

EHR-SeqSQL : A Sequential Text-to-SQL Dataset For Interactively Exploring Electronic Health Records

Anonymous ACL submission

Abstract


Text-to-SQL parsing allows non-expert users to obtain information from databases using natural language. There are several important yet under-explored objectives in this field: interactivity, compositionality, and efficiency. In this paper, we present EHR-SeqSQL, a sequential text-to-SQL dataset, specifically designed for Electronic Health Record (EHR) databases. We demonstrate the benefits of multi-turn setting over single-turn setting with respect to compositionality, and provide a new data split and an additional test set to evaluate compositional generalization. Furthermore, we introduce unique special tokens in SQL queries to enhance execution efficiency. By addressing all these objectives above, our research aims to bridge the gap between industry needs and academic research in the text-to-SQL domain.

1 Introduction


Text-to-SQL provides a practical opportunity for non-experts to explore databases, even without prior knowledge of the database operations. Electronic Health Records (EHRs) are large-scale relational databases (RDBs) storing vast and comprehensive patient data (Johnson et al., 2016; Pollard et al., 2018). Medical experts often ask questions that require highly complex reasoning across multiple tables and access to a vast number of records within a single query (Lee et al., 2022). Handling such complexity in large-scale databases remains a significant challenge in the current text-to-SQL research (Li et al., 2023).

Several efforts have been made to construct text-to-SQL datasets for EHRs. MIMIC-SQL (Wang et al., 2020) is the first text-to-SQL dataset that targets a subset of MIMIC-III (Johnson et al., 2016), one of the widely-used open-source EHR databases, consisting of 26 tables. DrugEHRQA (Bardhan et al., 2022) provides the first question answering dataset that incorporates both structured tables and


EHRSQL

 Q. What is the name of the **drug** that has been prescribed **two times** in the **current hospital encounter** to **patient 14467**?

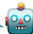
```
SQL SELECT t1.drug
FROM ( SELECT prescriptions.drug, COUNT( prescriptions.startdate ) AS c1
      FROM prescriptions
      WHERE prescriptions.hadm_id IN (
        SELECT admissions.hadm_id
        FROM admissions
        WHERE admissions.subject_id = 14467
        AND admissions.dischtime IS NULL )
      GROUP BY prescriptions.drug) AS t1
WHERE t1.c1 = 2
```


A : Potassium chloride 

EHR-SeqSQL (Ours)


 Q1. Give me the hospital stay id for **patient 14467** on the **current hospital visit**.


```
SQL† SELECT admissions.hadm_id FROM admissions
WHERE admissions.subject_id = 14467 AND admissions.dischtime IS NULL
```

A1 : 213008 


 Q2. What are the **drugs** linked to **A1**?


```
SQL† SELECT prescriptions.drug FROM prescriptions
WHERE prescriptions.hadm_id IN ( PREV_RESULT1 )
```

A2 : Potassium chloride, Zolpidem tartrate 

 Q3. How many times has each drug been prescribed in **A2**?

```
SQL† SELECT prescriptions.drug, COUNT( prescriptions.startdate ) AS c1
FROM prescriptions WHERE prescriptions.hadm_id IN ( PREV_RESULT1 )
GROUP BY prescriptions.drug ) AS t1
```

A3 : Potassium chloride: 2, Zolpidem tartrate: 1 

 Q4. What is the drug that appears **two times** in **A3**?

```
SQL† SELECT t1.drug FROM ( PREV_QUERY3 ) AS t1 WHERE t1.c1 = 2
```

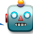
A4 : Potassium chloride 

Figure 1: **EHRSQL vs. EHR-SeqSQL** EHR-SeqSQL is a dataset that adapts the single-turn setting of EHRSQL into a multi-turn setting. The SQL queries in EHR-SeqSQL include the special tokens to refer to the previous context, which can be executed in the database with simple post-processing.

unstructured notes from EHRs. EHRSQL (Lee et al., 2022) is a dataset curated based on a survey from various medical experts, reflecting the diverse information needs of the actual medical field.

Still, there are important yet under-explored objectives for the practical application of text-to-SQL

041
042
043
044
045
046

models. **First, the text-to-SQL system should incorporate interactivity.** Most of the existing text-to-SQL research assumes a single-turn scenario, whereas the process of exploring information in real-world situations is often continuous (Iyyer et al., 2017; Yu et al., 2019a,b). For better user usability, it is desirable for information retrieval solutions to be interactive as well. **Second, the text-to-SQL system should learn compositionality.** Realistically, datasets cannot capture all the varied needs of users. Therefore, it is crucial that the model can handle a wider range of queries especially when these unseen queries comprise components (*i.e.* sub-queries) that the model has explicitly seen during training. This is the compositional generalization ability, which is quite challenging even for the highly proficient language models (Qiu et al., 2022). **Lastly, the text-to-SQL system should embrace efficiency.** Real-world databases are often significantly larger than academic ones (Hazoom et al., 2021). Efficient query execution is crucial, given the magnitude of the databases in actual hospitals (*e.g.* the original MIMIC-III dataset includes the medical history of 46k patients). This becomes even more significant when we consider a real-time interaction scenario between a text-to-SQL model and a user (Li et al., 2023).

In this work, we introduce EHR-SeqSQL which addresses all three objectives. Our contribution is as follows:

- **Construction of a sequential text-to-SQL dataset for exploring EHR data:** As illustrated in Figure 1, EHR-SeqSQL is designed for multi-step interactions, built by decomposing the diverse and complex queries of EHRSQL and setting each question as the interaction goal. To the best of our knowledge, EHR-SeqSQL is the first multi-turn text-to-SQL dataset targeting structured medical records. We make our dataset public to encourage future research.
- **Validating the effectiveness of EHR-SeqSQL in the compositional generalization:** We designed two experiments to evaluate the compositional generalization ability. Through our experiments, it was found that decomposing questions enables the model to better generalize to unseen interaction goals during training. Furthermore, it demonstrated the potential for generalization on longer sequences not encountered during training.

- **Proposing special tokens for SQL:** We introduce a highly effective way to enhance query execution efficiency in multi-step interactions with the use of the novel special tokens for SQL. We observed that these tokens significantly increase time efficiency during query execution as well as improve the performance of text-to-SQL models. Moreover, the efficiency grows significantly as the database size increases.

2 Related Work

2.1 Multi-turn Text-to-SQL

There are only a few datasets that are specifically designed for context-dependent text-to-SQL tasks. ATIS (Hemphill et al., 1990) was the first to include a series of user questions aimed at interacting with a flight database. Recently, Yu et al. (2019a,b) proposed cross-domain context-dependent text-to-SQL datasets, based on the questions from Spider (Yu et al., 2018) serving as interaction goals. While their motivation aligns closely with ours, their questions and the grounding databases are predominantly simplistic. In real-world scenarios, however, users typically have much more complex requirements, and databases often contain a significantly larger number of rows (Lee et al., 2022; Li et al., 2023). EHR-SeqSQL is designed upon this objective, based on the questions collected from an actual survey and grounding the real-world database, MIMIC-III (Johnson et al., 2016). To the best of our knowledge, this is the first multi-turn text-to-SQL dataset in the healthcare domain.

2.2 Compositional Generalization

Compositional generalization, a major challenge in semantic parsing, enables a model to manage diverse and unfamiliar user queries using components recognized from training. It is typically achieved by creating different training and test splits of interest. Lake and Baroni (2018) discovered that sequence-to-sequence models struggle to learn compositional structures, as demonstrated through their primitive and length splits. More recently, template splits (Finegan-Dollak et al., 2018) and TMCD splits (Shaw et al., 2020) were proposed to evaluate compositional generalization on text-to-SQL datasets. In this work, we further extend the concept of compositional generalization into an interactive multi-turn setting within the text-to-SQL domain.

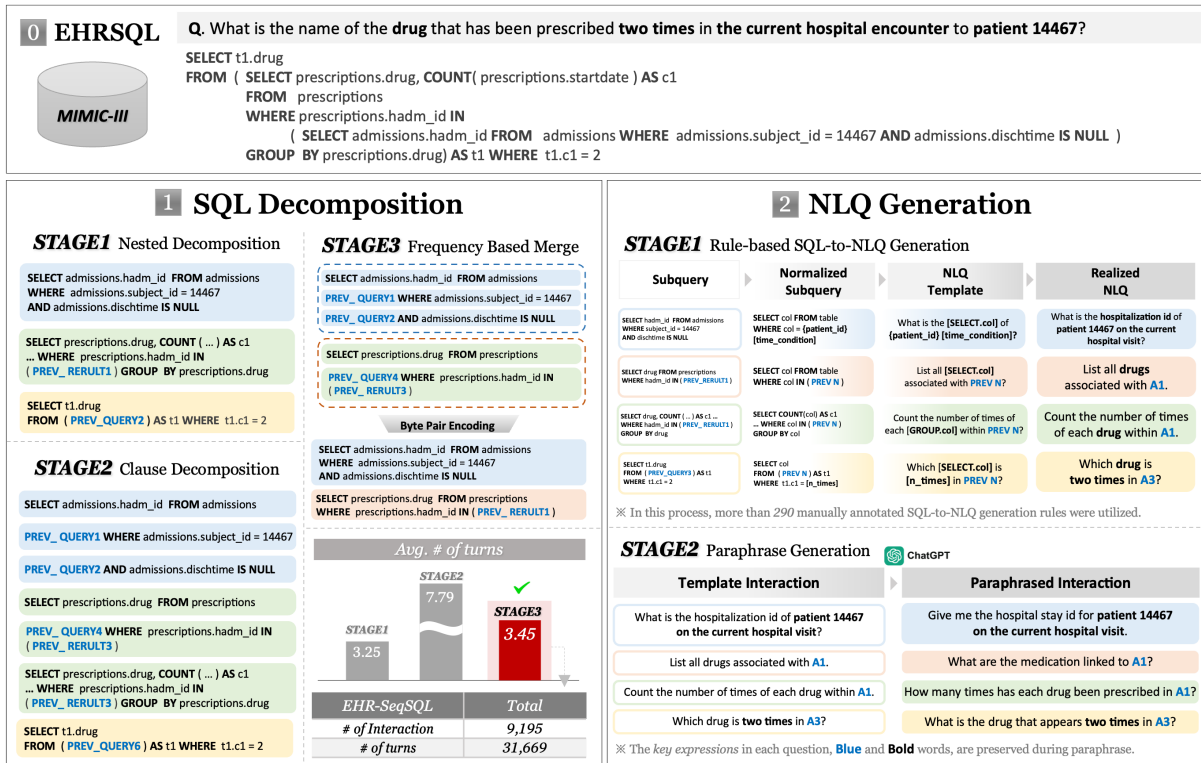


Figure 2: Overview of the dataset construction process. We transform a single text-SQL pair from EHRSQL into multi-turn text-SQL pairs for EHR-SeqSQL. We first decompose the original SQL query into subqueries in Stage 1 and 2, then apply the BPE algorithm to merge frequent subquery patterns in Stage 3. We then generate natural language questions (NLQ) for each subquery based on manually created NLQ templates. Lastly, NLQs are paraphrased using ChatGPT to create the final dataset.

2.3 SQL Efficiency

In domains that utilize large and complex databases such as EHRs, the efficiency of SQL becomes increasingly critical. The majority of existing research focuses on improving SQL efficiency through post-hoc query optimization (Zhou et al., 2021; Li et al., 2023) or efficient database indexing (Zhou, 2022). In contrast to these approaches, we propose special tokens for SQL that specifically target SQL efficiency in a multi-turn setting. Our approach not only improves execution efficiency but also enhances performance, which is further discussed in Section 5.3.

3 Data Construction

We aim to convert text-SQL pairs of EHRSQL to multi-turn text-SQL pairs that embrace interactivity and query execution efficiency. Note that EHRSQL contains text-SQL samples for two different EHR sources, namely MIMIC-III (Johnson et al., 2016) and eICU (Pollard et al., 2018). We use the MIMIC-III version of EHRSQL, given its complex schema and wider adoption by the NLP community.

Our data construction process consists of two parts: SQL decomposition and natural language question (NLQ) generation. In SQL decomposition, we first break down each SQL query of EHRSQL into a sequence of subqueries to create more granular meanings of queries. Next, we create a corresponding NLQ for each subquery. Compared to NLQ decomposition, SQL decomposition enables systematic breakdown of a question without any loss of information, due to the strict, rule-based nature of SQL syntax. The overall process is shown in Figure 2 and we elaborate each step in the following sections.

3.1 SQL Decomposition

Stage 1. Decomposing Nested Query As SQL queries in EHRSQL follow a nested structure, we decompose the EHRSQL queries using their nested structure. In this way, subqueries can contain more granular meanings of the original query, which in turn allows us to create multiple turns of interaction. The inner query is asked first, and then a question corresponding to the outer query is asked while referring to the result of the inner query. This

	ATIS	SParC	CoSQL	EHR-SeqSQL
# of Interactions	1,658	4,298	3,007	9,195
# of Turns	11,653	12,726	15,598	31,669
Avg. # of Turns	7.0	3.0	5.2	3.5
Avg. # of Tables / DB	27	5.1	5.1	26
Domain	Airline	Cross	Cross	Medical
Compositional Split	✗	✗	✗	✓
SQL Execution Efficiency	✗	✗	✗	✓

Table 1: Comparison of EHR-SeqSQL with other multi-turn text-to-SQL datasets.

decomposition approach guarantees the executability of queries at each turn and allows anaphoric expressions in the subsequent turns.

Meanwhile, we devised two special tokens for SQL to leverage the previous turn’s subqueries and results: `prev_query` and `prev_result`. These tokens are accompanied by a specific turn index (e.g. `prev_query1`, `prev_result2`). Specifically, `prev_query` token refers to the generated query of the specified turn, whereas the `prev_result` token refers to the execution result of the specified turn. Two tokens are substituted with either the referred query or its execution result before execution. The effectiveness of these tokens is later discussed in Section 5.3.

Stage 2. Decomposing SQL Clauses After Stage 1, multi-turn questions are obtained, where each can be asked with reference to the answer to previous questions. Still, these questions often convey multiple conditions that users might ask in separate questions. For example, the question “*What was the last hospitalized admission time that patient 17694 was admitted via transfer from emergency room?*” contains the specific conditions about the admission time, the patient, and the admission route. Therefore, we further decompose queries on the clause-level: WHERE, ORDER BY, and HAVING clauses and aggregation functions (e.g. MAX, MIN, SUM) in the SELECT clause. In cases where multiple clauses appear together, we parse in the order of SQL execution. Exceptions were made only when the original semantics of the query became ambiguous after being decomposed. Please see further details in Appendix A.1.

Stage 3. Merging Subqueries by Frequency

We noted that frequently appearing consecutive subqueries are most likely decomposed due to the specific database schema characteristics or SQL structures rather than each being queries users

would naturally ask. In Stage 3, we employ a recursive application of the Byte Pair Encoding (BPE) algorithm to merge frequent consecutive subquery pairs. This recursive approach allows for the amalgamation of not only two but up to three or more subqueries, thereby simplifying complex patterns and enhancing query interpretability. Further details are given in Appendix A.2.

3.2 NLQ Generation

Stage 1. Rule-based SQL-to-NLQ Generation

While previous SQL-to-NLQ studies (Gan et al., 2022; Wu et al., 2021) decomposed SQL at the clause level for sub-questions and concatenated clauses later, such simple concatenation has a risk of unnaturalness. Thus we create NLQ templates for each corresponding SQL subqueries for the quality of NLQ. To efficiently annotate NLQs for each corresponding SQL subqueries, we first normalize subqueries by replacing specific table names, column names, and condition values to abstract terms such as `table` and `col` (see Figure 2). We then manually create NLQ templates for each normalized subqueries. These templates include slots for table names, column names, condition values, SQL functions (e.g. GROUP BY, AVG), and time expressions. Then the actual question is generated from the NLQ templates by the slot-filling process. More details are in Appendix A.3.

Stage 2. Paraphrase Generation

We further paraphrase the NLQ templates to enhance linguistic variability within the dataset, using ChatGPT for its superior performance in understanding and generating text (Guo et al., 2023). We ensured that all the slot values and turn indices were preserved after paraphrasing. Additionally, we employed the self-consistency method and a duplicate question detection model to ensure the quality of the paraphrases, following Lee et al. (2022). On average, we obtained 10.39 paraphrases for each NLQ tem-

268 plate. Finally, we randomly assign the paraphrased
269 template to each question and fill the slots with the
270 original condition values.

271 **Quality Check** In the process of decomposing
272 SQL queries, we strictly adhered to traditional SQL
273 syntax to ensure the intactness. However, generat-
274 ing NLQs based on templates necessitate meticu-
275 lous quality checks. On the turn level, we examined
276 the consistency of NLQ templates of each question.
277 Our primary concerns were 1) *Completeness*: We
278 make sure that all information in the SQL is also
279 explicitly stated in the NLQ, and 2) *Naturalness*:
280 We ensured the question templates were as natural
281 as possible after the masked information is realized.
282 On the interaction level, we evaluated whether the
283 interaction accurately represents the original intent
284 of the EHRSQL question. We also carefully an-
285 alyzed how turns are connected naturally within
286 an interaction through coreference. We performed
287 the turn-level check for every question template
288 and both turn-level and interaction-level quality
289 checks for 1,000 sampled interactions from the fi-
290 nal data. When unnatural template or interaction
291 flow is found, we created additional templates or
292 edited the template to be clear until all the sam-
293 pled instances conformed to our quality standards
294 in both turn and interaction level.

295 3.3 Data Statistics

296 Table 1 presents the statistics of EHR-SeqSQL
297 compared to other multi-turn text-to-SQL datasets.
298 EHR-SeqSQL has the largest number of interac-
299 tions and turns compared to all existing multi-turn
300 text-to-SQL datasets. Notably, this is the first work
301 to introduce special tokens within SQL queries to
302 improve execution efficiency. Furthermore, we also
303 provide a new split and additional test set to evalu-
304 ate the compositional generalization, as detailed in
305 Section 5.

306 Questions in our dataset are categorized into four
307 types: *independent*, *referential*, *filtering*, and *modi-*
308 *fying*. A question can belong to multiple categories
309 as they are not mutually exclusive. *Independent*
310 questions lack prior context. *Referential* questions
311 refer to previous questions or answers while *filter-*
312 *ing* questions narrow the previous question’s scope
313 by adding conditions. *Modifying* questions are par-
314 ticularly challenging as they only mention the al-
315 tered condition and omit all the same conditions.
316 Table 6 shows the distribution of the questions.

4 Experimental Setup 317

4.1 Task 318

319 The objective of the model is to generate a SQL
320 query given the interaction history and current ques-
321 tion. We experiment with two versions of interac-
322 tion history: one that includes only the previous
323 questions, denoted as QQ , and another that includes
324 both the previous questions and their corresponding
325 SQL queries, denoted as QS . For the QS setting,
326 the model is trained with the interaction history
327 using the ground-truth queries, while during the
328 inference phase, the model’s own predictions are
329 used in order to simulate a real-world application.

4.2 Baselines 330

4.2.1 Fine-tuning Models 331

332 We employ both a fine-tuning approach and an in-
333 context learning approach for our baselines. For
334 the fine-tuning approach, we use the T5 models
335 (Raffel et al., 2020), the general-purpose sequence-
336 to-sequence models as our baseline models. We did
337 not use the state-of-the-art models for SPaC (Yu
338 et al., 2019b) or CoSQL (Yu et al., 2019a) due to
339 their SQL grammar being confined to Spider which
340 is not compatible with our dataset.

4.2.2 In-context Learning Models 341

342 We use ChatGPT for our in-context learning mod-
343 els. Instead of a zero-shot approach, we employ
344 few-shot prompting. We retrieve similar exem-
345 plars for each test instance using the BM25 algo-
346 rithm to use as prompt. For every experiment with
347 in-context learning setting, we use 20-shot. For
348 EHRSQL, which is single-turn, computing the sim-
349 ilarity for each question is sufficient. However,
350 for EHR-SeqSQL, which is multi-turn and context-
351 dependent, it’s necessary to consider both the cur-
352 rent question and the interaction history. Therefore,
353 we use 10 retrieved examples related to the cur-
354 rent question and another 10 related to the entire
355 interaction history. More details such as prompt
356 and few-shot retrieval in our in-context learning
357 baseline are in Appendix E.

4.3 Evaluation 358

359 The two commonly used metrics in text-to-SQL
360 are Exact Match Accuracy (EM) and Execution
361 Accuracy (EX). However, the EM metric can some-
362 times be overly strict because it doesn’t take into
363 account predictions that have different SQL syntax

but yield the same execution result as the ground-truth query. Thus we use the EX score to measure the model performance based on query execution results. Before execution, the special tokens are replaced with the generated queries or their execution result through post-processing, thus any errors from the referenced turn will be propagated.

Additionally, for the multi-turn setting, we adopt Interaction Match (IM) and Question Match (QM) following Yu et al. (2019b) while computing EX. IM measures the accuracy of the entire interaction while QM measures the accuracy of each turn. We utilize appropriate metrics regarding the purpose of each experiment.

5 Experiments

We now empirically demonstrate the benefits of EHR-SeqSQL in terms of two types of compositional generalization (§ 5.1, § 5.2), as well as the effectiveness of the special tokens (§ 5.3).

5.1 Generalization to Unseen Interaction Goals

Compositional Split We first aim to evaluate whether training models in a multi-turn setting can lead to the acquisition of compositional generalization abilities. To explore this aspect, we split our dataset in a compositional manner, namely *compositional split*. This split differs from the *random split* in EHRSQL where the distributions of SQL structures in the training and test sets are nearly identical. In contrast, *compositional split* includes unseen SQL structures in the test set, though these structures can be decomposed into smaller parts that are all present in the training data.

More concretely, we first define the terms *compositions* and *components* and use the concepts to automatically split the dataset. *Compositions* represent SQL templates in EHRSQL, where condition values, SQL functions (i.e. SUM, AVG, etc), time expressions (i.e. subqueries constraining *last year*, *in 2023*, etc) are masked. *Components* refer to the decomposed SQL template clauses derived from Stage 2, as explained in Section 3.1. Each composition, which corresponds to an interaction goal, contains a set of components that exist as individual SQL subqueries in the dataset. Table 11 provides concrete examples of components and compositions. We employ a greedy algorithm to split the dataset, similar to Shaw et al. (2020). Table 7 provides the statistics of the random and compositional

splits. More details are given in Appendix B.

Metric We only measure IM to compare the performance of a model trained on EHRSQL and a model trained on EHR-SeqSQL. Given that an interaction in EHR-SeqSQL corresponds to a question in EHRSQL, correctly predicting every question within an interaction is equal to correctly predicting a single question in EHRSQL.

Result The experimental results are shown in Table 2. In a random split, all models exhibit strong performance with both EHRSQL and EHR-SeqSQL. However, in the compositional split, models typically show superior generalization ability when trained with EHR-SeqSQL. T5 models, in particular, show a significant performance increase, ranging from 7.63%p (from 67.8 to 75.43), to 30.34%p (from 52.15 to 82.49). T5-base finetuned in *QQ* setting has the best generalization ability. It suggests that a task-specific fine-tuning enables effective extraction of necessary information only from NLQs without being distracted by other contextual factors. On the other hand, ChatGPT suffers a performance drop in the *QQ* setting, which is suspected to be due to the absence of predicted SQL queries within the prompt, which leads ChatGPT to have less information on target representation. Still, a multi-turn setting leads to a performance increase of 7.63% in the *QS* setting. It’s worth noting that ChatGPT has better compositional generalization ability in a single-turn setting than fine-tuned models. Overall, our experiments show that training with EHR-SeqSQL allows the models to generalize well even with unseen questions.

5.2 Generalization to Longer Interactions

Longer Interaction Generation In this section, we aim to explore whether the models trained with EHR-SeqSQL can comprehend longer interactions, by achieving another type of compositional generalization ability. To test this, we created a new test set ($Test_L$) composed of multiple related interactions that simulate follow-up questions. This mirrors real-world situations where multiple related questions naturally arise as a user often seeks to familiarize themselves with the knowledge they are curious about (Yu et al., 2019b). For each test set of *random split* and *compositional split*, we created 100 longer interaction instances by connecting related interactions, denoting them as $Test_{LR}$ and $Test_{LC}$. The statistics for $Test_L$ can be found in Table 7. These interactions are, on average, five

Model	Random			Compositional		
	EHRSQL	EHR-SeqSQL		EHRSQL	EHR-SeqSQL	
		<i>QQ</i>	<i>QS</i>		<i>QQ</i>	<i>QS</i>
T5-Base	94.25 \pm 0.54	94.69 \pm 0.01	95.38 \pm 1.48	52.15 \pm 0.85	82.49 \pm 1.80	75.38 \pm 1.94
T5-3B	90.68 \pm 9.47	93.07 \pm 3.71	92.38 \pm 4.9	54.90 \pm 3.51	71.39 \pm 1.05	70.07 \pm 1.85
ChatGPT	91.22	80.74	91.37	67.80	60.50	75.43

Table 2: Model performances in two different splits. For the fine-tuning models, we trained each model with three different random seeds. We report the average score and the standard deviation.

Model		Test _{LR}		Test _{LC}	
		QM \uparrow	IFF \uparrow	QM \uparrow	IFF \uparrow
T5-Base	<i>QQ</i>	54.13	5.70 (15.17)	50.58	5.08 (16.64)
	<i>QS</i>	71.33	4.51 (8.03)	61.58	3.72 (8.31)
ChatGPT	<i>QQ</i>	92.94	11.24 (15.17)	73.66	5.98 (16.64)
	<i>QS</i>	94.57	12.21 (15.17)	79.67	6.81 (16.64)

Table 3: Model performances on longer interactions. For the IFF score, the perfect score for each setting is denoted in parentheses.

times longer than those in the training data. Notably, $Test_{LR}$ poses a challenge with regard to generalization for longer sequences, while $Test_{LC}$ additionally provides a challenge in generalization for unseen questions, resulting in a more challenging but more practical scenario. Details about the longer interaction generation process can be found in Appendix C.1.

Metric In this experiment, as can be seen in Table 3, we investigate whether a model trained on an average of approximately *three turns* can generalize to interactions of more than an average of 14 turns. Therefore, IM metric is deemed excessively strict. Instead, we report QM to measure the ratio of correctly answered turns. We further measure the Index of the First Failure (IFF), to capture where the model first generates an incorrect response. The methodology for calculating the IFF score is detailed in Appendix C.2.

Result The experiment results are presented in Table 3. Similarly with the findings in Section 5.1, scores from $Test_{LR}$ outperform that from $Test_{LC}$. This is likely because the model should generalize to the unseen interaction goals as well as longer interactions in $Test_{LC}$. Interestingly, IFF scores from all models are higher than the average length

of interactions in the training data. This suggests that the models can learn interaction-level compositionality and generalize to longer turns after being trained with EHR-SeqSQL. However, for the *QS* setting, T5-base model had a disadvantage in taking long interactions due to its input length constraint, which led to the lowest IFF score. On the other hand, ChatGPT outperformed the fine-tuned T5 models in both QM and IFF scores, highlighting its adaptability to longer interactions. Especially, in the random split, ChatGPT achieved an IFF score of 12.21 in the *QS* setting, quite close to the ideal score of 15.17.

5.3 Effects of Special Tokens for SQL

In this section, we evaluate the effect of two special tokens—`prev_query` and `prev_result`. These two tokens allow models to easily reference either the previous query or its execution result and thereby alleviate the decoding burden. Additionally, `prev_result` can further reduce the execution overhead by preventing duplicated sub-queries to be executed multiple times. We assess the utility of these tokens from two aspects: model performance and query execution efficiency.

5.3.1 Effects on Model Performance

Details First, we evaluate the impact of the special tokens on model performance. We compare a model trained with queries with the original EHR-SeqSQL, which include the special tokens, to another trained with its standard SQL query version. We ensure that all other training factors, such as model architecture, learning rate, or optimizer, remain consistent between the two models. To prevent the test set from being too simple, which could potentially undermine the impact of the special tokens, we use a compositional split. We report both QM and IM to assess the performance of the models at both the question and the interaction level.

Model		SQL		SQL [†]	
		QM [↑]	IM [↑]	QM [↑]	IM [↑]
T5-Base	QQ	75.37	52.49	89.97	82.49
	QS	82.10	64.35	86.46	75.38
ChatGPT	QQ	72.07	45.14	75.15	60.50
	QS	88.22	72.55	88.23	75.43

Table 4: Comparison of model performance trained with the original SQL queries and queries with our special tokens (SQL[†]).

Patient	Range	SQL		SQL [†]	
		avg.↓	med.↓	avg.↓	med.↓
1k	Full	0.22	0.01	0.18	0.01
	ST	0.20	0.09	0.09	0.01
10k	Full	9.01	0.05	7.96	0.03
	ST	3.25	1.17	1.15	0.15
46k	Full	164.69	0.19	139.44	0.08
	ST	31.57	4.60	6.00	0.52

Table 5: Execution time of SQL queries measured in databases of different sizes. We report the statistics for the entire queries (Full) and for the queries that contain the special token, `prev_result` (ST). Units are deciseconds (10^{-1}).

Result As shown in Table 4, the incorporation of special tokens consistently enhances the performance across all settings, regardless of the model variant or the type of interaction history. In the *QQ* setting, where the absence of prior query history makes the contextual questions more challenging, the special tokens contribute to a substantial performance increase. This is because those tokens simplify referencing previous questions. Specifically, at the question level, the use of these special tokens leads to a performance increase of 14.6%p (from 75.37 to 89.97), and 30%p (from 52.49 to 82.49) at the interaction level. In the *QS* setting, on the other hand, the reference is more straightforward because the previous SQL information is given. Still, special tokens enhance model performance by reducing the complexity of the target representation. ChatGPT demonstrates a robust performance even without the special tokens in the *QS* setting. We speculate this is because ChatGPT is trained on diverse data, which likely includes the standard SQL. Therefore, it might be familiar with standard text-to-SQL tasks, which is also consistent with the recent finding (Liu et al., 2023).

5.3.2 Effects on Execution time

Details We further evaluate the impact of our special tokens with respect to query execution ef-

iciency. Specifically, we compare the execution time of the queries from EHR-SeqSQL and that of their standard SQL version. The original database consists of the medical records of 1,000 patients and has a size of 95MB. We additionally constructed two larger databases following Lee et al. (2022), one with the medical records of 10,000 patients and another with all 46,520 patients that MIMIC provides. The size of these databases is 921MB and 5.06GB respectively. The results are reported on the queries from the test set in a random split, which covers all types of SQL queries in EHR-SeqSQL¹.

Result The result is shown in Table 5. We observed an 18% decrease in average execution time, in the original database. On a single query basis, the execution time is reduced at most to 99.89%, where the original SQL query has five nested queries inside. The effectiveness of the special token tends to increase with the size of the database. In the largest database, the special token yielded an average time reduction of 81.0%. Given that real-world databases typically contain huge amounts of data, we anticipate that the special token will yield a significant practical impact in multi-turn text-to-SQL environments.

6 Conclusion

We present EHR-SeqSQL, the first sequential text-to-SQL dataset designed for Electronic Health Record databases, aimed at improving interactivity, compositionality, and efficiency in text-to-SQL parsing. In terms of compositionality, our experiments show that the decomposition of questions significantly enhances the model’s ability to generalize to unseen interaction goals. Furthermore, it exhibited the potential for generalizing longer interactions that were not encountered during training. Additionally, we propose novel SQL-like tokens that enhance execution efficiency by reusing query execution results from the previous turn, as well as improve the performance of sequential text-to-SQL tasks by reducing the decoding burden on the model. As a step forward to interactive settings in text-to-SQL, we believe our dataset will serve as a valuable testbed for assessing contextual and interactive text-to-SQL tasks in EHRs and bridge the gap between industry needs and academic research.

¹We excluded the execution time that is nearly zero.

601 Limitations

602 Firstly, EHR-SeqSQL is generated through SQL
603 query decomposition from EHRSQL, which may
604 not fully represent real-world use cases. For ex-
605 ample, real-world questions are often ambiguous
606 and may contain unanswerable questions given the
607 database. Although handling ambiguous and unan-
608 swerable questions is not within the scope of this
609 work, it is an important topic that we would like
610 to explore in the future. Secondly, our dataset as-
611 sumes users possess a certain level of understand-
612 ing of the database schema. For example, when
613 asking about a patient’s records, users must first
614 identify the specific hospitalization ID for that pa-
615 tient. This assumption could limit the applicability
616 of our dataset in scenarios where users lack this
617 level of comprehension or knowledge. We note,
618 however, that in domain-specific settings, users are
619 typically aware of the database’s structure to cer-
620 tain extent.

621 Ethics Statement

622 Our main concerns are patient privacy and ensuring
623 that the proposed dataset as well as the construction
624 process are both compliant with the privacy regu-
625 lations. The MIMIC-III database, which requires
626 credentialed access via PhysioNet, is a database
627 that contains de-identified medical records of in-
628 tensive care patients. Lee et al. (2022) further de-
629 identified MIMIC-III for EHRSQL by shifting time
630 and values across the database. We use the database
631 constructed from EHRSQL, which carefully pre-
632 vents the re-identification of medical records from
633 the questions. We utilized the ChatGPT api to para-
634 phrase our manually generated question templates.
635 Each question template corresponds to the decom-
636 posed queries from EHRSQL and would unlikely
637 contain the sensitive medical information.

638 References

639 Jayetri Bardhan, Anthony Colas, Kirk Roberts, and
640 Daisy Zhe Wang. 2022. Drugehrqa: A question an-
641 swering dataset on structured and unstructured elec-
642 tronic health records for medicine related queries.
643 *arXiv preprint arXiv:2205.01290*.

644 Catherine Finegan-Dollak, Jonathan K Kummerfeld,
645 Li Zhang, Karthik Ramanathan, Sesh Sadasivam,
646 Rui Zhang, and Dragomir Radev. 2018. Improving
647 text-to-sql evaluation methodology. *arXiv preprint*
648 *arXiv:1806.09029*.

649 Yujian Gan, Xinyun Chen, Qiuping Huang, and
650 Matthew Purver. 2022. Measuring and improving
651 compositional generalization in text-to-sql via com-
652 ponent alignment. *arXiv preprint arXiv:2205.02054*.

653 Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang,
654 Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng
655 Wu. 2023. How close is chatgpt to human experts?
656 comparison corpus, evaluation, and detection. *arXiv*
657 *preprint arXiv:2301.07597*.

658 Jiaqi Guo, Zecheng Zhan, Yan Gao, Yan Xiao, Jian-
659 Guang Lou, Ting Liu, and Dongmei Zhang. 2019. To-
660 wards complex text-to-sql in cross-domain database
661 with intermediate representation. *arXiv preprint*
662 *arXiv:1905.08205*.

663 Moshe Hazoom, Vibhor Malik, and Ben Bogin. 2021.
664 Text-to-sql in the wild: a naturally-occurring dataset
665 based on stack exchange data. *arXiv preprint*
666 *arXiv:2106.05006*.

667 Charles T Hemphill, John J Godfrey, and George R
668 Doddington. 1990. The atis spoken language sys-
669 tems pilot corpus. In *Speech and Natural Language:*
670 *Proceedings of a Workshop Held at Hidden Valley,*
671 *Pennsylvania, June 24-27, 1990*.

672 Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. 2017.
673 Search-based neural structured learning for sequen-
674 tial question answering. In *Proceedings of the 55th*
675 *Annual Meeting of the Association for Computational*
676 *Linguistics (Volume 1: Long Papers)*, pages 1821–
677 1831.

678 Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H
679 Lehman, Mengling Feng, Mohammad Ghassemi,
680 Benjamin Moody, Peter Szolovits, Leo Anthony Celi,
681 and Roger G Mark. 2016. MIMIC-III, a freely accessi-
682 ble critical care database. *Scientific data*, 3(1):1–9.

683 Brenden Lake and Marco Baroni. 2018. Generalization
684 without systematicity: On the compositional skills
685 of sequence-to-sequence recurrent networks. In *In-*
686 *ternational conference on machine learning*, pages
687 2873–2882. PMLR.

688 Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu
689 Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo,
690 Jong-Yeup Kim, and Edward Choi. 2022. Ehrrsql: A
691 practical text-to-sql benchmark for electronic health
692 records. *Advances in Neural Information Processing*
693 *Systems*, 35:15589–15601.

694 Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiayi
695 Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu
696 Cao, Ruiying Geng, et al. 2023. Can llm already
697 serve as a database interface? a big bench for large-
698 scale database grounded text-to-sqls. *arXiv preprint*
699 *arXiv:2305.03111*.

700 Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S
701 Yu. 2023. A comprehensive evaluation of chat-
702 gpt’s zero-shot text-to-sql capability. *arXiv preprint*
703 *arXiv:2303.13547*.

704 Tom J Pollard, Alistair EW Johnson, Jesse D Raffa,
705 Leo A Celi, Roger G Mark, and Omar Badawi. 2018.
706 The eicu collaborative research database, a freely
707 available multi-center database for critical care re-
708 search. *Scientific data*, 5(1):1–13.

709 Linlu Qiu, Peter Shaw, Panupong Pasupat, Tianze Shi,
710 Jonathan Herzig, Emily Pitler, Fei Sha, and Kristina
711 Toutanova. 2022. Evaluating the impact of model
712 scale for compositional generalization in semantic
713 parsing. *arXiv preprint arXiv:2205.12253*.

714 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
715 Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
716 Wei Li, and Peter J Liu. 2020. Exploring the limits
717 of transfer learning with a unified text-to-text trans-
718 former. *The Journal of Machine Learning Research*,
719 21(1):5485–5551.

720 Peter Shaw, Ming-Wei Chang, Panupong Pasupat, and
721 Kristina Toutanova. 2020. Compositional general-
722 ization and natural language variation: Can a seman-
723 tic parsing approach handle both? *arXiv preprint*
724 *arXiv:2010.12725*.

725 Ping Wang, Tian Shi, and Chandan K Reddy. 2020.
726 Text-to-sql generation for question answering on elec-
727 tronic medical records. In *Proceedings of The Web*
728 *Conference 2020*, pages 350–361.

729 Kun Wu, Lijie Wang, Zhenghua Li, Ao Zhang, Xinyan
730 Xiao, Hua Wu, Min Zhang, and Haifeng Wang.
731 2021. Data augmentation with hierarchical sql-to-
732 question generation for cross-domain text-to-sql pars-
733 ing. *arXiv preprint arXiv:2103.02227*.

734 Tao Yu, Rui Zhang, He Yang Er, Suyi Li, Eric Xue,
735 Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze
736 Shi, Zihan Li, et al. 2019a. Cosql: A conversa-
737 tional text-to-sql challenge towards cross-domain nat-
738 ural language interfaces to databases. *arXiv preprint*
739 *arXiv:1909.05378*.

740 Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,
741 Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingn-
742 ing Yao, Shanelle Roman, et al. 2018. Spider: A
743 large-scale human-labeled dataset for complex and
744 cross-domain semantic parsing and text-to-sql task.
745 *arXiv preprint arXiv:1809.08887*.

746 Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern
747 Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li,
748 Bo Pang, Tao Chen, et al. 2019b. Sparc: Cross-
749 domain semantic parsing in context. *arXiv preprint*
750 *arXiv:1906.02285*.

751 Rong Zhou. 2022. Research on key performance index
752 prediction of distributed database based on machine
753 learning algorithm. In *International Conference on*
754 *Cognitive based Information Processing and Appli-*
755 *cations*, pages 563–567. Springer.

756 Xuanhe Zhou, Guoliang Li, Chengliang Chai, and Jian-
757 hua Feng. 2021. A learned query rewrite system
758 using monte carlo tree search. *Proceedings of the*
759 *VLDB Endowment*, 15(1):46–58.

Category	Example	# of questions
Independent	Q1. What is the specific item id of the hemoglobin lab test?	14,005 (44.22%)
Dependent		
Referential	Q1. Display the hospitalization ids of patient 76173. Q2. How many times each drug was prescribed <u>within A1</u> since 2101?	11,560 (36.50%)
Filtering	Q1. What are the calcium, total lab test values tested during A1? Q2. Retrieve <u>the last tested case in A2</u> .	5,947 (18.78%)
Modifying	Q1. During A1, what was the last measured value of A2? Q2. <u>What about the first measured case?</u>	507 (1.16%)

Table 6: Category of the questions within EHR-SeqSQL.

A.1 Details in Stage 2 in Section 3.1

761

In stage 2, we decompose the subqueries based on the SQL clauses. Decomposed clauses are parsed according to the logical order of execution of an SQL statement - WHERE, ORDER BY, HAVING, and SELECT. However, there are certain cases that are excluded from this decomposition. First, if a WHERE clause contains any aliases of table or column names without specifying the original name, it is not decomposed for clarity. For example, see `SELECT DISTINCT T1.C1 FROM (PREV_QUERY5) AS T1 WHERE T1.valuenum = 73.0`. If decomposed, the meaning of T1 would be ambiguous since we do not know the meaning of T1, which is not a column name in MIMIC-III. Thus we ensured the conditions or clauses with table aliases are always used with the subquery where the alias is defined. Second, we do not decompose the clauses for optional data cleansing SQL clauses. For example, see `SELECT microbiology events.org_name FROM microbiologyevents WHERE microbiologyevents.spec_type_desc = 'foot culture' AND microbiologyevents.org_name IS NOT NULL`. In this SQL query, the bolded subquery is intended to ensure all the selected rows have contents, excluding any NULL values. It is an optional condition and a SQL writing style choice which is not included in the original question, so we choose not to decompose such clauses. Also, the ten shortest question templates from EHRSQL are maintained without further splitting or modifications, which are not likely to be asked through multiple sentences. These include questions asking about the drug intake method, the cost of a lab test, or the number of current patients. We do not decompose GROUP BY but combine it either with SELECT or HAVING, since they are not explicitly expressed in the natural language questions (Guo et al., 2019; Wu et al., 2021).

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

A.2 Details in Stage 3 in Section 3.1

782

The BPE algorithm was applied to merge the subqueries from both Stage 1 and Stage 2. Specifically, we first derived SQL templates by masking condition values in the SQL queries obtained in Stage 2. Each template was treated as a token in the BPE algorithm, and the syntactically mergeable pairs among the most frequent token bigrams were repeatedly merged. A syntactically mergeable pair refers to a pair of queries that were originally a single SQL query but were separated due to our decomposition strategy. Since BPE algorithm calculates frequency without considering the relationship of token bigrams, we added this constraint to make every resultant SQL query to be executable. Then, we sampled half of each randomly to maximize the complexity of each subquery and the diversity of each interaction. Considering the total number of turns in each stage, we merged bigrams that appear more than 100 times in the subqueries from Stage 1 and bigrams that appear more than 150 times in the subqueries from Stage 2. The final SQL queries for EHR-SeqSQL are acquired throughout three stages.

783

784

785

786

787

788

789

790

791

792

793

Split	Random			Compositional		
	Train	Test	Test _L	Train	Test	Test _L
# of Question Temp	167	166	166	121	46	46
# of Interactions	8,546	649	100	6,375	2,820	100
# of Turns	29,438	2,231	1,417	22,134	9,535	1,564
Avg. # of Turns	3.44	3.44	14.17	3.47	3.38	15.64
Max # of Turns	9	9	25	8	8	27

Table 7: Statistics for two different splits.

A.3 List of NLQ templates

In Section 3.2, we mentioned the process for rule-based SQL-to-NLQ generation. Detailed steps of this process can be found in Table 8. We used specific SQL templates paired with NLQ templates for this process, as presented in Table 9. Each placeholder in the NLQ template is determined by the DB schema and the values present in the actual SQL query. A few examples showcasing the pairing of DB schema and NL expressions are available in Table 10. The placeholders related to *value*, *operation*, and *time* are used in the same way as described by Lee et al. (2022).

A.4 Prompt for paraphrasing

Figure 3 shows the prompt for paraphrasing the questions at an interaction level. For each NLQ template, key expressions (such as condition values, reference indices, etc.) that must be strictly preserved were indicated in order to maintain consistency. During the data construction process, we utilized the ChatGPT API to paraphrase template questions. Prior to paraphrasing, all specific values within each template are replaced with representative, generic values for their respective slots. Once paraphrased, these generic values are then realized back to their original form. This process naturally prevents any sensitive information from being sent to the ChatGPT server.

B Details on the Compositional Split

The concept of compositions and components is based on the SQL query since it is the common factor across an EHRSQL instance and its corresponding EHR-SeqSQL instance. Table 11 provides three examples of composition and components. There are EHRSQL questions which is an interaction goal of the corresponding interaction in EHR-SeqSQL, and their templates where the same anonymizing logic in deriving composition is applied.

We use a greedy algorithm to automatically split our data into training and test sets. Starting with all compositions assigned to the training set, we iteratively allocate a composition that has the maximum number of unique components to the test set while constraining all the components in the composition to exist in the training set, until no composition can be further assigned to the test set. You can find that the components in the last row are all present in the components of the first two rows. Thus, according to

Prompt for paraphrasing

---NLQ of text-to-SQL sequence.
 <Our NLQ template> e.g., Retrieve only the last taken case today from result3.

Our objective is to substitute each natural language query (NLQ) turn in a contextual text-to-SQL sequence with a paraphrased sentence that maintains the same meaning using alternate vocabulary or word order. It is essential to ensure that the semantic coherence of each turn with SQL is preserved. Also, each sentence should not be too long. Please write concisely and clearly. Please make 20 different NLQs for the question. You never ever change these words: <Key expression> e.g., ['last', 'today', 'result3']

Figure 3: Prompt for paraphrasing.

interaction goal	What were the top four frequent drugs that patients were prescribed within the same month after having been prescribed with nateglinide last year?			
idx	SQL Query	SQL Template	NLQ Template	Generated NLQ
1	SELECT admissions.subject_id, prescriptions.startdate FROM prescriptions JOIN admissions ON prescriptions.hadm_id = admissions.hadm_id WHERE prescriptions.drug = 'nateglinide' [time_filter_global]	SELECT table.column, table.column FROM table JOIN table ON table.column = table.column WHERE table.column = [val_placeholder] [time_filter_global]	List all [SELECT.col:admission.subject_id] and their [SELECT.col:prescriptions.startdate] associated with [val_placeholder:nateglinide] [time_filter_global:last year].	List all patient ids and their prescription time associated with nateglinide last year.
2	SELECT admissions.subject_id, prescriptions.drug, prescriptions.startdate FROM prescriptions JOIN admissions ON prescriptions.hadm_id = admissions.hadm_id WHERE [time_filter_global]	SELECT table.column, table.column FROM table JOIN table ON table.column = table.column WHERE [time_filter_global]	List all [SELECT.col:admission.subject_id] and their [SELECT.col:prescriptions.drug] and [SELECT.col:prescriptions.startdate] [time_filter_global:last year].	List all patient ids and their drugs and prescription time last year.
3	SELECT T2.drug, DENSE_RANK() OVER (ORDER BY COUNT(*) DESC) AS C1 FROM ([PREV_QUERY1]) AS T1 JOIN ([PREV_QUERY2]) AS T2 ON T1.subject_id = T2.subject_id WHERE T1.startdate <T2.startdate [time_filter_within] GROUP BY T2.drug	SELECT table.column, DENSE_RANK() OVER (ORDER BY COUNT(*) DESC) AS column FROM ([PREV1]) AS table JOIN ([PREV2]) AS table ON table.column = table.column WHERE table.column <table.column [time_filter_within] GROUP BY table.column	List the frequency rankings of [PREV:2] that patients received [time_filter_within:within the same month] after the [PREV:1].	List the frequency rankings of A2 that patients received within the same month after A1.
4	SELECT T3.drug FROM ([PREV_QUERY3]) AS T3 WHERE T3.C1 <= [n_rank]	SELECT table.column FROM ([PREV3]) AS table WHERE table.column <= [n_rank]	List the top [n_rank:four] [SELECT.col:prescriptions.drug] in [PREV:3].	List the top four drugs in A3.

Table 8: SQL-to-NLQ Generation Process.

SQL template	NLQ template
SELECT ([PREV0]) - ([PREV1])	What is the difference between the [PREV0] and [PREV1]?
SELECT ([PREV0]) [comparison] ([PREV1])	Is [PREV0] [comparison] than [PREV1]?
SELECT [agg_function](table.column) FROM ([PREV]) AS table	What is the [SELECT.col] of [PREV]?
SELECT COUNT(DISTINCT table.column) FROM ([PREV]) AS table	Count the number of patients in [PREV].
SELECT COUNT(DISTINCT table.column) FROM table #cond_parsed	Count the number of [PREV-1].
SELECT COUNT(DISTINCT table.column) FROM table WHERE [time_filter_global1]	Count the number of [PREV-1].
SELECT COUNT(*) FROM table WHERE table.column = ([PREV])	Count the number of [SELECT.col] associated with [PREV].
SELECT COUNT(*) FROM table WHERE table.column = [val_placeholder]	Count the number of [val_placeholder].
SELECT COUNT(*) FROM table WHERE table.column IN ([PREV])	Count the number of [SELECT.col] associated with [PREV].
SELECT COUNT(*)>0 FROM table WHERE table.column = [val_placeholder]	Has [val_placeholder] been admitted to the hospital?
SELECT COUNT(*)>0 FROM table WHERE table.column IN ([PREV])	Are there any [SELECT.col] in [PREV]?
SELECT SUM(table.column) FROM table WHERE table.column IN ([PREV])	What is the [SELECT.col] associated with [PREV]?
SELECT SUM(table.column) FROM table WHERE table.column IN ([PREV])	What is the total amount of [PREV-1]?
SELECT table.column FROM ([PREV]) AS table WHERE table.column [n_times]	Which [SELECT.col] is [n_times] in [PREV]?
SELECT table.column FROM ([PREV]) AS table WHERE table.column <= [n_rank]	List top [n_rank] [SELECT.col] in [PREV].
SELECT table.column FROM table	List all [SELECT.col] from [FROM.table].
SELECT table.column FROM table WHERE [age_group]	List all [SELECT.col] associated with patients aged [age_group].
SELECT table.column FROM table WHERE table.column = ([PREV])	What is the [SELECT.col] of [PREV]?
SELECT table.column FROM table WHERE table.column = [val_placeholder]	List all [SELECT.col] of [val_placeholder].
SELECT table.column FROM table WHERE table.column IN ([PREV])	List all [SELECT.col] associated with [PREV].
SELECT table.column, table.column FROM table	List all [SELECT.col.0] and [SELECT.col.1].
[PREV] [time_filter_exact1]	What was the [time_filter_exact2] measured case from [PREV-1]?
[PREV] [time_filter_global1_dec1]	Retrieve only the cases [time.verb] [time_filter_global1_dec1] from [PREV].
[PREV] [time_filter_global1_dec2]	Retrieve only the cases [time.verb] [time_filter_global1_dec2] from [PREV].
[PREV] [time_filter_global1]	Retrieve only the cases [time.verb] [time_filter_global1] from [PREV].
[PREV] AND [age_group]	Retrieve only the cases associated with patients aged [age_group] from [PREV_QEURY].
[PREV] AND table.column = ([PREV])	Retrieve only the cases associated with [PREV] from [PREV].
[PREV] AND table.column IN ([PREV])	Retrieve only the cases associated with [PREV] from [PREV].
[PREV] AND table.column IS NULL	What is the current one in [PREV]?
[PREV] WHERE [age_group]	Retrieve only the cases associated with patients aged [age_group] from [PREV].
[PREV] WHERE [time_filter_global1]	Retrieve only the cases [time.verb] [time_filter_global1] from [PREV].
[PREV] WHERE table.column = ([PREV])	Retrieve only the cases associated with [PREV] from [PREV].
[PREV] WHERE table.column IN ([PREV])	Retrieve only the cases associated with [PREV] from [PREV_QEURY].

Table 9: SQL & NLQ template.

DB schema	NL expression
admissions.admittime	admission time
admissions.dob	date of birth
admissions.dod	date of death
admissions.subject_id	patient
chartevents.charttime	chart time
chartevents.itemid	vital sign item id
chartevents.valuenum	value of vital sign
cost.hadm_id	hospital stay
diagnoses_icd.charttime	time of diagnosis
diagnoses_icd.icd9_code	diagnosis
diagnoses_icd.icd9_code	diagnosis ICD-9 code
inpuvents_cv.amount	volume of intake
inpuvents_cv.itemid	input event item id
labevents.itemid	lab test
labevents.itemid	lab test item id
labevents.valuenum	value of lab test
microbiologyevents.org_name	organism name
microbiologyevents.spec_type_desc	microbiology test
outputevents.itemid	output event item id
prescriptions.drug	drug
prescriptions.startdate	prescription time
procedures_icd.charttime	time of procedure
procedures_icd.hadm_id	hospital stay
procedures_icd.icd9_code	procedure ICD-9 code
procedures.icd9_code	procedure

Table 10: Examples of DB schema & NL expression pairs.

Interaction Goal	Template	Composition	Set of Components
Is the value of glucose of patient 71192 last measured on the first hospital visit less than the second to last value measured on the first hospital visit?	Is the value of {lab_name} of patient {patient_id} [time_filter_exact2] measured [time_filter_global2] [comparison] than the [time_filter_exact1] value measured [time_filter_global1]?	SELECT (SELECT labevents.valuenum FROM labevents WHERE labevents.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = {patient_id} [time_filter_global1]) AND labevents.itemid IN (SELECT d_labitems.itemid FROM d_labitems WHERE d_labitems.label = {lab_name}) [time_filter_exact1]) < (SELECT labevents.valuenum FROM labevents WHERE labevents.hadm_id in (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = {patient_id} [time_filter_global2]) AND labevents.itemid IN (SELECT d_labitems.itemid FROM d_labitems WHERE d_labitems.label = {lab_name}) [time_filter_exact2]))	1. SELECT admissions.hadm_id FROM admissions 2. [PREV_QUERY] WHERE admissions.subject_id = {patient_id} 3. [PREV_QUERY] AND [time_filter_global1] 4. SELECT d_labitems.itemid FROM d_labitems 5. [PREV_QUERY] WHERE d_labitems.label = {lab_name} 6. SELECT labevents.valuenum FROM labevents 7. [PREV_QUERY] WHERE labevents.hadm_id IN ([PREV_RESULT]) 8. [PREV_QUERY] AND labevents.itemid IN ([PREV_RESULT]) 9. [PREV_QUERY] [time_filter_exact1] 10. [PREV_QUERY] [time_filter_exact2] 11. SELECT ([PREV_RESULT]) [comparison] ([PREV_RESULT])
What was the maximum arterial bp [diastolic] of patient 18866 yesterday?	What was the [agg_function] of patient {patient_id} [time_filter_global1]?	SELECT [agg_function](chartevents.valuenum) FROM chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = {patient_id})) AND chartevents.itemid IN (SELECT d_items.itemid FROM d_items WHERE d_items.label = {vital_name} AND d_items.linksto = 'chartevents') [time_filter_global1]	1. SELECT admissions.hadm_id FROM admissions 2. [PREV_QUERY] WHERE admissions.subject_id = {patient_id} 3. SELECT icustays.icustay_id FROM icustays 4. [PREV_QUERY] WHERE icustays.hadm_id IN ([PREV_RESULT]) 5. SELECT d_items.itemid FROM d_items 6. [PREV_QUERY] WHERE d_items.label = {vital_name} AND d_items.linksto = 'chartevents' 7. SELECT chartevents.valuenum FROM chartevents 8. [PREV_QUERY] WHERE chartevents.icustay_id IN ([PREV_RESULT]) 9. [PREV_QUERY] AND chartevents.itemid IN ([PREV_RESULT]) 10. [PREV_QUERY] [time_filter_global1] 11. SELECT [agg_function](chartevents.valuenum) FROM chartevents WHERE chartevents.icustay_id IN ([PREV_RESULT]) AND chartevents.itemid IN ([PREV_RESULT]) [time_filter_global1]
Is the arterial bp [diastolic] of patient 25461 last measured on the last icu visit greater than the second to last value measured on the last icu visit?	Is the {vital_name} of patient {patient_id} [time_filter_exact2] measured [time_filter_global2] [comparison] than the [time_filter_exact1] value measured [time_filter_global1]?	SELECT (SELECT chartevents.valuenum from chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id from icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id from admissions WHERE admissions.subject_id = {patient_id}) [time_filter_global1]) AND chartevents.itemid IN (SELECT d_items.itemid from d_items WHERE d_items.label = {vital_name} AND d_items.linksto = 'chartevents') [time_filter_exact1]) > (SELECT chartevents.valuenum from chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id from icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id from admissions WHERE admissions.subject_id = {patient_id}) [time_filter_global2]) AND chartevents.itemid IN (SELECT d_items.itemid from d_items WHERE d_items.label = {vital_name} AND d_items.linksto = 'chartevents') [time_filter_exact2]))	1. SELECT admissions.hadm_id FROM admissions 2. [PREV_QUERY] WHERE admissions.subject_id = {patient_id} 3. SELECT icustays.icustay_id FROM icustays 4. [PREV_QUERY] WHERE icustays.hadm_id IN ([PREV_RESULT]) 5. [PREV_QUERY] AND [time_filter_global1] 6. SELECT d_items.itemid FROM d_items 7. [PREV_QUERY] WHERE d_items.label = {vital_name} AND d_items.linksto = 'chartevents' 8. SELECT chartevents.valuenum FROM chartevents 9. [PREV_QUERY] WHERE chartevents.icustay_id IN ([PREV_RESULT]) 10. [PREV_QUERY] AND chartevents.itemid IN ([PREV_RESULT]) 11. [PREV_QUERY] [time_filter_exact1] 12. [PREV_QUERY] [time_filter_exact2] 13. SELECT ([PREV_RESULT]) [comparison] ([PREV_RESULT])

Table 11: Examples of compositions and components.

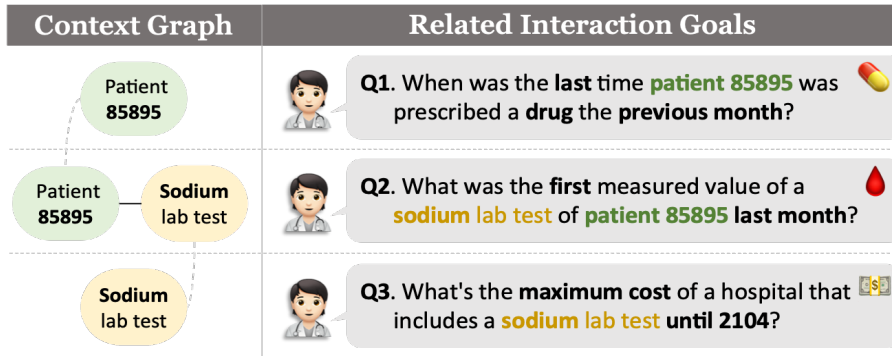


Figure 4: Related interaction goals and their context graph.

our split algorithm, if both the first two compositions are in the training set, the last composition can be assigned to the test set.

C Longer Interaction

C.1 Generation Process

As given in Figure 4, we define the concept of related interaction goal by using the context graph. Each interaction goal has its own context graph, whose nodes are defined by the specific condition values (e.g. conditions for patient, drug, or lab tests) in the original EHRSQL. Two questions are deemed related if they have overlapping condition values. Q1 and Q2 are related because they share the same patient 85895. However, Q1 and Q3 are not directly related without Q2 as a bridge. While we used three independent questions from the EHRSQL to demonstrate this concept in Figure 4, a longer and context-dependent interaction used in Section 5.2 is depicted in Figure 5. In this example for longer interaction, five EHRSQL questions are connected, each paired with an original SQL query. In contrast, EHR-SeqSQL involves context-dependent interaction spanning twelve turns, each paired with a query with our special tokens for SQL.

C.2 IFF Score Calculation

$$IFF = \begin{cases} n + 1, & \text{if all turns are correct} \\ k, & \text{otherwise} \end{cases}$$

k denotes the specific turn number in the interaction where the first incorrect response occurs. n represents the total number of turns in the interaction. The final IFF score for a test set is calculated as the average IFF score across all interactions. Hence, if the model achieves a perfect score for every interaction, the IFF score becomes one plus the average number of turns of all interactions.

Note that the perfect IFF score for QQ and QS setting is different in Table 3. This is because we removed the test sample that exceeds the maximum token length of T5 when the interaction history is concatenated with the current question for a fair comparison.

D Fine-tuning Baseline

D.1 Configuration

We fine-tuned the T5-Base and T5-3B using the Adam optimizer, with a global batch size of 32. The learning rate were set to 1e-4 for T5-Base and 5e-5 for T5-3B, respectively. For the other hyperparameter configurations, we followed the settings used in EHRSQL. All experiments were carried out on either a single A100 80G GPU or a A6000 48G GPU. The training process typically took around 10 hours for T5-Base and around 24 hours for T5-3B.

EHRSQL



Q. What was the name of the specimen test that patient 85895 has gotten the first the last month?

```
SQL SELECT microbiologyevents.spec_type_desc
FROM microbiologyevents
WHERE microbiologyevents.hadm_id IN (
  SELECT admissions.hadm_id
  FROM admissions
  WHERE admissions.subject_id = 85895 )
AND datetime(microbiologyevents.charttime,'start of month') =
datetime(current_time,'start of month','-1 month')
ORDER BY microbiologyevents.charttime ASC LIMIT 1
```



Q. When was the last time patient 85895 was prescribed a drug the previous month?

```
SQL SELECT prescriptions.startdate
FROM prescriptions
WHERE prescriptions.hadm_id IN (
  SELECT admissions.hadm_id
  FROM admissions
  WHERE admissions.subject_id = 85895 )
AND datetime(prescriptions.startdate,'start of month') =
datetime(current_time,'start of month','-1 month')
ORDER BY prescriptions.startdate DESC LIMIT 1
```



Q. What was patient 85895's first output on last month/05?

```
SQL
SELECT outputevents.itemid
FROM outputevents
WHERE outputevents.icustay_id IN (
  SELECT icustays.icustay_id
  FROM icustays
  WHERE icustays.hadm_id IN (
    SELECT admissions.hadm_id
    FROM admissions
    WHERE admissions.subject_id = 85895 ) )
AND datetime(outputevents.charttime,'start of month') =
datetime(current_time,'start of month','-1 month')
AND strftime('%d',outputevents.charttime) = '05'
ORDER BY outputevents.charttime ASC LIMIT 1
```



Q. What was the first measured value of a sodium, whole blood lab test of patient 85895 last month?

```
SQL
SELECT labevents.valuenum
FROM labevents
WHERE labevents.hadm_id IN (
  SELECT admissions.hadm_id
  FROM admissions
  WHERE admissions.subject_id = 85895 )
AND labevents.itemid IN (
  SELECT d_labitems.itemid
  FROM d_labitems
  WHERE d_labitems.label = 'sodium, whole blood' )
AND datetime(labevents.charttime,'start of month') =
datetime(current_time,'start of month','-1 month')
ORDER BY labevents.charttime ASC LIMIT 1
```



Q. What's the maximum cost of a hospital that includes a sodium, whole blood lab test until 2104?

```
SQL
SELECT MAX(T1.C1) FROM (
  SELECT SUM(cost.cost) AS C1
  FROM cost
  WHERE cost.hadm_id IN (
    SELECT labevents.hadm_id FROM labevents
    WHERE labevents.itemid IN (
      SELECT d_labitems.itemid FROM d_labitems
      WHERE d_labitems.label = 'sodium, whole blood' ) ) )
AND strftime('%Y',cost.chargetime) <= '2104'
GROUP BY cost.hadm_id ) AS T1
```

EHR-SeqSQL (Ours)



Q1. Give me the hospitalization ids of patient 85895.

```
SQL† SELECT admissions.hadm_id
FROM admissions
WHERE admissions.subject_id = 85895
```



Q2. Last month, what was the first specimen to be tested in A1?

```
SQL† SELECT microbiologyevents.spec_type_desc
FROM microbiologyevents WHERE microbiologyevents.hadm_id
IN ( PREV_RESULT1 ) AND datetime(microbiologyevents.charttime,'start of
month') = datetime(current_time,'start of month','-1 month')
ORDER BY microbiologyevents.charttime ASC LIMIT 1
```



Q3. What was the prescription time for the last A1 entry last month?

```
SQL†
SELECT prescriptions.startdate
FROM prescriptions
WHERE prescriptions.hadm_id IN ( PREV_RESULT1 )
AND datetime(prescriptions.startdate,'start of month') =
datetime(current_time,'start of month','-1 month')
ORDER BY prescriptions.startdate DESC LIMIT 1
```



Q4. List all ICU stay ids associated with patient 85895.

```
SQL† SELECT icustays.icustay_id FROM icustays
WHERE icustays.hadm_id IN ( PREV_RESULT1 )
```



Q5. Give me a list of output event item ids from outputevents table.

```
SQL† SELECT outputevents.itemid FROM outputevents
```



Q6. Show me only the cases that correspond to A4 in A5.

```
SQL† PREV_QUERY5 WHERE outputevents.icustay_id IN ( PREV_RESULT4 )
```



Q7. Provide me with the case that was first recorded on last month/05 from A6.

```
SQL† PREV_QUERY6 AND datetime(outputevents.charttime,'start of month') =
datetime(current_time,'start of month','-1 month')
AND strftime('%d',outputevents.charttime) = '05'
ORDER BY outputevents.charttime ASC LIMIT 1
```



Q8. Which labels are assigned to A7?

```
SQL† SELECT d_items.label FROM d_items
WHERE d_items.itemid IN ( PREV_RESULT7 )
```



Q9. I'm looking for the specific item id associated with the sodium, whole blood lab test. can you help?

```
SQL† SELECT d_labitems.itemid FROM d_labitems
HERE d_labitems.label = 'sodium, whole blood'
```



Q10. Last month, what was the first value that was tested for A9 during the first test of A1?

```
SQL† SELECT labevents.valuenum FROM labevents
WHERE labevents.hadm_id
IN ( PREV_RESULT1 ) AND labevents.itemid IN ( PREV_RESULT9 )
AND datetime(labevents.charttime,'start of month') =
datetime(current_time,'start of month','-1 month')
ORDER BY labevents.charttime ASC LIMIT 1
```



Q11. List the hospitalization ids that correspond to A9.

```
SQL† SELECT labevents.hadm_id FROM labevents
WHERE labevents.itemid IN ( PREV_RESULT9 )
```



Q12. What is the maximum price for hospitalization associated with A11 until 2104?

```
SQL† SELECT MAX(T1.C1)
FROM (
  SELECT SUM(cost.cost) AS C1 FROM cost
  WHERE cost.hadm_id IN ( PREV_RESULT11 )
  AND strftime('%Y',cost.chargetime) <= '2104'
  GROUP BY cost.hadm_id ) AS T1
```

Patient 85895



Patient 85895

&

sodium, whole blood lab test



sodium, whole blood lab test

Figure 5: Example of Longer Interaction.

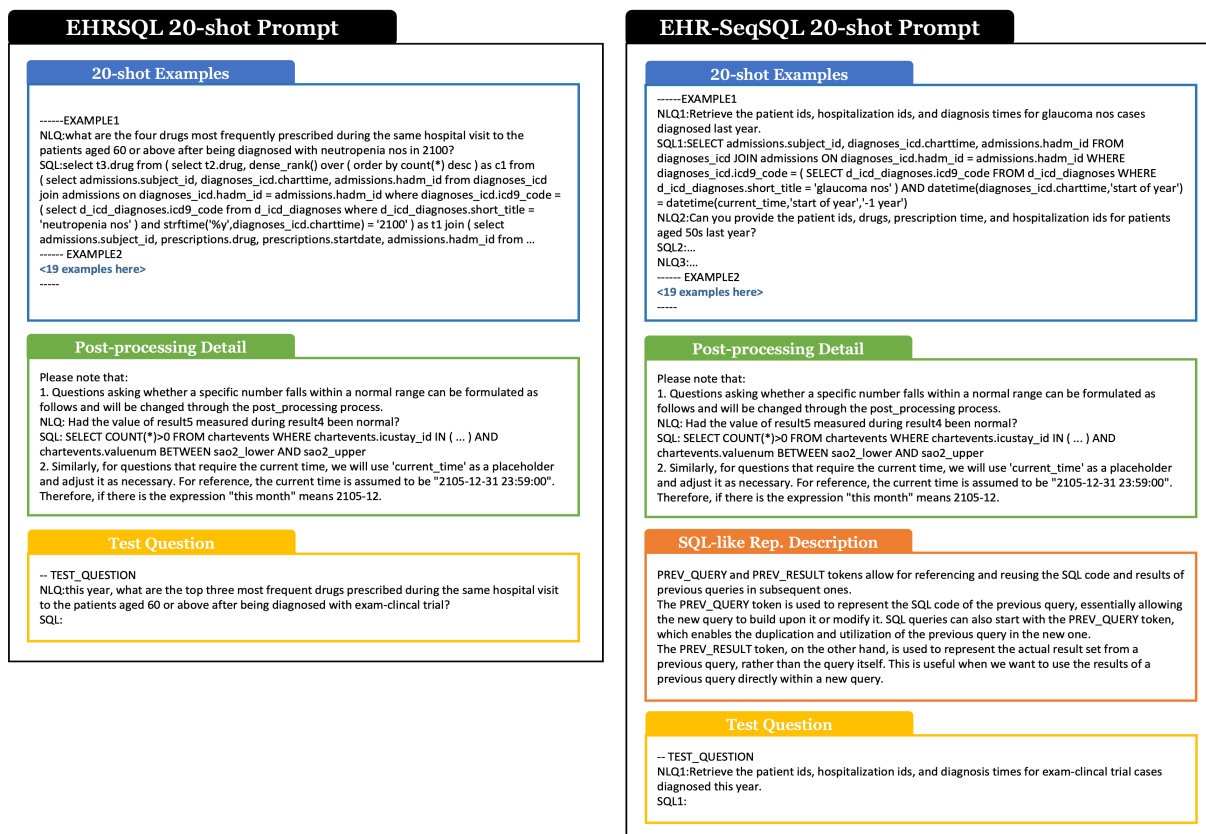


Figure 8: Prompt configuration for chatGPT.

E In-context Learning Baseline

850

E.1 Prompt Configuration

851

Figure 8 demonstrates how we configured the prompt of ChatGPT for both the EHRSQL and EHR-SeqSQL. Unlike EHRSQL where the target representation is standard SQL, EHR-SeqSQL include the special tokens, which are first introduced in this work. Thus, we include a simple description for `PREV_QUERY` and `PREV_RESULT` tokens in the prompt. In the *QS* setting during the inference process, for the current turn, the interaction history includes both previous questions and queries generated by ChatGPT.

852

853

854

855

856

857

E.2 Few-Shot Learning Approach for EHR-SeqSQL

858

We developed a new method for few-shot learning for the multi-turn, context-dependent setting of EHR-SeqSQL. This method is detailed in Figure ???. To begin, we created two corpora: one for interaction-level training data and another for turn-level training data. For our few-shot learning approach, we retrieved examples, with half based on the interaction history and the other half based on the current question. In our experiments, we used a total of 20 examples for few-shot learning.

859

860

861

862

863

F Use Cases of the Special Tokens

864

Table 12 illustrates an example of the interactions in EHR-SeqSQL where the special token significantly reduce the length of target representation as well as the execution time. Due to readability and spatial issues, the standard SQL version of each turn has been omitted from the table. You can see the standard SQL would be quite lengthy by looking at the SQL[†] that the used special token refers to.

865

866

867

868

Index	Question	Target Representation (SQL [†])	Time (SQL [†])	Standard SQL	Time (SQL)
1	Which icu stay ids are associated with patient 30826 on the current hospital visit?	SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)	0.334	SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)	0.334
2	Could you tell me the item id for weight, please?	SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents'	0.334	SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents'	0.333
3	During result1, what was the last value of result2 that was measured?	SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (PREV_RESULT1) AND chartevents.itemid IN (PREV_RESULT2) ORDER BY chartevents.charttime DESC LIMIT 1	20.352	SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)) AND chartevents.itemid IN (SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents') ORDER BY chartevents.charttime DESC LIMIT 1	31.362
4	What about the first measured case?	SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (PREV_RESULT1) AND chartevents.itemid IN (PREV_RESULT2) ORDER BY chartevents.charttime ASC LIMIT 1	19.684	SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)) AND chartevents.itemid IN (SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents') ORDER BY chartevents.charttime ASC LIMIT 1	30.695
5	What is the variation between result3 and result4?	SELECT (PREV_RESULT3) - (PREV_RESULT4)	0.000	SELECT (SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)) AND chartevents.itemid IN (SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents') ORDER BY chartevents.charttime DESC LIMIT 1) - (SELECT chartevents.valuenum FROM chartevents WHERE chartevents.icustay_id IN (SELECT icustays.icustay_id FROM icustays WHERE icustays.hadm_id IN (SELECT admissions.hadm_id FROM admissions WHERE admissions.subject_id = 30826 AND admissions.dischtime IS NULL)) AND chartevents.itemid IN (SELECT d_items.itemid FROM d_items WHERE d_items.label = 'admit wt' AND d_items.linksto = 'chartevents') ORDER BY chartevents.charttime ASC LIMIT 1)	60.722

Table 12: Example of an interaction in EHR-SeqSQL where the target representations (SQL[†]) contain special tokens. We also report average execution times in milliseconds(10^{-3}), where the queries are executed three times. Their standard SQL versions are also reported for comparison.