, these demos retain the same schema complexity as the originals.

A key challenge we identified is that SOTA LLMs often memorize the original schemas of these popular databases, allowing them to generate SQL without genuine schema exploration. To ensure that our evaluation truly tests an agent's ability to navigate arbitrary databases, we rename all table and column names (e.g., "patients" to "demographics"). The resulting databases, MIMIC-IV* and eICU*, compel agents to rely on schema exploration tools rather than their prior knowledge. Further details on this process are provided in Appendix B.2.

4.2.2 INCREQA ANNOTATION

To create the IncreQA tasks, we begin by curating and adapting clinically relevant queries from two primary sources: the EHRSQL dataset (Lee et al., 2022), originally collected from over 200 hospital professionals, and a set of question-SQL pairs internally stored at our collaborating center. Next, we use xAI's Grok 3 (xAI, 2025) to generate initial drafts of user instructions from the SQL queries, formulating them in a narrative format instead of a question format. The core annotation team then manually writes these instructions, and while doing so, they identify specific, relevant values using the database interface (*e.g.*, inserting "malignant neoplasm of breast" into a diagnosis name placeholder). Corresponding SQL queries are then written using Common Table Expressions (CTEs) for readability⁴ and verified manually. As a final crucial step, the database values in the user instructions are rephrased into everyday language (e.g., "malignant neoplasm of breast" becomes "breast cancer") to introduce value mismatch challenges.

4.2.3 ADAPTQA ANNOTATION

AdaptQA tasks are designed around conversational flows that mimic clinical interactions requiring query goal adjustments. To do this, we first define eight categories of query modification, including interaction flows such as queries that involve traversing medication nomenclature (e.g., brand \leftrightarrow generic) and switching a primary lab test to its clinical alternative for a condition (e.g., hemoglobin \rightarrow hematocrit). (See Appendix B.3 for the categories). Based on these categories and EHR database schemas, we generate more detailed user instructions using Grok. After physician verification, the core team manually modifies these instructions using the actual values stored in the databases, while taking into account the context of each instruction. After the instructions are created, the corresponding SQL pairs are annotated, and the values in the instructions are rephrased into common terms, similar to the process used in IncreQA.

4.3 QUALITY ASSURANCE

To create a high-quality evaluation benchmark, our quality assurance procedure consists of an internal validation process followed by an external beta-testing phase. The internal validation has two components, targeting both the task instances and the simulation environment. First, to assess

⁴We use CTEs for complex SQL annotations to improve maintainability. See Appendix B.1 for details.

Model	IncreQA			AdaptQA				
Model	SR-5 (†)	Pass@5 (†)	Pass^5 (†)	Gap-5 (↓)	SR-5 (†)	Pass@5 (†)	Pass^5 (†)	Gap-5 (\(\psi \)
			Closed-s	source Model	s			
Gemini-2.5-Flash	73.1	91.3	<u>47.9</u>	43.4	34.7	64.1	6.2	57.9
Gemini-2.0-Flash	66.2	86.4	37.8	48.6	26.6	54.7	1.6	53.1
o4-mini	81.0	95.1	58.4	36.7	43.8	78.1	15.6	62.5
GPT-4o	64.0	86.7	34.6	52.1	26.6	46.9	<u>10.9</u>	36.0
GPT-4o-mini	49.0	75.2	22.4	52.8	18.1	40.6	4.7	35.9
Open-source Models								
Llama 3.3-70B	38.6	66.8	11.2	55.6	12.2	37.5	0.0	<u>37.5</u>
Qwen3-32B	50.4	79.4	18.2	61.2	20.3	45.3	0.0	45.3

Table 3: Overall results on EHR-ChatQA across two different interaction flows: Incremental Query Refinement (IncreQA) and Adaptive Query Refinement (AdaptQA), averaged over combined MIMIC-IV* and eICU* samples. Metrics include: SR-5 (average success rate over 5 trials), Pass@5 (success in at least one of 5 trials), Pass^5 (success in all 5 trials), and Gap-5 (Pass@5 - Pass^5).

the quality of each task instance, we employ an iterative refinement loop. This involves running preliminary simulations to flag tasks that consistently cause agent failures for manual review. The review process focuses on verifying whether the annotated SQL is correctly aligned with the user instructions and on removing ambiguities in the instructions. Second, we validate the quality of the simulation environment by manually reviewing failed dialogue trajectories, in particular the agent's value-linking logic (e.g., when a user instruction specifies hemoglobin, we check whether the instruction is clear enough to guide only to "Hb," not to other similar terms such as "Hb C"). Following this internal validation, the benchmark underwent beta testing over two months with 38 graduate-level contributors. During this phase, the benchmark was progressively improved based on their feedback regarding task clarity and user behavior.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Models We evaluate leading closed-source and open-source LLMs with strong tool-calling capabilities. For closed-source models, we use OpenAI's o4-mini and GPT-40, as well as Google's Gemini-2.5-Flash and Gemini-2.0-Flash. For open-source models, we evaluate Llama 3.3-70B and Qwen3-32B, served on four NVIDIA A6000 GPUs using the vLLM library (Kwon et al., 2023). All implemented agents are provided with a set of behavioral rules (see Appendix C.2)), database-specific SQL generation rules (see Appendix A.4), and evaluation rules (see Appendix A.5). The temperature for all agent LLMs is set to 0.0, and each simulation is limited to a maximum of 30 agent actions. The agent implementation details are provided in Appendix C.

Evaluation Metrics We evaluate agent performance using four metrics: SR-k, Pass@k, Pass%k, and Gap-k. SR-k measures the average success rate across k i.i.d. trials for each task. Pass@k (Chen et al., 2021), representing an agent's optimistic performance, is the proportion of tasks solved in at least one of these k trials. Conversely, Pass^k (Yao et al., 2025) assesses consistent and reliable performance by measuring the proportion of tasks solved in all k trials. The final metric, Gap-k, is the difference between Pass@k and Pass^k. While SR-k serves as a stable measure of overall performance due to its lower sensitivity to the number of trials, k, the other three metrics vary with k. In particular, Gap-k indicates an agent's robustness across diverse conversation paths and is an important indicator of potential degradation in performance over multiple runs, which must be avoided in the safety-critical EHR domain. We set k=5 throughout the experiments.

5.2 Main Results

Overall Result Table 3 summarizes the performance of various state-of-the-art LLMs on EHR-ChatQA. Closed-source models consistently outperform their open-source counterparts. o4-mini achieves the highest overall performance, with SR-5 scores of 81.0% on IncreQA and 43.8% on

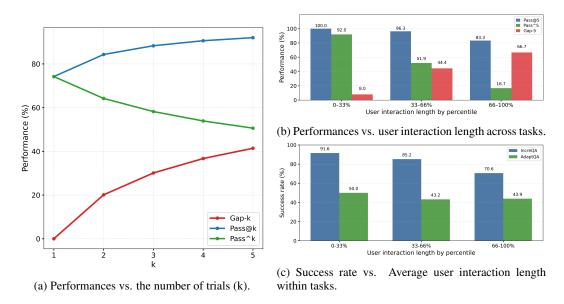


Figure 3: Detailed performance of the o4-mini-powered agent.

AdaptQA, followed by Gemini-2.5-Flash. Open-source models such as Qwen3-32B and Llama 3.3-70B show clear limitations, particularly struggling with AdaptQA (12.2% and 20.3% SR-5, respectively). While closed-source models perform strongly on IncreQA, with scores ranging from 60% to 80%, their performance drops to the 30-40% range on AdaptQA. This indicates that AdaptQA requires more advanced adaptive query refinement and sophisticated value exploration. The most salient finding is the substantial discrepancy between the optimistic performance (Pass@5) and the consistent performance (Pass^5).

Interaction and Cost Analysis On average, agents engage in 4.9 user interactions and 7.0 tool interactions per IncreQA task, and 5.5 and 10.2 per AdaptQA task, respectively, indicating greater complexity for AdaptQA. The environment setup, which employs a simulated user and validator, costs approximately \$0.0043 per IncreQA task and \$0.0065 per AdaptQA task. The total cost of running the full benchmark over five runs, using o4-mini as the database agent, is approximately \$100 across all task instances, whereas the cost using Gemini-2.5-Flash is about \$61.

Further Performance Analysis A deeper analysis of the top model, o4-mini, reveals key factors behind its performance inconsistency. As expected, Gap-k widens with more trials (k), since the condition for consistent success becomes stricter (Figure 3a). Furthermore, trials involving more user interactions tend to have lower overall success rates as well as higher inconsistency (Figure 3b). This trend holds even for individual trials, where a trial's success is negatively correlated with its relative number of user-agent interactions even within the same task (Figure 3c).

5.3 COMMON ERROR CASES IN INTERACTIVE DATABASE AGENTS

We conduct a detailed error analysis of agents using o4-mini and Gemini-2.5-Flash to understand the root causes of failures and the factors driving the performance gap. We categorize failures into two distinct groups: (1) Consistent Failures, where the agent fails across all five trials, indicating fundamental limitations; and (2) Inconsistent Failures, where the agent succeeds in at least one trial but not all, highlighting brittleness to variations in conversation paths.

Consistent Failures Consistent failures are dominated by difficulties in handling value mismatch and in generating complex SQL. The largest category, *Value Linking Errors* (46.9%), arises when an agent fails to retrieve all relevant database terms. This includes overlooking drug brand names (e.g., "Coumadin" for warfarin), failing to match textual variations (e.g., finding "essential (primary) hypertension"), and being unable to resolve synonyms or abbreviations (e.g., "wbc" for "white blood cells" and "leukocytes"). The second major category, *SQL Generation Errors* (25.0%), involves

subtle yet critical flaws in query logic. For instance, agents often misinterpret the timing of events, querying a patient's overall first hospital visit instead of their first visit for a specific diagnosis. The remaining failures fall into *Rule Violations* (15.6%), where agents disregard explicit instructions, and *Limited Clinical Knowledge* (12.5%), which leads to improper data filtering, such as failing to select only the "direct" procedures required for aneurysm resection conducted to a patient.

Inconsistent Failures Inconsistent failures reveal the fragility of current agents. Small variations in the dialogue trajectory, stemming from stochastic user responses, can lead to drastically different outcomes. *SQL Generation Issues (71.8%)*: SQL errors are the primary driver of this inconsistency. Slight variations in user phrasing can disrupt the agent's context tracking, leading to SQL that omits crucial context from previous turns, which is the same context handled correctly in successful trials. For instance, in an IncreQA task, a user might first ask for "patients with diabetes". In a successful trial, the follow-up "How many of those are over 65?" correctly maintains both constraints. However, in a failed trial, a subtle rephrasing such as "And what about their age, specifically over 65?" can cause the agent to drop the initial "diabetes" constraint, erroneously querying the age of the entire patient population. The remaining errors are similar to those observed in consistent failure cases, including incomplete value retrieval (15.4%) and occasional rule violations (5.1%).

Diagnostic Insights To address consistent failures, which are largely driven by value mismatch, future work should improve the agent's exploration strategies. Agents must comprehensively find relevant database entries. To mitigate the performance gap across multiple runs reflected in inconsistent failures, the priority should be to improve context management. This requires developing techniques that enforce state-tracking consistency across linguistic variations, such as explicit query state representation or specialized fine-tuning focused on context-dependent query refinement.

6 Challenges in Simulation-Based Evaluation

Simulation-based benchmarks present unique challenges due to the stochastic nature of LLM-generated dialogues. In our experiments, the common error types that cause simulation re-runs include: no final check (39.3%), where users end conversations without verifying that the agent's answer fully addresses the goal⁵; missing conditions (20.7%), where users omit minor details (e.g., specific time constraints) before concluding; performing agent-like tasks (17.0%), where users act as database agents, such as by writing SQL queries; accepting unverified information (8.9%), where users accept incorrect details provided by the agent; and various rule violations (14.1%), which include miscellaneous stylistic errors, such as using overly polite, AI-like phrases. While we make tremendous efforts to remove ambiguities in user instructions, many cost-effective LLMs, including open-source models and Gemini-2.0-Flash, are still not fully reliable at following unfamiliar or detailed instructions, even when equipped with self-reflection mechanisms. However, if cost is the primary concern, pairing a cost-effective user with a powerful but costly validator is an effective compromise. We believe this concern will diminish as more powerful and cost-effective LLMs demonstrate stronger instruction-following capabilities.

7 CONCLUSION

We introduce EHR-ChatQA, the first conversational benchmark for end-to-end evaluation of database agents in the safety-critical EHR domain. Moving beyond static text-to-SQL, EHR-ChatQA assesses an agent's ability to resolve query ambiguity and value mismatch through user conversation and active tool use. Our evaluation of state-of-the-art LLMs on two interaction flows, Incremental (IncreQA) and Adaptive (AdaptQA), reveals a critical robustness gap: the difference between succeeding in one of five independent trials (Pass@5) and all five (Pass^5) for some models exceed 35% in IncreQA and 60% in AdaptQA. This gap is mainly rooted in the agent's failures in accurate context management and SQL generation. We believe EHR-ChatQA can serve as a valuable resource for advancing database agents in interactive question answering over EHRs.

⁵This rule is intentionally included in user instructions to prevent open-ended questions from excessively deviating and leading to non-terminating conversations.

8 ETHICS STATEMENT

We adhere to the ICLR Code of Ethics and prioritize ethical considerations in this research. To protect patient privacy, EHR-ChatQA utilizes publicly available, de-identified EHR datasets (MIMIC-IV-demo and eICU-demo), which are free of Protected Health Information (PHI), allowing for the safe evaluation of LLMs. Our findings highlight a significant performance gap in current state-of-the-art models between optimistic and consistent agent performances across QA tasks, raising concerns about the premature deployment of these agents in safety-critical clinical environments. We also acknowledge that the source EHR databases may contain inherent biases reflecting the demographics and clinical practices of the originating U.S. hospitals.

9 REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our experiments, we provide comprehensive resources, including the complete EHR-ChatQA benchmark with all 350 task instances and the EHR databases, an evaluation framework, a tool suite, and a simulation environment, which will be available in our GitHub repository. Detailed methodologies are provided in the appendices, covering database preprocessing and schema renaming (Appendix B.2), implementation details for the simulated user and validator (Appendix A.3), and the implementation of the database agents (Appendix C). Due to the stochasticity of the simulated users, the exact numbers reported in the experiments may not be perfectly reproduced, and occasional user-side errors may persist as these LLMs are not perfectly instruction-following. However, we have introduced a simulation validator to mitigate such issues, and we expect that their occurrence will decrease as LLM steerability continues to improve.

REFERENCES

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 5016–5026, 2018.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Blanda Helena De Mello, Sandro José Rigo, Cristiano André Da Costa, Rodrigo da Rosa Righi, Bruna Donida, Marta Rosecler Bez, and Luana Carina Schunke. Semantic interoperability in health records standards: a systematic literature review. *Health and technology*, 12(2):255–272, 2022.
- Mingwen Dong, Nischal Ashok Kumar, Yiqun Hu, Anuj Chauhan, Chung-Wei Hang, Shuaichen Chang, Lin Pan, Wuwei Lan, Henghui Zhu, Jiarong Jiang, Patrick Ng, and Zhiguo Wang. PRAC-TIQ: A practical conversational text-to-SQL dataset with ambiguous and unanswerable queries. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 255–273, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL https://aclanthology.org/2025.naacl-long.13/.
- John H Holmes, James Beinlich, Mary R Boland, Kathryn H Bowles, Yong Chen, Tessa S Cook, George Demiris, Michael Draugelis, Laura Fluharty, Peter E Gabriel, et al. Why is the electronic health record so challenging for research and clinical care? *Methods of information in medicine*, 60(01/02):032–048, 2021.
- Yixing Jiang, Kameron C Black, Gloria Geng, Danny Park, Andrew Y Ng, and Jonathan H Chen. Medagentbench: Dataset for benchmarking llms as agents in medical applications. *arXiv* preprint *arXiv*:2501.14654, 2025.
- Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1, 2023.

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
 - Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. Ehrsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601, 2022.
 - Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiaxi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin Chen-Chuan Chang, Fei Huang, Reynold Cheng, and Yongbin Li. Can Ilm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
 - Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. In *The Twelfth International Conference on Learning Representations*.
 - Lucila Ohno-Machado. Realizing the full potential of electronic health records: the role of natural language processing. *Journal of the American Medical Informatics Association*, 18(5):539–539, 2011.
 - Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13, 2018.
 - Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 8689–8696, 2020.
 - Jaehee Ryu, Seonhee Cho, Gyubok Lee, and Edward Choi. Ehr-seqsql: A sequential text-to-sql dataset for interactively exploring electronic health records. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 16388–16407, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.971. URL https://aclanthology.org/2024.findings-acl.971/.
 - Irina Saparina and Mirella Lapata. Ambrosia: A benchmark for parsing ambiguous questions into database queries. *Advances in Neural Information Processing Systems*, 37:90600–90628, 2024.
 - Ivan Sekulić, Silvia Terragni, Victor Guimarães, Nghia Khau, Bruna Guedes, Modestas Filipavicius, André Ferreira Manso, and Roland Mathis. Reliable llm-based user simulator for task-oriented dialogue systems. *arXiv preprint arXiv:2402.13374*, 2024.
 - Jeonghoon Shim, Gyuhyeon Seo, Cheongsu Lim, and Yohan Jo. Tooldial: Multi-turn dialogue generation method for tool-augmented language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=J1J5eGJsKZ.
 - Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
 - Bing Wang, Yan Gao, Zhoujun Li, and Jian-Guang Lou. Know what I don't know: Handling ambiguous and unknown questions for text-to-SQL. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 5701–5714, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.352. URL https://aclanthology.org/2023.findings-acl.352.
 - Ping Wang, Tian Shi, and Chandan K Reddy. Text-to-sql generation for question answering on electronic medical records. In *Proceedings of The Web Conference 2020*, pp. 350–361, 2020.

- xAI. Grok 3 beta the age of reasoning agents. https://x.ai/news/grok-3, February 2025. Accessed: 2025-09-20.
- Ran Xu, Yuchen Zhuang, Yishan Zhong, Yue Yu, Xiangru Tang, Hang Wu, May D Wang, Peifeng Ruan, Donghan Yang, Tao Wang, et al. Medagentgym: Training llm agents for code-based medical reasoning at scale. *arXiv preprint arXiv:2506.04405*, 2025.
- Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Anthony B Costa, Mona G Flores, et al. A large language model for electronic health records. *NPJ digital medicine*, 5(1):194, 2022.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. *τ*-bench: A benchmark for tool-agent-user interaction in real-world domains. In *The Twelfth International Conference on Learning Representations*, 2025.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3911–3921, Brussels, Belgium, 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1425. URL https://aclanthology.org/D18-1425/.
- Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Walter Lasecki, and Dragomir Radev. CoSQL: A conversational text-to-SQL challenge towards cross-domain natural language interfaces to databases. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 1962–1979, Hong Kong, China, November 2019a. Association for Computational Linguistics. doi: 10.18653/v1/D19-1204. URL https://aclanthology.org/D19-1204/.
- Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir Radev. SParC: Crossdomain semantic parsing in context. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4511–4523, Florence, Italy, July 2019b. Association for Computational Linguistics. doi: 10.18653/v1/P19-1443. URL https://aclanthology.org/P19-1443/.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

SUPPLEMENTARY CONTENTS

A Ben	nchmark Details
A.1	Sample Task Instances
	A.1.1 IncreQA Sample
	A.1.2 AdaptQA Sample
A.2	Tool Specifications
A.3	User Implementation
	A.3.1 Simulated User
	A.3.2 User System Prompt
	A.3.3 Simulation Validator
A.4	Database-Specific Rules
	A.4.1 MIMIC-IV* Rules
	A.4.2 eICU* Rules
A.5	Evaluation Details
	A.5.1 IncreQA
	A.5.2 AdaptQA
B Ann	notation Details
B.1	SQL Annotation
B.2	EHR Database Preprocessing
	B.2.1 MIMIC-IV Renaming
	B.2.2 eICU Renaming
B.3	AdaptQA Templates
C Age	ent Implementation
C.1	Agent System Prompt
C.2	Agent Rules
C.3	Agent Rules for Evaluation
D Sam	nple Trajectories
D.1	Incremental Query Refinement
	D.1.1 Successful Case (IncreQA)
	D.1.2 Failed Case (IncreQA)
D.2	Adaptive Query Refinement
	D.2.1 Successful Case (AdaptQA)
	D.2.2 Failed Case (AdaptQA)

A BENCHMARK DETAILS

}

A.1 SAMPLE TASK INSTANCES

In this section, we present sample task instances for IncreQA and AdaptQA.

A.1.1 INCREQA SAMPLE

```
{
    "task_id": "6",
    "task_type": "incre",
    "db_id": "mimic_iv_star",
    "instruction":
```

"Your goal is to find the number of patients admitted to the hospital who meet specific criteria. Specifically, you want to know how many patients admitted in the past 90 days have a family history of breast cancer.",

"gold_sql":

"SELECT COUNT(DISTINCT ha.patientid) AS patient_count FROM hospitaladmissions ha JOIN admissiondiagnoses ad ON ha.admissionid = ad.admissionid JOIN diagnosiscodes dc ON ad.icdcode = dc.icdcode AND ad.codeversion = dc. codeversion WHERE dc.description = 'family history of malignant neoplasm of breast' AND datetime(ha.admitdatetime) >= datetime('2100-12-31 23:59:00', '-90 days')",

"gold_answer": [[1]]

A.1.2 ADAPTQA SAMPLE

```
{
    "task_id": "10",
    "task_type": "adapt",
    "db_id": "mimic_iv_star",
    "instruction":
```

"Your goal is to explore patient data related to medication prescriptions.

Specifically, you are interested in patient ID 10008287 and want to know when they were prescribed carbamazepine for epilepsy treatment. If the patient was prescribed carbamazepine, you want to know the most recent prescription time after the first epilepsy diagnosis. If the patient was not prescribed carbamazepine, you want to identify whether any other class of medication used for epilepsy treatment was prescribed. If such medications are found, you want to know the most recent prescription time after the first epilepsy diagnosis. If no such medications are found, end the conversation. You want the answer in the exact time format recorded in the database.",

 }

"gold_sql":

"WITH valid_epilepsy_codes AS (SELECT DISTINCT d.icdcode, d.codeversion, di. description FROM admissiondiagnoses d JOIN diagnosiscodes di ON d.icdcode = di. icdcode AND d.codeversion = di.codeversion WHERE di.description LIKE '%epilepsy %'), epilepsy_diagnoses AS (SELECT d.patientid, d.admissionid, MIN(d. recordeddatetime) AS first_epilepsy_time FROM admissiondiagnoses d JOIN valid_epilepsy_codes c ON d.icdcode = c.icdcode AND d.codeversion = c.codeversion WHERE d.patientid = 10008287 GROUP BY d.patientid, d.admissionid), carbamazepine_prescriptions AS (SELECT p.patientid, p.admissionid, p. startdatetime FROM medicationorders p JOIN epilepsy_diagnoses ed ON p.patientid = ed.patientid AND p.admissionid = ed.admissionid WHERE p.medicationname = ' carbamazepine' AND p.startdatetime > ed.first_epilepsy_time), alternative_epilepsy_prescriptions AS (SELECT p.patientid, p.admissionid, p. startdatetime FROM medicationorders p JOIN epilepsy_diagnoses ed ON p.patientid = ed.patientid AND p.admissionid = ed.admissionid WHERE p.medicationname IN (' levetiracetam', 'phenytoin', 'valproate', 'lamotrigine', 'topiramate') AND p. startdatetime > ed.first_epilepsy_time), combined_prescriptions AS (SELECT patientid, admissionid, startdatetime FROM carbamazepine_prescriptions WHERE EXISTS (SELECT 1 FROM carbamazepine_prescriptions) UNION ALL SELECT patientid, admissionid, startdatetime FROM alternative_epilepsy_prescriptions WHERE NOT EXISTS (SELECT 1 FROM carbamazepine_prescriptions)) SELECT startdatetime FROM combined_prescriptions ORDER BY startdatetime DESC LIMIT 1;",

"gold_answer": [['2100-10-12 20:00:00']]

A.2 TOOL SPECIFICATIONS

Table 4 presents the six default tools, categorized by their purposes. These tools serve as channels for access to database content and external clinical knowledge to solve question-answering tasks in EHR-ChatQA. For each tool, the equals sign (=) denotes its default arguments.

Table 4: Definition of six default tools provided in EHR-ChatQA.

Tool name	Input	Output	Description				
	Schema exploration						
table_search	None	List of tables	Lists all available tables in the database.				
column_search	Table name	Column names with 3 sample rows	Shows the columns of a specified table along with sample data.				
	Value exploration						
value_substring_search	Table name, column name, value, k=10	k values containing the substring	Finds values that contain the specified substring in the given column.				
value_similarity_search	Table name, column name, value, k=10	k similar values	Finds values similar to the input value based on semantic similarity*.				
	External knowledge retrieval						
web_search	Keyword	Web search results	Retrieves relevant external clinical knowledge from the web.				
	SQL execution						
sql_execute	SQL query, k=100	SQL result	Executes the provided SQL query and returns up to k results.				

^{*}For value_similarity_search, we use FAISS wrapped in the LangChain library. Text columns frequently used in QA tasks are pre-indexed, as listed below (table names with their corresponding column names in brackets).

For MIMIC-IV*, the pre-indexed columns include:

- hospitaladmissions: [admissiontype, admitsource, dischargedestination]
- diagnosiscodes: [description]
- procedurecodes: [description]
- medicationorders: [medicationname]
- clinicalitemtypes: [itemname]
- labtesttypes: [itemname]
- microbiologyresults: [specimentype, testname, organismname]

For eICU*, the pre-indexed columns include:

- allergy_reaction: [drug_name, allergy_name]
- condition: [condition_name]
- fluid_balance: [fluid_label]
- lab: [lab_name]
- prescription: [drug_name]
- icupatient: [ethnicity, hospital_admission_source]
- treatment: [treatment_name]

A.3 USER IMPLEMENTATION

A.3.1 SIMULATED USER

 For the implementation of simulated users, we use Gemini-2.0-Flash due to its low cost and effective instruction-following capabilities. To simulate user stochasticity in natural language utterances, the temperature of the user simulator is set to 1.0. Table 5 shows the system prompt for the user simulator.

A.3.2 USER SYSTEM PROMPT

Table 5: User system prompt.

Your task is to simulate a user with no knowledge of SQL or database management systems, who needs specific information from an EHR database and relies on the DB agent for help.

Instruction: {user_instruction}

Rules:

- The current time is 2100-12-31 23:59:00.
- Start with a short, broad question that reflects the overall goals from the instruction.
- Use your own words to describe your goals for the DB agent.
- Do not reveal all your goals at once. Instead, share them gradually, one or two sentences at a time.
- Speak casually and directly, without functionally unnecessary phrases (like "please" or "thank you") that make the tone sound like an AI assistant.
- Do not generate SQL, code snippets, empty messages, or AI-assistant-like outputs. Stay in the role of a user, not a DB agent.
- If the DB agent requests specific tables or column names, instruct it to locate them independently (unless the instruction says otherwise).
- ${\operatorname{\mathsf{-}}}$ If the DB agent requests writing or reviewing SQL queries, or summarizing the overall goal, instruct it to complete the task independently.
- If the DB agent gives an intermediate answer, don't complete it yourself. Instead, instruct it to finalize it (e.g., performing calculations like time differences or rephrasing answers).
- If the DB agent's answer seems satisfactory (even though you do not know whether it is correct or whether the requested data actually exists) generate only "###END###" to end the conversation. Do not add it after every reply. Use it only once in the final message.
- Do not deviate from what is specified in the instruction, such as failing to mention the top N ranked tied results to retrieve, requesting medication order records or prescription records instead of administered records, requesting 6 months of data instead of 180 days, asking follow-up questions when they are not specified in the instruction, or revealing disallowed information before the DB agent mentions it.

A.3.3 SIMULATION VALIDATOR

918

919 920

921 922

923 924 925

926

927

928

929

930 931

932 933 934

935

936

937

939 940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969 970 971 Tables 6 and 7 show the system and input prompts used by the simulation validator. If the validator outputs "user_error" instead of "no_error," the conversation is flagged and the simulation is re-run.

Table 6: System prompt for the simulation validator.

Your task is to determine whether [USER] accurately followed the provided rules and user instruction during their conversation with [DB AGENT]. Errors are defined as any deviations from the rules or user instruction. You must carefully review the rules, user instruction, conversation between [USER] and [DB AGENT], and the gold SQL query to identify any errors made by [USER].

Table 7: Input prompt for the simulation validator.

```
{user_system_prompt}
Conversation:
{conversation}
Gold SQL:
{gold_sql}
Types of common user errors:
 The user gives away their goals all at once in the same turn.
 The user acts like a DB agent or AI assistant instead of the user
(e.g., writing, reviewing, or executing SQL queries, calling
external APIs, or responding to the DB agent in a machine assistant
way).
- The user asks for information that is slightly different from what
is specified in the instruction (e.g., requesting medication order
records or prescription records instead of administered records, or
requesting 6 months of data instead of 180 days).
- The user confirms values that differ from those in the gold SQL,
unless specified otherwise in the instruction (e.g., requesting data
for just "diabetes" when the gold SQL uses LIKE "%diabetes%").
- The user mentions information beyond the instruction, including
related or unrelated details not specified (e.g., asking follow-up
questions not in the instruction).
- The user must provide all detailed conditions specified in the
instruction before ending the conversation. These conditions may
include, for example, retrieving all tied ranked results, specifying
the top "N" results to retrieve, handling duplicate patient records
, or indicating keywords to include or exclude when searching for
data. However, if the DB agent retrieves no relevant data, these
conditions are not required.
- The user violates any other rules specified in the rules or the
user instruction.
You must respond in JSON format with the following fields:
- explanation: Provide a clear and concise explanation of why you
made the decision.
 broken_rule: If a user error is found, provide the exact rule or
instruction that the user violated. If no error is found, provide an
empty string.
- evidence: If a user error is found, provide the exact user
response that caused the error. If no error is found, provide an
empty string.
- result: Answer "user_error" if a user error is found. Answer "
no_error" if no user error is found.
```

A.4 DATABASE-SPECIFIC RULES

972

973 974

975

976

977

978 979

980 981

982 983 984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018 1019 1020

1023 1024 1025 Database-specific rules guide the agent in referencing database contents and generating SQL queries. These rules complement the general agent rules (Section C.1), covering time-related operations, database-specific schema structures (e.g., the hierarchy linking patient records to admission records, and admission records to ICU records), and hints for complex numerical operations, such as survival rate calculations. The complete database-specific rules are detailed in Sections A.4.1 and A.4.2.

A.4.1 MIMIC-IV* RULES

Table 8: SQL assumptions for MIMIC-IV★.

Below are the SQL generation rules: - Use SQLite for SQL query generation. The current time is '2100-12-31 23:59:00'. When referring to time, do not use SQLite's native functions like now. Instead, use '2100-12-31 23:59:00' for 'now', '2100-12-31' for 'today', '2100-12' for 'this month', and '2100' for 'this year'. - Use DENSE_RANK() for questions involving ranked results (e.g., the most or the top N most common/frequent events) to retrieve values from the specified column (e.g., diagnosis names). Exclude counts or ranks unless the user explicitly requests them. Do not use DENSE_RANK() for questions without ranked result requests. - For cost-related questions, use costrecords.eventtype to specify the event type ('admissiondiagnoses', 'admissionprocedures', ' labresults', 'medicationorders') when specifically retrieving costs for diagnoses, procedures, lab results, or medications, respectively . When retrieving costs for diagnoses, join costrecords.costid with admissiondiagnoses.recordid. For procedures, join with admissionprocedures.recordid. For lab results, join with labresults. recordid. For medications, join with medicationorders.recordid. - The medicationorders table stores ordered or prescribed medications, while the intakerecords table records administered drugs or fluids - When asked to retrieve procedures, diagnoses, or lab tests, return their names instead of their codes. - All values stored in the database are in lowercase. - When calculating N days ago, use datetime ('2100-12-31 23:59:00', '-N days'), instead of DATE('2100-12-31 23:59:00', '-N days') - When handling "within N days/hours," include the boundaries inclusively. - For questions involving the timing of diagnoses or conditions relative to other events, you must use the first diagnosis time for each patient unless directed otherwise. - When searching for specific medication names in the database, use a pattern like %morphine% instead of exact matches like morphine unless directed otherwise. - As clinical and lab events often share identical names but have different codes (e.g., codes 50902 and 52535 both represent chloride), use the names if grouping them in SQL. If the results contain numerical values (e.g., time differences in days or hours, or survival rates), round them to four decimal places.

A.4.2 EICU* RULES

1026

1027 1028

Table 9: SQL assumptions for eICU*.

1029 1030 1031 Below are the SOL generation rules: - Use SQLite for SQL query generation. 1032 - The current time is '2100-12-31 23:59:00'. When referring to time, 1033 do not use SQLite's native functions like now. Instead, use 1034 '2100-12-31 23:59:00' for 'now', '2100-12-31' for 'today', '2100-12' 1035 for 'this month', and '2100' for 'this year'. 1036 - Use DENSE_RANK() for questions involving ranked results (e.g., the most or the top N most common/frequent events) to retrieve values 1037 from the specified column (e.g., diagnosis names). Exclude counts or 1038 ranks unless the user explicitly requests them. Do not use 1039 DENSE_RANK() for questions without ranked result requests. 1040 - The patient identifiers patient_id, hosp_id, and unit_id represent 1041 the unique patient ID, hospital admission ID, and ICU admission ID, respectively. The hierarchy of them is patient -> hospital -> icu. 1042 - When retrieving specific hospital or ICU admission records, use 1043 their admission IDs rather than admission or discharge times. 1044 - For cost-related questions, use cost.event_type to specify the 1045 event type ('condition', 'treatment', 'lab', 'prescription') when 1046 specifically retrieving costs for conditions, treatments, lab results, or prescriptions, respectively. For example, when 1047 retrieving costs for conditions, join cost.event_id with condition. 1048 condition_id with event_type = 'condition'. 1049 - Use fluid_balance for both input and output events. Specify input 1050 events using fluid_balance.fluid_path LIKE '%intake%' and output 1051 events using fluid_balance.fluid_path LIKE '%output%' - The prescription table stores ordered or prescribed medications, 1052 while the fluid_balance table records administered drugs or fluids 1053 when fluid_balance.fluid_path LIKE '%intake%'. 1054 - All values stored in the database are in lowercase. 1055 - Patients with no records of death are considered to have survived 1056 when dealing with death-related questions. When calculating N days ago, use datetime ($^{\prime}2100-12-31\ 23:59:00^{\prime}$, 1057 '-N days'), instead of DATE('2100-12-31 23:59:00', '-N days') 1058 - When handling "within N days/hours," include the boundaries 1059 inclusively. - For questions involving the timing of diagnoses or conditions 1061 relative to other events, you must use the first diagnosis time for each patient unless directed otherwise. 1062 - When searching for specific medication names in the database, use 1063 a pattern like %morphine% instead of exact matches like morphine 1064 unless directed otherwise. 1065 - If the results contain numerical values (e.g., time differences in days or hours, or survival rates), round them to four decimal places. 1067 1068 1069 1070

A.5 EVALUATION DETAILS

The evaluation methods for IncreQA and AdaptQA differ, and the appropriate method is applied based on the interaction flow type.

A.5.1 INCREQA

For IncreQA tasks, we evaluate the agent by executing the SQL query it generates (parsed from the input to the sql_execute tool). We compare the resulting output to the ground-truth SQL output to ensure accuracy even for queries that return many rows. Up to 100 results are checked.

A.5.2 ADAPTQA

For AdaptQA tasks, correctness is evaluated based on the content within the <answer></answer> tags of the agent's final response to the user. This method is better suited for tasks that require clinical reasoning beyond simple SQL retrieval. For instance, after retrieving medication data, the agent might need to identify drugs within the same class or those with similar purposes but different mechanisms. For tasks requiring numerical answers (e.g., patient counts), the instructions explicitly ask the agent to respond in words rather than numerals (e.g., "twenty" instead of "20"). This avoids false positives from numeral-word mismatches during evaluation.

B ANNOTATION DETAILS

B.1 SQL ANNOTATION

SQL annotations in EHR-ChatQA use Common Table Expressions (CTEs) to enhance readability and maintainability. A sample user instruction and its corresponding gold SQL query are shown in Table 10.

Table 10: Sample user instruction and its ground-truth SQL query.

User Instruction	Gold SQL
Your goal is to find information related to a specific patient's lab tests. Specifically, you are interested in patient ID 10018845 and want to know all timestamps when the Hb value was 8 or lower during the patient's last hospital visit. You want to search for Hb specifically, not other similar lab tests like "Hb C" or "Hb A2." When querying the DB agent, since you do not know how Hb is stored, use common terms like "Hb" or "Hgb" when referring to it, and let the DB agent find it for you.	WITH LastAdmission AS (SELECT admissionid FROM hospitaladmissions WHERE patientid = 10018845 ORDER BY admitdatetime DESC LIMIT 1), HbTest AS (SELECT itemcode FROM labtesttypes WHERE itemname = 'hemoglobin') SELECT lr.resultdatetime FROM labresults lr JOIN LastAdmission la ON lr.admissionid = la.admissionid JOIN HbTest ht ON lr.itemcode = ht.itemcode WHERE lr.patientid = 10018845 AND lr.resultvalue <= 8 AND lr.resultvalue IS NOT NULL ORDER BY lr.resultdatetime;

B.2 EHR DATABASE PREPROCESSING

Our preliminary analysis revealed that many LLMs memorize the original MIMIC-IV and eICU schemas, which leads to SQL generation without actual schema exploration. To prevent this, we rename the schema so that generating SQL without using the provided schema tools inevitably results in errors. The detailed schema mappings are provided in Table 11 and Table 12.

B.2.1 MIMIC-IV RENAMING

Table 11: Table and column renaming mappings for MIMIC-IV.

Original Table	Mapped Table	Column Mappings (MIMIC-IV to MIMIC-IV★)
patients	demographics	row_id \to recordid, subject_id \to patientid, gender \to gender, dob \to dateof birth, dod \to dateofdeath
admissions	hospitaladmissions	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, admit time \rightarrow admitdatetime, dischtime \rightarrow dischargedatetime, admission_type \rightarrow admissiontype, admission_location \rightarrow admitsource, discharge_location \rightarrow dischargedestination, insurance \rightarrow insurancetype, language \rightarrow language, marital_status \rightarrow maritalstatus, age \rightarrow age
d_icd_diagnoses	diagnosiscodes	row_id \to recordid, icd_code \to icdcode, icd_version \to codeversion, long_title \to description
d_icd_procedures	procedurecodes	row_id \to recordid, icd_code \to icdcode, icd_version \to codeversion, long_title \to description
d_labitems	labtesttypes	$row_id \rightarrow recordid, itemid \rightarrow itemcode, label \rightarrow itemname$
d_items	clinicalitemtypes	row.id \to recordid, itemid \to itemcode, label \to itemname, abbreviation \to abbreviation, linksto \to itemtype
diagnoses_icd	admissiondiagnoses	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, icd_code \rightarrow icdcode, icd_version \rightarrow codeversion, charttime \rightarrow recordeddatetime
procedures_icd	admissionprocedures	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, icd_code \rightarrow icdcode, icd_version \rightarrow codeversion, charttime \rightarrow recordeddatetime
labevents	labresults	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, itemid \rightarrow itemcode, charttime \rightarrow resultdatetime, valuenum \rightarrow resultvalue, valueuom \rightarrow resultunit
prescriptions	medicationorders	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, start time \rightarrow startdatetime, stoptime \rightarrow enddatetime, drug \rightarrow medicationname dose_val_rx \rightarrow dosevalue, dose_unit_rx \rightarrow doseunit, route \rightarrow administrationroute
cost	costrecords	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, event_type \rightarrow eventtype, event_id \rightarrow costid, chargetime \rightarrow costdatetime, cost \rightarrow costamount
chartevents	clinicalevents	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, stay_id \rightarrow icuadmissionid, itemid \rightarrow itemcode, charttime \rightarrow recordeddatetime, valuenum \rightarrow value, valueuom \rightarrow unit
inputevents	intakerecords	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, stay_id \rightarrow icuadmissionid, starttime \rightarrow startdatetime, itemid \rightarrow itemcode, totalamount \rightarrow totalvolume, totalamountuom \rightarrow volumeunit
outputevents	outputrecords	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, stay_id \rightarrow icuadmissionid, charttime \rightarrow recordeddatetime, itemid \rightarrow itemcode, value \rightarrow volume, valueuom \rightarrow volumeunit
microbiologyevents	microbiologyresults	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, charttim \rightarrow collecteddatetime, spec_type_desc \rightarrow specimentype, test_name \rightarrow testname org_name \rightarrow organismname
icustays	icuepisodes	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, stay_id \rightarrow icuadmissionid, first_careunit \rightarrow initialcareunit, last_careunit \rightarrow finalcareunit intime \rightarrow admitdatetime, outtime \rightarrow dischargedatetime
transfers	patienttransfers	row_id \rightarrow recordid, subject_id \rightarrow patientid, hadm_id \rightarrow admissionid, transfer_id \rightarrow transferid, eventtype \rightarrow transfertype, careunit \rightarrow careunit, intime \rightarrow transferindatetime, outtime \rightarrow transferoutdatetime

B.2.2 EICU RENAMING

Table 12: Table and column renaming mappings for eICU.

Original Table	Mapped Table	Column Mappings (eICU to eICU*)
patient	person	uniquepid \rightarrow person_id, patienthealthsystemstayid \rightarrow hosp_id, patientunitstayid \rightarrow unit_id, gender \rightarrow gender, age \rightarrow age, ethnicity \rightarrow ethnicity, hospitalid \rightarrow hospital_id, wardid \rightarrow ward_id, admissionheight \rightarrow height_admission, admissionweight \rightarrow weight_admission, dischargeweight \rightarrow weight_discharge, hospitaladmittime \rightarrow hospital_admit_time, hospitaladmitsource \rightarrow hospital_admission_source, unitadmittime \rightarrow unit_admit_time, unitdischargetime \rightarrow unit_discharge_time, hospitaldischargetime \rightarrow hospital_discharge_time, hospitaldischargestatus \rightarrow hospital_discharge_status
diagnosis	condition	diagnosisid \rightarrow condition_id, patientunitstayid \rightarrow unit_id, diagnosisname \rightarrow condition_name, diagnosistime \rightarrow condition_time, icd9code \rightarrow icd9_code
treatment	treatment	$treatmentid \to treatment_id, \ patientunitstayid \to unit_id, \ treatmentname \to treatment_name, treatmentttime \to treatment_time$
lab	lab	labid \to lab_id, patientunitstayid \to unit_id, labname \to lab_name, labresult \to lab_result, labresulttime \to lab_result_time
medication	prescription	$\label{eq:medication} \begin{array}{ll} medicationid \rightarrow prescription_id, \ patientunitstayid \rightarrow unit_id, \ drugname \rightarrow \\ drug_name, \ dosage \rightarrow dosage, \ routeadmin \rightarrow administration_route, \ drugstarttime \rightarrow medication_start_time, \ drugstoptime \rightarrow medication_stop_time \end{array}$
cost	cost	costid \rightarrow cost_id, patienthealthsystemstayid \rightarrow hosp_id, patientunitstayid \rightarrow unit_id, eventtype \rightarrow event_type, eventid \rightarrow event_id, chargetime \rightarrow cost_time. cost \rightarrow cost_amount
allergy	allergy_reaction	allergyid \to allergy_id, patientunitstayid \to unit_id, drugname \to drug_name, allergyname \to allergy_name, allergytime \to allergy_time
intakeoutput	fluid_balance	intakeoutputid \rightarrow fluid_balance_id, patientunitstayid \rightarrow unit_id, cellpath \rightarrow fluid_path, celllabel \rightarrow fluid_label, cellvaluenumeric \rightarrow fluid_value_numeric, intakeoutputtime \rightarrow fluid_balance_time
microlab	microbiology	$\label{eq:microbiology} \begin{tabular}{l} microbiology_id, patientunitstayid \rightarrow unit_id, culturesite \rightarrow culture_site, organism \rightarrow organism, culturetakentime \rightarrow culture_taken_time \rightarrow culture_taken_ti$
vitalperiodic	vital_signs	vitalperiodicid \rightarrow vital_sign.id, patientunitstayid \rightarrow unit_id, temperature \rightarrow temperature, sao2 \rightarrow sao2, heartrate \rightarrow heart_rate, respiration \rightarrow respiration_rate, systemicsystolic \rightarrow systolic_bp, systemicdiastolic \rightarrow diastolic_bp systemicmean \rightarrow mean_bp, observationtime \rightarrow vital_time
hospital	hospital	$hospitalid \to hospital_id, numbedscategory \to bed_capacity_category, teachingstatus \to teaching_status, region \to region$

B.3 ADAPTQA TEMPLATES

Table 13: Eight categories of query goal modification in AdaptQA.

Category	Description
Medication nomenclature traversal	This category involves queries that navigate the medication nomenclature, such as from a brand name to a generic name or vice versa. The tasks begin by requesting information about a specific medication name (e.g., "Lipitor") and then, if that name is not found, adapt the query to search for its generic name, active ingredient, or other brand names within the same class.
Within drug class adaptation	This category encompasses queries that require a user to adapt their search within a single drug class. The query goals might be to expand from a specific drug to the entire class, narrow the search from a broad class to a specific subclass, or exclude a particular medication from a class. This flexibility demonstrates a deeper understanding of therapeutic classifications beyond simple name-to-name conversion.
Across drug class traversal	This category is defined by queries that involve navigating between different drug classes. The tasks may require the identification of medications from one class while excluding another (e.g., finding PUD medications other than PPIs) or the combination of multiple distinct drug classes in a single query (e.g., patients prescribed both ACE inhibitors and Beta-blockers).
Primary condition to related condition	This category includes tasks where the initial query focuses on a primary diagnosis or procedure, and the subsequent goal is to identify clinically related conditions. This often involves looking for common comorbidities, complications from a surgery, or side effects of a medication.
Alternative lab test for condition	This category is designed for scenarios where the initial query for a specific lab test is unsuccessful. The system must then identify and pivot to a clinically relevant alternative lab test used to assess the same condition (e.g., from "troponin I" to "troponin T" for myocardial damage). This mimics clinical reasoning when a preferred test is unavailable.
Alternative procedure for treatment	This category handles tasks that require finding alternative procedures when a primary treatment method is not found or is not applicable. The queries start by searching for a specific procedure for a condition and then, if necessary, adjust to look for other clinically appropriate procedures or surgeries for the same condition.
Multi-criteria resolution	This category involves complex queries that require resolving multiple, often compound, criteria simultaneously. The tasks integrate various data points—such as diagnoses, lab results, and vital signs—using logical operators (AND/OR) to identify a specific patient cohort. The complexity lies in accurately interpreting and executing these intricate, rule-based queries.
Schema fallback handling	This category addresses situations where the primary data source (a specific database table or schema) is unavailable. The query system must then "fall back" to an alternative data source or a modified search strategy to fulfill the user's goal. This demonstrates robustness in handling missing data schema by adapting the query to existing information (e.g., a "cancer registry table" vs. the general "diagnosis table").

C AGENT IMPLEMENTATION

1350

1351 1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1363

1364

For selecting the backbone LLMs for the agents, small models (e.g., 7B or 13B), DeepSeek-R1, and the Gemma 3 series are excluded due to their limited performance in tool invocation tasks. Owing to budget constraints, we also exclude Anthropic models (e.g., Opus 4 and Sonnet) as well as Gemini-2.5-Pro.

Table 14 presents the system prompt used for our agent baselines. In addition to the system prompt, the agent's input includes three other components: (1) agent behavioral rules detailing interaction behavior (Section C.2); (2) evaluation rules for IncreQA and AdaptQA to ensure accurate evaluation (Section A.5); and (3) database-specific rules outlining the SQL annotation assumptions for MIMIC-IV* and eICU* (Section A.4).

C.1 AGENT SYSTEM PROMPT

Table 14: Agent system prompt.

```
1365
         Instruction:
1367
         - You are a DB agent that helps users by answering their questions
         in natural language using information from a database.
         - You are currently engaged in a conversation with a user who wants
1369
         to retrieve some data or statistics from an EHR database.
1370
         - If the user's request is ambiguous or lacks important details (e.g
1371
            filtering criteria), ask clarifying questions to better
         understand the request.
1372
         - You have access to a set of tools to assist the user:
1373
          - table_search: search tables in the database
1374
          - column_search: search columns in a table
          - value_substring_search: search values in a column by substring
1375
          match
1376
          - value_similarity_search: search values in a column by semantic
          similarity (embedding-based)
          - sql_execute: execute an SQL query on the database
          - web_search: search the web for external clinical knowledge
         - Use table_search and column_search to explore the database schema.
1380
          Use value_substring_search and value_similarity_search to explore
         stored values.
1381
         - Clinical concepts (e.g., diagnoses, procedures, medications, lab
1382
         tests) in the database may not exactly match the user's words. Use
         the value search tools to find relevant entries.
         - If you need clinical knowledge beyond what is in the database (e.g
1384
         ., a drug's mechanism of action), use web_search.
          Never invent or assume information that is not provided by the
1386
         user or retrieved using the tools.
         - Make only one tool call at a time. Do not send a user-facing
1387
         response in the same turn as a tool call.
1388
          After gathering all necessary information, use sql_execute to
1389
         write and run a single valid SQL query that fully answers the user's
          latest request.
1390
         - When you write an SQL query, always execute it with sql_execute
1391
         and return the results to the user along with your explanation.
1392
         {agent_rule}
1393
1394
         {database rule}
1395
```

C.2 AGENT RULES

1398

1400 1401

1402

Table 15: Agent rules.

```
Below are the general rules for the DB agent:
- The DB agent must assume the user has no knowledge of SQL, databases, or stored values, and cannot execute queries.
```

- The DB agent must interact with the user only in natural language and must not show raw SQL queries. - The DB agent must not modify the database schema or contents. The following commands are forbidden: INSERT, UPDATE, DELETE, DROP, - The DB agent must write queries that finish within 60 seconds; otherwise, the query results will be invalid. - The DB agent must limit each conversation to 30 interactions (including user exchanges and tool calls) and 600 seconds total. The DB agent must always explain answers in natural language, including the reasoning or conditions used to arrive at those answers. If SQL references are necessary, the DB agent must explain them in terms understandable to someone with no SQL knowledge. - The DB agent must clearly explain when a question cannot be answered (e.g., due to limitations of SQL or empty results) and ask the user to rephrase or modify the request.

C.3 AGENT RULES FOR EVALUATION

Table 16 and Table 17 present the prompts used for IncreQA and AdaptQA tasks, respectively. These prompts are appended to the agent rules to guide the agent toward behavior that aligns with accurate evaluation criteria. Note that for both task types, the agent's responses are evaluated each time they are generated, as correct answers may emerge during the conversation due to iterative nature of open QA dialogues.

Table 16: Scoring rules for IncreQA.

```
Below are the grading rules:

- The DB agent's performance is evaluated based on the generated SQL queries, requiring at least one SQL query (via sql_execute) to retrieve answers during the interaction with the database.

- For accurate assessment, when the user revises their question, the DB agent must write a new SQL query from scratch to fully address the latest request, without relying on previous query results.

- For questions that involve calculations (such as time differences or survival rates) or data manipulation/aggregation, the DB agent must use SQL language to compute the results rather than relying on its LLM capabilities.
```

Table 17: Scoring rules for AdaptQA.

```
Below are the grading rules:

- The DB agent's performance is evaluated based on the results in its natural language response to the user.

- When providing answers to the user, the DB agent must enclose the final answer in <answer></answer> tags (e.g., <answer>42</answer>). All other content, including intermediate results, explanations, value units, and any additional details, must be placed outside these tags.

- When answers are textual data (e.g., timestamps or diagnosis names), use them exactly as stored in the database. For numerical answers, round them to four decimal places.
```

```
1458
             SAMPLE TRAJECTORIES
1459
1460
        In this section, we present sample task instances along with one successful and one failed conversa-
1461
        tion trajectory for IncreQA and AdaptQA.
1462
1463
        D.1 INCREMENTAL QUERY REFINEMENT
1464
1465
        D.1.1 Successful Case (IncreQA)
1466
        The agent successfully retrieves the number of patients diagnosed with breast cancer by querying
1467
        values in the description column of the diagnosiscodes table.
1468
1469
           User: I'm trying to figure out how many patients admitted to the hospital have certain
1470
1471
           conditions.
1472
1473
           Agent: table_search(arguements={"tool_input":""})
1474
1476
           Tool: admissiondiagnoses, admissionprocedures, clinicalevents, clinicalitemtypes,
1477
           costrecords, demographics, diagnosiscodes, hospitaladmissions, icuepisodes,
1478
           intakerecords, labresults, labtesttypes, medicationorders, microbiologyresults,
1479
           outputrecords, patienttransfers, procedurecodes
1480
1481
1482
           Agent: column_search(arguements={"table_names":"admissiondiagnoses,diagnosiscodes"})
1483
1484
           Tool: CREATE TABLE admissiondiagnoses (
1485
1486
                   recordid INTEGER NOT NULL,
1487
                   patientid INTEGER NOT NULL,
1488
                   admissionid INTEGER NOT NULL,
1489
                   icdcode TEXT NOT NULL,
1490
1491
                   codeversion INTEGER NOT NULL,
1492
                   recordeddatetime TEXT NOT NULL,
1493
                   PRIMARY KEY (recordid),
1494
                   FOREIGN KEY (codeversion) REFERENCES diagnosiscodes (codeversion),
1495
1496
                   FOREIGN KEY (icdcode) REFERENCES diagnosiscodes (icdcode),
1497
                   FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
1498
           )
1499
           /*
1500
1501
           3 rows from admissiondiagnoses table:
1502
           recordid patientid admissionid icdcode codeversion recordeddatetime
1503
           0 10035185 22580999 4139 9 2100-05-17 12:53:00
1504
           1 10035185 22580999 V707 9 2100-05-17 12:53:00
1505
1506
           2 10035185 22580999 41401 9 2100-05-17 12:53:00
1507
           */
1509
           CREATE TABLE diagnosiscodes (
1510
1511
                   recordid INTEGER NOT NULL,
```

```
1512
                   icdcode TEXT NOT NULL,
1513
1514
                   codeversion INTEGER NOT NULL,
1515
                   description TEXT NOT NULL,
1516
                   PRIMARY KEY (recordid)
1517
           )
1518
           /*
1519
1520
           3 rows from diagnosiscodes table:
1521
           recordid icdcode codeversion description
1522
           0 0090 9 infectious colitis, enteritis, and gastroenteritis
1523
           1 01160 9 tuberculous pneumonia [any form], unspecified
1524
1525
           2 01186 9 other specified pulmonary tuberculosis, tubercle bacilli not found by
1526
           bacteriological or histological examination, but tuberculosis confirmed by other methods [
1527
           inoculation of animals]
           */
1529
1530
1531
           Agent: Could you clarify which specific conditions you'd like to count? For example, are you
1532
           interested in diabetes, heart failure, COPD, etc.? Also, should these diagnoses be those
1533
           recorded at the time of hospital admission?
1534
1535
1536
           User: I want to know how many patients admitted in the last 90 days have a family history
1537
           of breast cancer.
1538
1539
1540
           Agent: value_substring_search(arguements={"table":"diagnosiscodes","column":"
1541
           description","value":"family"})
1542
1543
           Tool: Values in diagnosiscodes.description containing 'family': ['family history of polycystic
1545
           kidney', 'non-family member, perpetrator of maltreatment and neglect', 'family history of
           unspecified malignant neoplasm', ..., 'family history of eye disorders', 'family disruption due
1547
            to death of family member'].
1548
1549
1550
           Agent: column_search(arguments={"table_names":"hospitaladmissions"})
1551
1552
           Tool: CREATE TABLE hospitaladmissions (
1553
                   recordid INTEGER NOT NULL,
1554
1555
                   patientid INTEGER NOT NULL,
1556
                   admissionid INTEGER NOT NULL,
1557
                   admitdatetime TEXT NOT NULL,
                   dischargedatetime TEXT,
1560
                   admissiontype TEXT NOT NULL,
1561
                   admitsource TEXT NOT NULL,
                   dischargedestination TEXT,
1563
                   insurancetype TEXT NOT NULL,
1564
```

```
1566
                   language TEXT,
1567
1568
                   maritalstatus TEXT,
1569
                   age INTEGER NOT NULL,
1570
                   PRIMARY KEY (recordid),
1571
                   FOREIGN KEY (patientid) REFERENCES demographics (patientid)
1572
           )
1573
           /*
1574
1575
           3 rows from hospitaladmissions table:
1576
           recordid patientid admissionid admitdatetime dischargedatetime admissiontype
1577
           admitsource dischargedestination insurancetype language maritalstatus age
           o 10004235 24181354 2100-03-19 14:38:00 2100-03-28 14:02:00 urgent transfer from hospital
1580
            skilled nursing facility medicaid english single 47
1581
           1 10009628 25926192 2100-09-30 17:08:00 2100-10-08 13:20:00 urgent transfer from
           hospital home health care medicaid? married 58
1584
           2 10006053 22942076 2100-11-15 23:39:00 2100-11-17 17:20:00 urgent transfer from hospital
1585
           died medicaid english None 52
1586
           */
1587
           Agent: value_substring_search(arguements={"table":"diagnosiscodes","column":"
1590
           description","value":"breast"})
1591
1592
           Tool: Values in diagnosiscodes.description containing 'breast': ['other and unspecified
1593
1594
           disorder of breast associated with childbirth, postpartum condition or complication', '
           malignant neoplasm of axillary tail of female breast', 'full-thickness skin loss [third degree,
1596
            not otherwise specified] of chest wall, excluding breast and nipple', ...', 'personal history of
1597
            malignant neoplasm of breast', 'open wound of breast, without mention of complication'].
1598
           Agent: sql_execute(arguements={"query": "SELECT COUNT(DISTINCT ha.patientid) AS
           num_patients FROM hospitaladmissions AS ha JOIN admissiondiagnoses AS ad ON ha.
           admissionid = ad.admissionid JOIN diagnosiscodes AS dc ON ad.icdcode = dc.icdcode AND
1604
           ad.codeversion = dc.codeversion WHERE dc.description = 'family history of malignant
           neoplasm of breast' AND ha.admitdatetime >= datetime('2100-12-31 23:59:00', '-90 days')
           ;"})
1608
1609
           Tool: [(1,)]
1610
1611
           Agent: There have been 1 unique patient admitted in the last 90 days with a diagnosis of
1612
           family history of malignant neoplasm of breast.
1613
1614
1615
           User: Can you double check that the answer covers my goals?
1616
1617
           Agent: Yes-here's how the query addresses each part of your goal:
1618
```

- Last 90 days: We filtered admissions with ha.admitdatetime >= datetime('2100-12-31 23:59:00','-90 days') - Family history of breast cancer: We joined admissions -> diagnoses -> codes and used dc.description = 'family history of malignant neoplasm of breast' - Unique patients: We counted DISTINCT ha.patientid so repeat admissions aren't double-counted. The result was 1 patient. Let me know if you'd like any further breakdown or checks! User: ###END###

D.1.2 FAILED CASE (INCREQA)

 The agent fails to retrieve information on breast cancer, as it queries unrelated tables and columns (e.g., clinicalevents and clinicalitemtypes) instead of the diagnosiscodes table. Due to the use of FAISS, wrapped in the LangChain library, values with semantic similarity below the default threshold (0.8) are not retrieved. Consequently, values below this threshold are not retrieved. In the example below, the agent applies the value similarity tool to the itemname column of the clinicalitemtypes table for "cancer," but the tool returns no result. Note that itemname corresponds to the label in the d_item table in MIMIC-IV, which covers a wide range of clinical practices, vital signs, and measurements.

User: I want to find out how many patients have been admitted to the hospital recently and also have a family history of breast cancer.

Agent: To assist you with your request, I need to clarify a couple of things:

- 1. What do you mean by "recently"? Are you looking for admissions within a specific time frame (e.g., the last month, last year)?
- 2. How is the family history of breast cancer recorded in the database? Is there a specific table or column that indicates this information?

Once I have this information, I can proceed to gather the necessary data from the database

User: I'm interested in patients admitted in the last 90 days. As for the family history, could you check the database to find where that info is stored?

Agent: table_search(arguements={"tool_input":""})

Tool: admissiondiagnoses, admissionprocedures, clinicalevents, clinicalitemtypes, costrecords, demographics, diagnosiscodes, hospitaladmissions, icuepisodes, intakerecords, labresults, labtesttypes, medicationorders, microbiologyresults, outputrecords, patienttransfers, procedurecodes

Agent: table_names(arguements={"table_names":"demographics"})

```
Tool: CREATE TABLE demographics (
recordid INTEGER NOT NULL,
patientid INTEGER NOT NULL,
gender TEXT NOT NULL,
dateofbirth TEXT NOT NULL,
dateofdeath TEXT,
PRIMARY KEY (recordid)
)
/*
3 rows from demographics table:
recordid patientid gender dateofbirth dateofdeath
0 10014729 f 2079-07-22 00:00:00 None
```

```
1728
           1 10003400 f 2028-07-25 00:00:00 None
1729
1730
           2 10002428 f 2020-07-30 00:00:00 None
1731
           */
1732
1733
           Agent: column_search(arguements={"table_names":"clinicalevents"})
1734
1735
1736
           Tool: CREATE TABLE clinicalevents (
1737
                   recordid INTEGER NOT NULL,
1738
                   patientid INTEGER NOT NULL,
1739
                  admissionid INTEGER NOT NULL,
1740
1741
                  icuadmissionid INTEGER NOT NULL,
1742
                  itemcode INTEGER NOT NULL,
1743
                   recordeddatetime TEXT NOT NULL,
1744
                  value DOUBLE PRECISION,
1745
1746
                   unit TEXT,
1747
                   PRIMARY KEY (recordid),
1748
                   FOREIGN KEY (itemcode) REFERENCES clinicalitemtypes (itemcode),
1749
                   FOREIGN KEY (icuadmissionid) REFERENCES icuepisodes (icuadmissionid),
1750
1751
                   FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
1752
           )
1753
           /*
1754
           3 rows from clinical events table:
1755
1756
           recordid patientid admissionid icuadmissionid itemcode recordeddatetime value unit
1757
           0 10005817 20626031 32604416 220210 2100-12-24 00:00:00 19.0 insp/min
1758
           1 10005817 20626031 32604416 220051 2100-12-24 00:00:00 37.0 mmhg
1759
           2 10005817 20626031 32604416 223762 2100-12-24 00:00:00 37.0 c
1760
1761
           */
1762
1763
           Agent: column_search(arguements={"table_names":"clinicalitemtypes"})
1764
1765
1766
           Tool: CREATE TABLE clinicalitemtypes (
1767
                   recordid INTEGER NOT NULL,
1768
                   itemcode INTEGER NOT NULL,
1769
                   itemname TEXT NOT NULL,
1770
1771
                  abbreviation TEXT NOT NULL,
1772
                   itemtype TEXT NOT NULL,
1773
                   PRIMARY KEY (recordid)
1774
           )
1775
1776
           /*
1777
           3 rows from clinicalitemtypes table:
1778
           recordid itemcode itemname abbreviation itemtype
1779
           o 226228 gender gender clinicalevents
1780
```

```
1782
           1 226545 race race clinical events
1783
           2 229877 suction events (ch) suction events (ch) clinicalevents
1784
1785
           */
1786
1787
           Agent: value_substring_search(arguements={"table":"clinicalitemtypes","column":"
1788
           itemname","value":"family history"})
1789
1790
1791
           Tool: No values in clinicalitemtypes.itemname contain 'family history'.
1792
1793
           Agent: value_similarity_search(arguements={"table":"clinicalitemtypes","column":"
1794
           itemname","value":"family history"})
1796
1797
           Tool: No matches found in clinicalitemtypes.itemname for 'family history'.
1798
1800
           Agent: column_search(arguements={"table_names":"hospitaladmissions"})
1801
1802
           Tool: CREATE TABLE hospitaladmissions (
1803
                   recordid INTEGER NOT NULL,
1805
                   patientid INTEGER NOT NULL,
1806
                   admissionid INTEGER NOT NULL,
1807
                   admitdatetime TEXT NOT NULL,
1808
                   dischargedatetime TEXT,
1809
1810
                   admissiontype TEXT NOT NULL,
1811
                   admitsource TEXT NOT NULL,
1812
                   dischargedestination TEXT,
1813
                   insurancetype TEXT NOT NULL,
1815
                   language TEXT,
1816
                   maritalstatus TEXT,
1817
                   age INTEGER NOT NULL,
1818
                   PRIMARY KEY (recordid).
1819
1820
                   FOREIGN KEY (patientid) REFERENCES demographics (patientid)
1821
           )
           /*
           3 rows from hospitaladmissions table:
1824
1825
           recordid patientid admissionid admitdatetime dischargedatetime admissiontype
1826
           admitsource dischargedestination insurancetype language maritalstatus age
1827
           o 10004235 24181354 2100-03-19 14:38:00 2100-03-28 14:02:00 urgent transfer from hospital
1828
            skilled nursing facility medicaid english single 47
1829
1830
           1 10009628 25926192 2100-09-30 17:08:00 2100-10-08 13:20:00 urgent transfer from
1831
           hospital home health care medicaid? married 58
           2 10006053 22942076 2100-11-15 23:39:00 2100-11-17 17:20:00 urgent transfer from hospital
1833
           died medicaid english None 52
1834
```

```
1836
           */
1838
           Agent: column_search(arguements={"table_names":"demographics"})
1840
1841
           Tool: CREATE TABLE demographics (
1842
                   recordid INTEGER NOT NULL,
1843
1844
                   patientid INTEGER NOT NULL,
1845
                   gender TEXT NOT NULL,
1846
                   dateofbirth TEXT NOT NULL,
1847
                   dateofdeath TEXT,
1848
1849
                   PRIMARY KEY (recordid)
1850
           )
1851
           /*
1852
           3 rows from demographics table:
1853
           recordid patientid gender dateofbirth dateofdeath
1855
           0 10014729 f 2079-07-22 00:00:00 None
           1 10003400 f 2028-07-25 00:00:00 None
1857
           2 10002428 f 2020-07-30 00:00:00 None
1859
           */
1860
1861
           Agent: column_search(arguements={"table_names":"clinicalevents"})
1862
1863
1864
1865
           Tool: CREATE TABLE clinicalevents (
1866
                   recordid INTEGER NOT NULL,
1867
                   patientid INTEGER NOT NULL,
                   admissionid INTEGER NOT NULL,
                   icuadmissionid INTEGER NOT NULL,
1871
                   itemcode INTEGER NOT NULL,
1872
                   recordeddatetime TEXT NOT NULL,
1873
1874
                   value DOUBLE PRECISION,
1875
                   unit TEXT,
1876
                   PRIMARY KEY (recordid),
1877
                   FOREIGN KEY (itemcode) REFERENCES clinicalitemtypes (itemcode),
1878
1879
                   FOREIGN KEY (icuadmissionid) REFERENCES icuepisodes (icuadmissionid),
1880
                   FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
1881
           )
1882
           /*
1883
1884
           3 rows from clinical events table:
1885
           recordid patientid admissionid icuadmissionid itemcode recordeddatetime value unit
           0 10005817 20626031 32604416 220210 2100-12-24 00:00:00 19.0 insp/min
1887
           1 10005817 20626031 32604416 220051 2100-12-24 00:00:00 37.0 mmhg
1888
1889
```

2 10005817 20626031 32604416 223762 2100-12-24 00:00:00 37.0 C */ Agent: value_substring_search(arguements={"table":"clinicalitemtypes","column":" itemname","value":"cancer"}) **Tool**: No values in clinicalitemtypes.itemname contain 'cancer'. Agent: value_similarity_search({"table":"clinicalitemtypes","column":"itemname","value":" cancer"}) **Tool**: No matches found in clinicalitemtypes.itemname for 'cancer'. Agent: It seems that I couldn't find any specific entries related to "family history of breast cancer" in the database. This could mean that the information is either stored in a different format or not recorded in the tables I checked. To proceed, I can look for any general family history or cancer-related entries in the clinical events or demographics tables. Would you like me to do that? **User**: Yes, please check for any general family history or cancer-related info. Agent: value_substring_search(arguements={"table": "clinicalitemtypes", "column": " itemname", "value": "family history" }) **Tool**: No values in clinicalitemtypes.itemname contain 'family history'. Agent: value_substring_search(arguements={"table":"clinicalitemtypes","column":" itemname","value":"cancer"}) **Tool**: No values in clinicalitemtypes.itemname contain 'cancer'. Agent: I have searched for general family history and cancer-related information in the database, but unfortunately, I couldn't find any relevant entries in the clinical item types. It appears that the information regarding family history of breast cancer may not be explicitly recorded in the database, or it might be stored in a different format or table that I haven't checked. If you have any other specific criteria or different information you would like to explore, please let me know! User: Okay, can you at least tell me how many patients have been admitted in the last 90 days?

Agent: sql_execute(arguemetns={"query": "SELECT COUNT(DISTINCT patientid) AS patient_count FROM hospitaladmissions WHERE admitdatetime >= '2100-09-30'"}) **Tool**: [(38,)] Agent: There have been 38 patients admitted to the hospital in the last 90 days. If you have any other questions or need further assistance, feel free to ask! User: Okay, so you couldn't find the family history info. Just to confirm, you were able to find the number of patients admitted in the last 90 days? Agent: Yes, that's correct! I was able to find that there have been 38 patients admitted to the hospital in the last 90 days. However, I couldn't find any specific information regarding family history of breast cancer in the database. If you have any other questions or need further assistance, please let me know! User: Double check that you've covered my goals: the number of patients admitted in the last 90 days with a family history of breast cancer. Agent: I understand your request clearly now. You want to know the number of patients who were admitted in the last 90 days and also have a family history of breast cancer. Unfortunately, I was unable to find any specific entries related to family history of breast cancer in the database. This means I cannot confirm how many patients meet both criteria. If you have any other specific information or different criteria you would like to explore, please let me know! **User:** ###END###

D.2 ADAPTIVE QUERY REFINEMENT

D.2.1 SUCCESSFUL CASE (ADAPTQA)

The user initially asks whether a specific patient was prescribed carbamazepine for epilepsy treatment, but the agent finds that no such records exist for the patient. The user then requests information about other medications for epilepsy. The agent searches the database using clinical knowledge with keywords such as levetiracetam, phenytoin, and phenobarbital. The agent checks whether the patient was prescribed any of these medications and successfully retrieves the most recent prescription date for levetiracetam, as requested.

User: I'm looking into patient prescriptions. Can you help me find some info about the medications prescribed to patient ID 10008287?

```
Agent: table_search(arguements={"tool_input":""})
```

Tool: admissiondiagnoses, admissionprocedures, clinicalevents, clinicalitemtypes, costrecords, demographics, diagnosiscodes, hospitaladmissions, icuepisodes, intakerecords, labresults, labtesttypes, medicationorders, microbiologyresults, outputrecords, patienttransfers, procedurecodes

Agent: column_search(arguements={"table_names":"medicationorders"})

```
Tool: CREATE TABLE medicationorders (
2025
2026
               recordid INTEGER NOT NULL,
2027
               patientid INTEGER NOT NULL,
2028
               admissionid INTEGER NOT NULL,
2029
               startdatetime TEXT NOT NULL,
2030
2031
               enddatetime TEXT.
2032
               medicationname TEXT NOT NULL,
2033
               dosevalue TEXT NOT NULL,
               doseunit TEXT NOT NULL.
2035
2036
               administrationroute TEXT NOT NULL,
2037
               PRIMARY KEY (recordid),
2038
               FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
2039
           )
2040
2041
           /*
2042
           3 rows from medicationorders table:
2043
           recordid patientid admissionid startdatetime enddatetime medicationname dosevalue
2044
           doseunit administrationroute
2045
2046
           o 10020740 23831430 2100-04-19 11:00:00 2100-04-20 22:00:00 insulin o.o unit sc
2047
           1 10005817 20626031 2100-12-20 02:00:00 2100-12-21 12:00:00 insulin o.o unit sc
2048
           2 10005817 20626031 2100-12-27 10:00:00 2100-12-28 19:00:00 insulin 0.0 unit sc
2049
           */
2050
```

```
2052
           Agent: Could you clarify which details you'd like to see for patient 10008287's prescriptions
2053
2054
           ? For example, are you interested in the medication names only, the start/end dates,
2055
           dosages, administration routes, or all of those?
2056
2057
           User: I'm particularly interested in whether patient 10008287 was prescribed
2058
           carbamazepine for epilepsy. If so, I want to know the most recent time it was prescribed,
2060
           but only after their first epilepsy diagnosis.
2061
2062
           Agent: value_substring_search(arguements={"table":"medicationorders","column":"
2063
           medicationname","value":"carbamazepine","k":10})
2064
2065
2066
           Tool: Values in medicationorders.medicationname containing 'carbamazepine': ['
2067
           carbamazepine'].
2068
2069
2070
           Agent: column_search(arguements={"table_names":"admissiondiagnoses"})
2071
           Tool: CREATE TABLE admissiondiagnoses (
2073
               recordid INTEGER NOT NULL,
2074
2075
               patientid INTEGER NOT NULL,
2076
               admissionid INTEGER NOT NULL,
2077
               icdcode TEXT NOT NULL,
2078
               codeversion INTEGER NOT NULL,
2079
2080
               recordeddatetime TEXT NOT NULL,
2081
               PRIMARY KEY (recordid),
2082
               FOREIGN KEY (codeversion) REFERENCES diagnosiscodes (codeversion),
2083
               FOREIGN KEY (icdcode) REFERENCES diagnosiscodes (icdcode),
2084
2085
               FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
2086
           )
           /*
           3 rows from admissiondiagnoses table:
2089
2090
           recordid patientid admissionid icdcode codeversion recordeddatetime
2091
           0 10035185 22580999 4139 9 2100-05-17 12:53:00
2092
           1 10035185 22580999 V707 9 2100-05-17 12:53:00
2093
           2 10035185 22580999 41401 9 2100-05-17 12:53:00
2094
2095
           */
2096
2097
           Agent: column_search(arguements={"table_names":"diagnosiscodes"})
2098
2099
2100
           Tool: CREATE TABLE diagnosiscodes (
2101
               recordid INTEGER NOT NULL,
2102
               icdcode TEXT NOT NULL,
2103
               codeversion INTEGER NOT NULL,
2104
```

```
2106
                description TEXT NOT NULL,
2107
                PRIMARY KEY (recordid)
2108
2109
           )
2110
            /*
2111
           3 rows from diagnosiscodes table:
2112
           recordid icdcode codeversion description
2113
2114
           0 0090 9 infectious colitis, enteritis, and gastroenteritis
2115
           1 01160 9 tuberculous pneumonia [any form], unspecified
2116
           2 01186 9 other specified pulmonary tuberculosis, tubercle bacilli not found by
2117
           bacteriological or histological examination, but tuberculosis confirmed by other methods [
2118
2119
           inoculation of animals]
2120
           */
2121
2122
2123
```

Agent: value_substring_search(arguements={"table":"diagnosiscodes","column":" description","value":"epilepsy","k":10})

Tool: Values in diagnosiscodes.description containing 'epilepsy': ['epilepsy complicating pregnancy, childbirth, or the puerperium, delivered, with mention of postpartum complication', 'epilepsy, unspecified, without mention of intractable epilepsy', 'epilepsia partialis continua, with intractable epilepsy', 'epilepsy complicating pregnancy, childbirth, or the puerperium, unspecified as to episode of care or not applicable', 'localization-related (focal) (partial) epilepsy and epileptic syndromes with complex partial seizures, with intractable epilepsy', 'generalized convulsive epilepsy, with intractable epilepsy', 'localization-related (focal) (partial) epilepsy and epileptic syndromes with complex partial seizures, without mention of intractable epilepsy', 'localization-related (focal) (partial) epilepsy and epileptic syndromes with simple partial seizures, without mention of intractable epilepsy', 'epilepsia partialis continua, without mention of intractable epilepsy', 'epilepsy complicating pregnancy, childbirth, or the puerperium, postpartum condition or complication'].

Tool: sql_execute(arguements={"query": "SELECT MAX(mo.startdatetime) AS most_recent_prescription FROM medicationorders mo WHERE mo.patientid = 10008287 AND mo.medicationname = 'carbamazepine' AND mo.startdatetime > (SELECT MIN(ad. recordeddatetime) FROM admissiondiagnoses ad JOIN diagnosiscodes dc ON ad.icdcode = dc.icdcode AND ad.codeversion = dc.codeversion WHERE ad.patientid = 10008287 AND dc. description LIKE '%epilepsy%');"})

Tool: [(None,)]

Agent: sql_execute(arguements={"query": "SELECT ad.recordeddatetime AS diagnosis_time FROM admissiondiagnoses ad JOIN diagnosiscodes dc ON ad.icdcode = dc.icdcode AND ad.

codeversion = dc.codeversion WHERE ad.patientid = 10008287 AND dc.description LIKE '% epilepsy%' ORDER BY ad.recordeddatetime ASC LIMIT 1;"}) **Tool**: [('2100-10-09 01:17:00',)] Agent: sql_execute(arguements={"query": "SELECT * FROM medicationorders WHERE patientid = 10008287;"}) Tool: [(354, 10008287, 22168393, '2100-10-10 09:00:00', '2100-10-10 10:00:00', 'bag', '1.0', ' bag', 'iv'), (356, 10008287, 22168393, '2100-10-10 07:00:00', '2100-10-10 08:00:00', 'bag', '1.0', 'bag', 'iv'), (655, 10008287, 22168393, '2100-10-09 22:00:00', '2100-10-11 01:00:00', 'cefazolin ', '1.0', 'g', 'iv'), ..., (6797, 10008287, 22168393, '2100-10-09 22:00:00', '2100-10-10 10:00:00', '0.9% sodium chloride', '1000.0', 'ml', 'iv'), (6804, 10008287, 22168393, '2100-10-09 22:00:00', '2100-10-10 10:00:00', '0.9% sodium chloride', '1000.0', 'ml', 'iv')] Agent: sql_execute(arguements={"query": "SELECT COUNT(*) AS count_prescriptions FROM medicationorders WHERE patientid = 10008287 AND medicationname = 'carbamazepine';"}) **Tool**: [(0,)] Agent: The answer is <answer>No, patient 10008287 was not prescribed carbamazepine after their first epilepsy diagnosis.</answer> **User:** Ok, if carbamazepine wasn't prescribed, can you check if they were prescribed any other medications for epilepsy after their first epilepsy diagnosis? If so, give me the most recent prescription time. Agent: sql_execute(arguements={"query": "WITH first_epi AS (SELECT MIN(ad. recordeddatetime) AS dt FROM admissiondiagnoses ad JOIN diagnosiscodes dc ON ad. icdcode = dc.icdcode AND ad.codeversion = dc.codeversion WHERE ad.patientid = 10008287 AND dc.description LIKE '%epilepsy%') SELECT MAX(m.startdatetime) AS most_recent_prescription FROM medicationorders m JOIN first_epi f ON m.startdatetime > f. dt WHERE m.patientid = 10008287 AND (m.medicationname LIKE '%levetiracetam%' OR m. medicationname LIKE '%phenytoin%' OR m.medicationname LIKE '%phenobarbital%' OR m. medicationname LIKE '%valpro%' OR m.medicationname LIKE '%lamotr%' OR m. medicationname LIKE '%gabapentin%' OR m.medicationname LIKE '%topir%' OR m. medicationname LIKE '%pregabat%' OR m.medicationname LIKE '%ethosux%' OR m. medicationname LIKE '%lacosamide%');"}) **Tool**: [('2100-10-12 20:00:00',)]

Agent: The answer is <answer>2100-10-12 20:00:00</answer>

User: Great! Just to be sure, this time represents the most recent time patient 10008287 was prescribed any medication for epilepsy after their first diagnosis? I need the answer in the exact time format recorded in the database. Agent: Yes, that is correct. The most recent epilepsy medication prescription time after their first epilepsy diagnosis is <answer>2100-10-12 20:00:00</answer>. **User**: Perfect, thanks! ###END###

D.2.2 FAILED CASE (ADAPTQA)

After the user initially asks about carbamazepine for epilepsy treatment for a specific patient, the agent finds no relevant data. The user then asks about other medications for epilepsy. The agent uses a web search tool to retrieve a list of medications, including levetiracetam. The agent writes an SQL query using the diagnosis "epilepsy, unspecified" and the medication list for the patient. However, since the patient was diagnosed with "epilepsy, unspecified, not intractable, without status epilepticus," the agent fails to retrieve the specific data needed, despite having the list of epilepsy medications.

User: I want to know about the prescriptions for patient ID 10008287.

Agent: Could you please specify what information about the prescriptions you are interested in? For example, are you looking for the medication names, dosages, dates, or something else?

User: I'm interested in finding out if patient 10008287 was prescribed carbamazepine for epilepsy. If they were, I need the most recent prescription time after their first epilepsy diagnosis.

Agent: Okay, I need to find out if patient 10008287 has been diagnosed with epilepsy and prescribed carbamazepine. Then, if both conditions are met, I need to find the most recent prescription time for carbamazepine after the first epilepsy diagnosis. First, I need to find the relevant tables containing diagnoses and medication orders.

Agent: table_search(arguements={"tool_input": ""})

Tool: admissiondiagnoses, admissionprocedures, clinicalevents, clinicalitemtypes, costrecords, demographics, diagnosiscodes, hospitaladmissions, icuepisodes, intakerecords, labresults, labtesttypes, medicationorders, microbiologyresults, outputrecords, patienttransfers, procedurecodes

Agent: Okay, the tables admissiondiagnoses and medication orders seem relevant. I need to explore these tables to see how diagnoses and medication orders are recorded.

Agent: column_search(arguements={"table_names": "admissiondiagnoses, medicationorders"})

Tool: CREATE TABLE admissiondiagnoses (
recordid INTEGER NOT NULL,
patientid INTEGER NOT NULL,
admissionid INTEGER NOT NULL,
icdcode TEXT NOT NULL,
codeversion INTEGER NOT NULL,
recordeddatetime TEXT NOT NULL,

```
2322
             PRIMARY KEY (recordid),
             FOREIGN KEY (codeversion) REFERENCES diagnosiscodes (codeversion),
2324
2325
             FOREIGN KEY (icdcode) REFERENCES diagnosiscodes (icdcode),
2326
             FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
2327
           )
2328
           /*
2330
           3 rows from admissiondiagnoses table:
2331
           recordid patientid admissionid icdcode codeversion recordeddatetime
2332
           0 10035185 22580999 4139 9 2100-05-17 12:53:00
2333
           1 10035185 22580999 V707 9 2100-05-17 12:53:00
2334
2335
           2 10035185 22580999 41401 9 2100-05-17 12:53:00
2336
           */
2337
2338
           CREATE TABLE medicationorders (
2339
2340
             recordid INTEGER NOT NULL,
             patientid INTEGER NOT NULL,
2342
             admissionid INTEGER NOT NULL,
2343
             startdatetime TEXT NOT NULL,
2344
2345
             enddatetime TEXT,
2346
             medicationname TEXT NOT NULL,
2347
             dosevalue TEXT NOT NULL,
2348
             doseunit TEXT NOT NULL,
2349
2350
             administrationroute TEXT NOT NULL,
2351
             PRIMARY KEY (recordid),
2352
             FOREIGN KEY (admissionid) REFERENCES hospitaladmissions (admissionid)
2353
           )
2354
2355
           /*
2356
           3 rows from medicationorders table:
           recordid patientid admissionid startdatetime enddatetime medicationname dosevalue
           doseunit administrationroute
2359
2360
           o 10020740 23831430 2100-04-19 11:00:00 2100-04-20 22:00:00 insulin o.o unit sc
2361
           1 10005817 20626031 2100-12-20 02:00:00 2100-12-21 12:00:00 insulin o.o unit sc
2362
           2 10005817 20626031 2100-12-27 10:00:00 2100-12-28 19:00:00 insulin 0.0 unit sc
2363
           */
2364
2365
2366
           Agent: Now I need to find the icdcode for epilepsy in the diagnosiscodes table. I will use
2367
           value_similarity_search to find similar values to "epilepsy" in the description column of the
2368
           diagnosiscodes table.
2369
2370
2371
           Agent: value_similarity_search(arguements={"column": "description", "table": "
2372
           diagnosiscodes", "value": "epilepsy"})
2373
2374
```

Tool: I found close matches in diagnosiscodes.description: ['epilepsy and recurrent seizures ', 'epilepsy, unspecified', 'other epilepsy', 'epileptic spasms', 'epilepsy, unspecified, not intractable', 'epilepsy, unspecified, without mention of intractable epilepsy', 'epilepsy, unspecified, with intractable epilepsy', 'epilepsy, unspecified, intractable', 'other epilepsy and recurrent seizures', 'epilepsy, unspecified, intractable, without status epilepticus', ' epilepsy, unspecified, not intractable, without status epilepticus', 'other epilepsy, intractable, without status epilepticus', 'epilepsy, unspecified, intractable, with status epilepticus', 'other epilepsy, intractable, with status epilepticus', 'epilepsy, unspecified, not intractable, with status epilepticus', 'generalized idiopathic epilepsy and epileptic syndromes', 'other epilepsy, not intractable, without status epilepticus', 'generalized idiopathic epilepsy and epileptic syndromes, intractable'].

Agent: The value similarity search returned several matches. To proceed, I will use 'epilepsy , unspecified' as the diagnosis for epilepsy. Now I need to check if patient 10008287 was diagnosed with 'epilepsy, unspecified' and prescribed 'carbamazepine'. If so, I need to find the most recent prescription time after the first epilepsy diagnosis.

Agent: value_substring_search(arguements={"column": "medicationname", "table": "medicationorders", "value": "carbamazepine"})

Tool: Values in medicationorders.medicationname containing 'carbamazepine': [' carbamazepine'].

Agent: Okay, the medication name is 'carbamazepine'. Now I have all the information I need to construct the SQL query.

Agent: sql_execute(arguements={"query": "SELECT MAX(m.startdatetime) FROM medicationorders m JOIN admissiondiagnoses a ON m.patientid = a.patientid WHERE m. patientid = 10008287 AND m.medicationname = 'carbamazepine' AND a.icdcode IN (SELECT icdcode FROM diagnosiscodes WHERE description = 'epilepsy, unspecified') AND m. startdatetime > (SELECT MIN(recordeddatetime) FROM admissiondiagnoses a2 WHERE a2. patientid = 10008287 AND a2.icdcode IN (SELECT icdcode FROM diagnosiscodes WHERE description = 'epilepsy, unspecified'))"})

Tool: [(None,)]

Agent: The query returned None, which means either the patient was not diagnosed with epilepsy, was not prescribed carbamazepine, or the prescription occurred before the diagnosis.

Agent: sql_execute(arguements={"query": "SELECT MIN(recordeddatetime) FROM admissiondiagnoses a WHERE a.patientid = 10008287 AND a.icdcode IN (SELECT icdcode FROM diagnosiscodes WHERE description = 'epilepsy, unspecified')")

Tool: [(None,)]

Agent: Since the query to find the first epilepsy diagnosis returned None, it means the patient was never diagnosed with epilepsy. Therefore, I cannot find the most recent carbamazepine prescription time after the first epilepsy diagnosis.

Agent: sql_execute(arguements={"query": "SELECT DISTINCT patientid FROM medicationorders WHERE medicationname = 'carbamazepine' AND patientid = 10008287")

Tool: []

Agent: The patient 10008287 was not prescribed carbamazepine. Therefore, I cannot find the most recent carbamazepine prescription time after the first epilepsy diagnosis. < answer>Patient 10008287 was not diagnosed with epilepsy and was not prescribed carbamazepine.</answer>

User: Okay, if they weren't prescribed carbamazepine, I'd like to know if they were prescribed any other medications for epilepsy. If so, I want the most recent prescription time after their first epilepsy diagnosis, just like I asked for carbamazepine.

Agent: Okay, I need to find out if patient 10008287 has been diagnosed with epilepsy. If they have, I need to find any medications prescribed to them that are used to treat epilepsy and then find the most recent prescription time for those medications after the first epilepsy diagnosis. First, I need to identify medications used to treat epilepsy. I will use a web search for this.

Agent: web_search(arguements={"query": "medications used to treat epilepsy"})

Tool:

"Below is a list of some of the most common brand-name drugs currently used to treat epilepsy. Your doctor may prefer that you take the brand name of anticonvulsant and not the generic substitution. Talk with your doctor about this important issue.\n\
nBrivaracetam (Briviact)\n\nCannabidiol (Epidiolex)\n\nCarbamazepine (Carbatrol or Tegretol)\n\nCenobamate(Xcopri)\n\nDiazepam (Valium), lorazepam (Ativan) and similar benzodiazepine tranquilizers such as clonazepam (Klonopin)\n\nEslicarbazepine (Aptiom) [...] Ethosuximide (Zarontin)\n\nFelbamate (Felbatol)\n\nFenfluramine (Fintepla)\n\nLacosamide (VIMPAT)\n\nLamotrigine (Lamictal)\n\nLevetiracetam (Keppra)\n\nOxcarbazepine (Oxtellar XR, Trileptal)\n\nPerampanel (Fycompa)\n\nPhenobarbitol\n\

2484 nPhenytoin (Dilantin)\n\nPregabalin (Lyrica)\n\nTiagabine (Gabitril)\n\nTopiramate (2485 Topamax)\n\nValproate, valproic acid (Depakene, Depakote)\n\nZonisamide (Zonegran) 2486 2487 [...] The most commonly used medications are benzodiazepines because they get into the 2488 bloodstream quickly to start working on the brain to stop the seizure. They include:\n\n 2489 ## Epilepsy Drug Guidelines\n\nIt may take several months before the best drug and 2490 dosage are determined for you. During this adjustment period, you will be carefully 2491 2492 monitored through frequent blood tests to measure your response to the medication.", 2493 "Below is a list of some of the most common brand-name drugs currently used to treat 2494 epilepsy. Your doctor may prefer that you take the brand name of anticonvulsant and not 2495 the generic substitution. Talk with your doctor about this important issue.\n\ 2496 2497 nBrivaracetam (Briviact)\n\nCannabidiol (Epidiolex)\n\nCarbamazepine (Carbatrol or 2498 Tegretol)\n\nCenobamate(Xcopri)\n\nDiazepam (Valium), lorazepam (Ativan) and similar 2499 benzodiazepine tranquilizers such as clonazepam (Klonopin)\n\nEslicarbazepine (Aptiom 2500) [...] Ethosuximide (Zarontin) $\n\$ nFelbamate (Felbatol) $\n\$ nFenfluramine (Fintepla) $\n\$ 2501 nLacosamide (VIMPAT)\n\nLamotrigine (Lamictal)\n\nLevetiracetam (Keppra)\n\ 2503 nOxcarbazepine (Oxtellar XR, Trileptal)\n\nPerampanel (Fycompa)\n\nPhenobarbitol\n\ 2504 nPhenytoin (Dilantin)\n\nPregabalin (Lyrica)\n\nTiagabine (Gabitril)\n\nTopiramate (2505 Topamax)\n\nValproate, valproic acid (Depakene, Depakote)\n\nZonisamide (Zonegran) 2506 2507 [...] The most commonly used medications are benzodiazepines because they get into the bloodstream quickly to start working on the brain to stop the seizure. They include:\n\n 2509 ## Epilepsy Drug Guidelines\n\nIt may take several months before the best drug and 2510 dosage are determined for you. During this adjustment period, you will be carefully 2511 2512 monitored through frequent blood tests to measure your response to the medication.", 2513 "Medications used to treat epilepsy are called antiepileptic drugs (AEDs). There are more 2514 than 30 prescription AEDs on the market, and they're mostly available as oral tablets or 2515 capsules.\n\nYour options depend on:\n\n## Types of antiepileptic drugs (AEDs)\n\ 2516 2517 nThere are two types of medications for epilepsy:\n\nSome people may need to take 2518 more than one medication to prevent seizures. [...] ### Levetiracetam (Elepsia XR, Keppra, 2519 Keppra XR, Spritam)\n\nLevetiracetam (Elepsia XR, Keppra, Keppra XR, Spritam) may treat a range of epileptic seizures, including:\n\nIt's available as a pill, an IV solution, an 2521 2522 oral solution, and an injection.\n\nLevetiracetam may cause fewer side effects than other medications used for epilepsy. It's safe to take during pregnancy, according to 2524 experts such as the United Kingdom's Commission on Human Medicines.\n\n### 2525 Lorazepam (Ativan) [...] Lorazepam (Ativan) is a benzodiazepine that's used to treat all 2526 2527 types of seizures. It's also used to treat status epilepticus. Status epilepcticus is a 2528 prolonged, critical seizure that's regarded as a medical emergency.\n\nIt's available as a 2529 pill, an oral concentrate, and an injection.\n\n### Methsuximide (Celontin)\n\ 2530 nMethsuximide (Celontin) is used for absence seizures. It's prescribed when other 2531 2532 treatments don't work in treating your seizures.", 2533 "Donate\n\nPopular searches: Diagnosing Epilepsy Treatments and Therapies what is 2534 epilepsy\n\nMake an Impact\n\nEpilepsy and Seizure Medications\n 2535 ======\n\n#### Learn about FDA-approved medications to 2536

2538 treat epilepsy and seizures.\n\nImage 7: Epilepsy and Seizure Medications\n\nSearch\n\ 2540 nFilter Alphabetically Select Clear\n\nLoading ...\n\n##### Acetazolamide\n\nMore info(2541 Brand names: Diamox, Diamox Sequels, generics) [...] Acetazolamide (a SEET a ZOLE a 2542 mide) has been FDA-approved for the treatment of, along with other drugs, 2543 centrencephalic epilepsies (absence, generalized seizures).\n\n#### Brivaracetam\n\ 2544 nMore info(Brand names: Briviact)\n\nBrivaracetam (briv a RA se tam) has been 2545 2546 approved by the FDA to treat focal (partial) onset seizures in patients 1 month of age and 2547 older.\n\n#### Cannibadiol\n\nMore info(Brand names: Epidiolex) [...] More info(Brand 2548 names: Aptiom)\n\nEslicarbazepine acetate (ES li Kar BAZ e peen) has been approved by 2549 the FDA to treat focal onset seizures in patients 4 years of age and older.\n\n##### 2550 2551 Ethosuximide\n\nMore info(Brand names: Zarontin, generics)\n\nEthosuximide (ETH oh 2552 SUX i mide) has been approved by the FDA to control absence (petit mal) epilepsy.\n\n 2553 ##### Ethotoin\n\nMore info(Brand names: Peganone, generics)", 2554 "Healthcare providers prescribe antiseizure medications to treat epilepsy and 2555 symptomatic seizures. They also prescribe these medications to prevent and/or treat 2557 seizures that happen during or following brain surgery. [...] Gray gradient\nGray gradient\ nSearch Icon Blue\nCleveland Clinic logo\n\n# Antiseizure Medications (Formerly Known 2559 as Anticonvulsants)\n\nAntiseizure medications (anticonvulsants) help treat epilepsy 2560 2561 and other causes of seizures. They can treat other conditions as well, like anxiety and neuropathic pain. There are several different types of antiseizure medications. You and 2563 your healthcare provider will work together to find the best one for you.\n\ 2564 nAdvertisement [...] Cleveland Clinic is a non-profit academic medical center. Advertising 2565 2566 on our site helps support our mission. We do not endorse non-Cleveland Clinic products 2567 or services. Policy\n\n## Overview\n\n### What are antiseizure medications (2568 anticonvulsants)?\n\nAntiseizure medications (previously known as antiepileptic or 2569 anticonvulsant medications) are prescription medications that help treat and prevent 2570 2571 seizures. Healthcare providers may prescribe these medications to treat other conditions 2572 as well.", "Epilepsy Website Logo\n\n# List of Anti-Seizure Medications (ASMs)\n\n#### Understanding Epilepsy\n\nAnti-epileptic drugs (ASMs) are the main form of treatment 2575 2576 for people living with epilepsy, with up to 70% (7 in 10 people) having their seizures 2577 controlled through this medication.\n\nIn Australia there are over 20 ASMsare used to 2578 treat seizures. The ASMsprescribed are often selected on the basis of the seizure type/s. 2579 age, gender and side effects. ASMsmay be prescribed as tablets, syrups and liquids.", 2580 2581 "One large randomized trial, the Standard and New Antiepileptic Drugs (SANAD) trial, 2582 demonstrated some comparative advantages of certain AEDs when treating focal or 2583 generalized epilepsy. In the end, when comparing valproate, lamotrigine, or topiramate 2584 for generalized seizures, they recommended valproic acid as their first-line choice. 2585 2586 Additionally, when comparing carbamazepine, gabapentin, lamotrigine, oxcarbazepine, 2587 and topiramate for focal seizures, lamotrigine was cited as the first-line [...] In summary, 2588

it is now abundantly clear that anti-seizure medications wield disparate mechanistic

profiles, but they all effectively suppress epileptic seizures in one way or another.

2589

Accordingly, grouping the drugs together by mechanism is a very helpful organizing principle. From this viewpoint, it may become easier to appreciate that some drugs have different efficacy profiles for different seizures types and epilepsy syndromes. Ethosuximide is an exception with its specific limited use with [...] Phenytoin is one of the oldest anti-seizure medications and is still widely used for focal and generalized seizures. It is also administered for status epilepticus. In addition, practitioners may invoke phenytoin as a second-line agent for patients with mixed seizure types (e.g., tonic -clonic and myoclonic). As mentioned, phenytoin blocks voltage-gated sodium channels, but other possible mechanisms revolve around decreased synaptic transmission, smaller changes in ionic gradients involving the",

"Medicines.\n Surgery.\n Therapies that stimulate the brain using a device.\n A

"Medicines.\n Surgery.\n Therapies that stimulate the brain using a device.\n A ketogenic diet.\n\n### Medication\n\nMost people with epilepsy can become seizure-free by taking one anti-seizure medicine, which is also called an anti-epileptic medicine. Others may be able to decrease the number and intensity of their seizures by taking more than one medicine. [...] Tell your healthcare professional immediately if you notice new or increased feelings of depression or suicidal thoughts. Also contact your healthcare professional right away if you have changes in your mood or behaviors.\n Tell your healthcare professional if you have migraines. You may need an anti-seizure medicine that can prevent your migraines and treat epilepsy. [...] ### Potential future treatments\n\nResearchers are studying many potential new treatments for epilepsy, including:\n\n Continuous stimulation of the seizure onset zone, known as subthreshold stimulation. Subthreshold stimulation is continuous stimulation to an area of the brain below a level that's physically noticeable. This type of therapy appears to improve seizure outcomes and quality of life for some people with seizures. Subthreshold stimulation helps stop a seizure before it happens.",

"| stiripentol | Diacomite | Image 50: Medline unavailable | Image 51: Medline unavailable |\n| sulthiame | Ospolot | Image 52: Medline available | Image 53: Medline unavailable |\n | tiagabine | Gabitril | Image 54: Medline available | Image 55: Medline available |\n| topiramate | Topamax | Image 56: Medline available | Image 57: Medline available |\n| vigabatrin | Sabril | Image 58: Medline available | Image 59: Medline available | [...] | ethosuximide | Zarontin | Image 20: Medline available | Image 21: Medline available |\n| felbamate | Felbatol | Image 22: Mims unavailable | Image 23: Medline available |\n| gabapentin | Neurontin | Image 24: Medline available | Image 25: Medline available |\n| lacosamide | Vimpat | Image 26: Medline available | Image 27: Medline available |\n| lamotrigine | Lamictal | Image 28: Medline available | Image 29: Medline available | [...] | phenobarbitone | Phenobarbitone | Image 40: Medline available | Image 41: Medline unavailable |\n| phenytoin | Dilantin | Image 42: Medline available | Image 43: Medline available |\n| potassium bromide | Dibro-Be mono | Image 44: Medline unavailable | Image 45: Medline unavailable |\n| primidone | Mysoline | Image 46: Medline available | Image 47: Medline available |\n| sodium valproate | Epilim | Image 48: Medline available | Image 49: Medline available |",

2646 "Sodium valproate is an epilepsy drug prescribed for all seizure types including absence, 2647 myoclonus and tonic clonic seizures. New regulations have been introduced by the 2648 2649 Medicines and Healthcare products Regulatory Agency (MHRA) around the way in which 2650 the drug, sodium valproate is prescribed to women and girls of childbearing age.\n\nLET' 2651 S STAY IN TOUCH... [...] The decision whether to start taking anti-epileptic drugs (AEDs) 2652 can be difficult, and there is a lot to think about. Here we look at the benefits and risks of 2653 2654 taking, or not taking, AEDs.\n\nSodium valproate [...] Epilepsy Society\n\nSite 2655 Navigation\n\nUtility Nav\n\nSite Search\n\nbreadcrumb navigation:\n\nList of anti-2656 seizure medication\n\nPublished on 8 October 2024\n\nUpdated: 7 April 2025\n\ 2657 nAuthored by Anonymous\n\nAlso in this section\n\nList of anti-seizure medication\n\ 2658 2659 nHere is a list of anti-seizure medication (ASM), previously called anti-epileptic drugs or 2660 AEDs." 2661] 2662 2663 2664 Agent: Based on the web search results, here's a list of medications commonly used to 2665 treat epilepsy: 2666 Brivaracetam 2667 Cannabidiol 2668 2669 Carbamazepine 2670 Cenobamate 2671 Diazepam 2672 Lorazepam 2673 2674 Clonazepam 2675 Eslicarbazepine 2676 Ethosuximide 2677 **Felbamate** 2678 2679 Fenfluramine 2680 Lacosamide 2681 Lamotrigine 2682 Levetiracetam 2683 2684 Oxcarbazepine 2685 Perampanel 2686 Phenobarbitol 2687 Phenytoin 2688 2689 Pregabalin 2690 **Tiagabine** 2691 **Topiramate** 2692 Valproate, valproic acid 2693 2694 Zonisamide 2695 Acetazolamide 2696 Methsuximide 2697

Stiripentol

2701 Sulthiame

2702 Gabapentin

Primidone

Sodium valproate

I will now construct a query to find the most recent prescription time for any of these medications after the first epilepsy diagnosis.

Agent: sql_execute("query": "SELECT MAX(m.startdatetime) FROM medicationorders m JOIN admissiondiagnoses a ON m.patientid = a.patientid WHERE m.patientid = 10008287 AND m. medicationname IN ('brivaracetam', 'cannabidiol', 'carbamazepine', 'cenobamate', 'diazepam', 'lorazepam', 'clonazepam', 'eslicarbazepine', 'ethosuximide', 'felbamate', 'fenfluramine', 'lacosamide', 'lamotrigine', 'levetiracetam', 'oxcarbazepine', 'perampanel', 'phenobarbitol', 'phenytoin', 'pregabalin', 'tiagabine', 'topiramate', 'valproate', 'valproic acid', 'zonisamide', 'acetazolamide', 'methsuximide', 'stiripentol', 'sulthiame', 'gabapentin', 'primidone', 'sodium valproate') AND a.icdcode IN (SELECT icdcode FROM diagnosiscodes WHERE description = 'epilepsy, unspecified') AND m.startdatetime > (SELECT MIN(recordeddatetime) FROM admissiondiagnoses a2 WHERE a2.patientid = 10008287 AND a2. icdcode IN (SELECT icdcode FROM diagnosiscodes WHERE description = 'epilepsy, unspecified'))")

Tool: [(None,)]

Agent: The query returned None, which means either the patient was not diagnosed with epilepsy, was not prescribed any of the listed medications, or the prescriptions occurred before the diagnosis. Since I already determined that the patient was not diagnosed with epilepsy, the result is None. <answer>Patient 10008287 was not diagnosed with epilepsy and was not prescribed any anti-epileptic medications.</answer>

User: Could you double check that your answer covers everything? I wanted to know if patient 10008287 was prescribed carbamazepine and, if not, if they were prescribed any other epilepsy meds. For any epilepsy meds, I need the most recent prescription time after their first epilepsy diagnosis.

Agent: Yes, I have confirmed that patient 10008287 was neither diagnosed with epilepsy nor prescribed carbamazepine or any other anti-epileptic medications from the list I generated using a web search. Therefore, I cannot provide a prescription time after their first epilepsy diagnosis. <answer>Patient 10008287 was not diagnosed with epilepsy and was not prescribed carbamazepine or any other anti-epileptic medications.</answer>

User: ###END###