
Schema Guidance as Managed Context for Enhanced Agentic SQL Generation in Clinical Databases

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models have substantially improved text-to-SQL across tasks
2 that utilize relational database systems, but applications to real-world clinical
3 queries remain unreliable. We argue that the bottleneck is not simply a lack
4 of context, but a failure to *manage* context. We introduce *schema guidance* as
5 a first-class representation of managed context for agentic SQL generation: a
6 structured, updateable object that encodes task-relevant schema knowledge such
7 as key filters, common joins, anti-patterns, example query patterns, and critical
8 notes. Unlike standard prompting approaches, we treat context as a structured
9 object that explicitly constrains the model’s search space. We formalize schema
10 guidance as both (i) a compressed representation that motivates the approach and
11 (ii) a search-space reduction mechanism that accounts for improved decoding
12 performance. We construct guidance via a fixed pipeline that is shared across
13 databases. We show empirically that schema guidance improves construction of
14 SQL queries in non-clinical SPIDER datasets, open EHRSQL MIMIC-III and eICU
15 datasets, and a closed institutional dataset with representative complex clinical
16 queries. These results suggest that learning the right context representation is a
17 critical and underexplored axis for reliable agentic SQL generation in real-world
18 clinical settings.

19 1 Introduction

20 Text-to-SQL has become one of the clearest demonstrations of how large language models (LLMs)
21 can turn natural language into executable programs. For electronic health records, which primarily
22 use relational database management systems, recent LLM-based text-to-SQL systems have shown
23 improvements using retrieval augmented generation [Ziletti and D’Ambrosi, 2024, 2025], ensemble
24 LLM approaches [Gundabathula and Kolar, 2024], and seq2seq based methods Wang et al. [2020].
25 Other approaches have focused on modifying database structure entirely, either creating their own
26 database structures Katz et al. [2024] or adopting graph databases [Thio et al., 2026].

27 However, reliability remains brittle in real clinical databases, where query generation must cope
28 with highly complex logic with deeply nested heterogeneous schemas, ambiguous or noisy database
29 structures, and domain-specific business rules. The result is a familiar pattern: the model often
30 identifies the general task correctly, but fails on the exact filter, wrong join path, spurious table, or
31 clinically invalid value mapping.

32 The central problem may not merely be *insufficient context*; it is *poorly managed context*. [Liu et al.,
33 2023, Modarressi et al., 2025] Most existing systems increase context quantity by retrieving more
34 schema text, values, or exemplars. [Hong et al., 2024, Caferoğlu et al., 2025] Some improve context
35 quality through retrieval augmentation. [Shen et al., 2023, Toteja et al., 2025] An important question
36 is what *representation of context* an agent should maintain as a stable object over time. This question

37 is especially important in clinical SQL generation, where raw schema dumps are noisy and expensive,
38 while clinically correct querying depends on compact but high-value knowledge such as: which table
39 is the canonical anchor for a cohort, which dimensions are mandatory to disambiguate a request,
40 which joins are common and which are anti-patterns, and which filters should never be omitted.

41 This paper proposes *schema guidance* as the missing abstraction layer. Schema guidance is a
42 structured representation attached to schema elements to guide query construction rather than merely
43 describe tables. A guidance entry can encode key filters, commonly joined tables, critical notes, and
44 `do_not` rules. This shifts the locus of optimization away from prompt wording alone and toward the
45 *context object* that shapes planning and decoding.

46 We operationalize guidance via a schema-guidance construction pipeline. Its distinctive design choice
47 is *information asymmetry*: a guidance optimizer can inspect full schemas and aggregate benchmark
48 failures, but the runtime SQL generator receives only compressed guidance. This makes guidance a
49 distillation target rather than just retrieved text. Conceptually, we view guidance as a compressed
50 sufficient structure for generation; mechanistically, we view it as a search-space shaping device that
51 prunes irrelevant hypotheses while aiming to keep clinically valid programs reachable.

52 **Clinical text-to-SQL Benchmarks** Single-domain benchmarks for text-to-SQL tasks in the clinical
53 domain have relied on question-answer pairs derived from surveys generated by providers Lee et al.
54 [2024], Wang et al. [2020]. These are built primarily on open source EHR databases (MIMIC-III
55 and eICU) and tailored to answer tasks requiring reliable text-to-SQL over electronic health records.
56 These do not reflect the broad range of clinical questions that may lie outside of ICU-heavy databases,
57 with queries that span operational as well as clinical tasks. We provide results on existing open
58 datasets but contribute a real-world closed institutional dataset, illustrating the potential of our
59 approach to improve performance across multiple complex clinical settings.

60 Our contributions are:

- 61 1. **Schema guidance as a database-invariant construction.** We formalize *schema guidance* as
62 a structured context object produced by a fixed three-stage construction (Section 3.3): input
63 normalization, a join-neighborhood graph, per-table generation, deterministic validation
64 with bounded repair, and budget-constrained summarization with a $(1-1/e)$ coverage
65 guarantee under the token budget. Only the per-database schema and configuration vary; the
66 construction itself is shared across databases. We call this property *construction invariance*
67 and adopt it as our operational notion of *database-agnostic*. The formal policy definition
68 and the submodular coverage argument are deferred to Appendices A and B.
- 69 2. **Information-asymmetric generation as search-space reduction.** The guidance optimizer
70 may inspect the full schema, but the SQL generator sees only the compressed, budget-
71 feasible guidance object. We frame this asymmetry as a search-space reduction at decoding
72 time. A recall-preserving formal statement appears in Appendix B.1.
- 73 3. **Empirical evaluation across non-clinical, open clinical, and closed clinical settings.** We
74 apply the same construction to multiple open-source clinical databases and a real-world
75 closed institutional benchmark, and compare against an unguided baseline given the raw
76 schema truncated to the same token budget. The held-out evaluation tests construction
77 invariance as the operational test of generalization: gains accrue to the shared guidance
78 construction rather than to static prompting or one-shot retrieval.

79 2 Problem Setting

80 Let $x \in \mathcal{X}$ denote a natural-language clinical query, $\mathcal{S} \in \mathfrak{S}$ a relational schema, and y^* a target SQL
81 program. A text-to-SQL system is a stochastic mapping

$$p_f(\hat{y} \mid x, c)$$

82 that produces SQL hypotheses conditioned on a context c . Two choices for c are of interest: an
83 unguided baseline that serializes the raw schema and truncates it to fit a fixed token budget B
84 (formally $\pi_{g_0}(\mathcal{S}) = \text{Trunc}_B(\mathcal{S})$; see Appendix A), and a derived guidance object $c = g$ obtained by
85 a structured transformation of \mathcal{S} . We write

$$\hat{y} \sim \text{Decode}(p_f(\cdot \mid x, c))$$

86 for a fixed decoding rule (e.g., greedy, beam, or best-of- k with execution feedback), held constant
87 across all conditions compared in this paper.

88 **Definition 2.1** (Schema-guidance policy). *Let Θ be a space of database-specific configuration*
89 *parameters (e.g., key/time naming heuristics, slowly-changing-dimension flag handling, sensitivity*
90 *rules). A schema-guidance policy is a parameterized mapping*

$$\pi_g : \mathcal{S} \times \Theta \rightarrow \mathcal{G}_B, \quad g = \pi_g(\mathcal{S}; \theta),$$

91 *producing a structured context object containing key filters, commonly joined tables, field notes,*
92 *example query patterns, and critical notes.*

93 We treat $\pi_g(\cdot; \theta)$ as a *procedure* that is invariant across databases: only the inputs \mathcal{S} and the
94 configuration θ vary, while the construction itself is shared. Section 3.3 specifies the procedure;
95 Section 4 reports its instantiation across the clinical databases studied as $g^{(d)} = \pi_g(\mathcal{S}^{(d)}; \theta^{(d)})$. We
96 refer to this property as *construction invariance* and treat it as our operational notion of *database-*
97 *agnostic*; we make no formal categorical claim.

98 **Empirical objectives.** We evaluate guidance policies on a held-out validation set

$$\mathcal{D}_{\text{val}} = \{(x_i, \mathcal{S}_i, y_i^*)\}_{i=1}^n.$$

99 Let $u(x, y^*, \hat{y}) \in [0, 1]$ be a bounded utility (e.g., execution match, Code F1, Table F1) with
100 $u(x, y^*, y^*) = 1$. The empirical utility of a policy is

$$\hat{J}(\pi_g) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\hat{y}_i \sim \text{Decode}(p_f(\cdot | x_i, \pi_g(\mathcal{S}_i; \theta)))} [u(x_i, y_i^*, \hat{y}_i)].$$

101 **Joint reporting** For each database d , we report $\hat{J}(\pi_g)$ for the guidance policy π_g , alongside an
102 unguided baseline π_{g_0} (raw schema, truncated to B). Improvements due to guidance are assessed
103 through pairwise comparisons.

104 3 Method

105 3.1 Schema guidance as runtime context

106 The operational role of schema guidance follows directly from the problem setting (Section 2): the
107 generator $p_f(\hat{y} | x, c)$ has two natural choices for c , namely the raw schema ($c = \mathcal{S}$) and a derived
108 guidance object ($c = g \in \mathcal{G}_B$). The unguided choice $c = \mathcal{S}$ is generally infeasible at clinical
109 scale once $|\mathcal{S}| > B$; truncation produces the baseline policy π_{g_0} used throughout our experiments.
110 Schema guidance is the alternative: a structured context object produced by a construction π_g that is
111 shared across databases (Section 3.3). Our approach augments runtime context by $g = \pi_g(\mathcal{S}; \theta)$, the
112 schema-derived guidance. The schema-derived component g is the only part of c produced by π_g .

113 3.2 Information asymmetry between construction and inference

114 We enforce an explicit asymmetry over which arguments each stage may inspect:

- 115 • the construction π_g may read the full schema \mathcal{S} and the configuration θ , both of which are
116 construction-time inputs and do not depend on \mathcal{D}_{val} ;
- 117 • the generator $p_f(\hat{y} | x, c)$ at initialization receives only the guidance object $c = g$ with $|g| \leq B$;
118 any further per-table information acquired during decoding is mediated by tool calls (Appendix D),
119 and the raw schema \mathcal{S} is never delivered as a flat artifact.

120 We do not feed validation outcomes from \mathcal{D}_{val} back into the construction of π_g or $\theta^{(d)}$; held-out
121 databases are constructed from $\mathcal{S}^{(d)}$ and $\theta^{(d)}$ alone, so the asymmetry concerns what each side may
122 inspect about the schema, not what either side knows about validation labels. The asymmetry results
123 in the generator’s admissible hypothesis space at inference time to be the support of $\text{Decode}(p_f(\cdot |$
124 $x, g))$. Because π_g is a procedure shared across databases (Section 3.3), improving θ or any stage of
125 ϕ generalizes across schemas without retraining f . We therefore treat π_g as the optimization surface
126 and f as fixed throughout the paper.

127 3.3 Schema-guidance construction pipeline

128 We construct $g^{(d)} = \pi_g(\mathcal{S}^{(d)}; \theta^{(d)})$ via a fixed three-stage pipeline whose only per-database inputs
129 are the raw schema and the configuration $\theta^{(d)}$. All stages are evaluated empirically. Full operator
130 signatures and the validator specification are given in Appendix B.

131 **Stage 1: Input normalization.** We canonicalize raw schema metadata for each database into
132 a uniform per-database JSON of tables, columns, data types, optional textual descriptions, and
133 foreign-key relationships (extracted from declarations, or inferred from naming conventions when
134 absent).

135 **Stage 2: Schema graph and join neighborhoods.** We build a foreign-key graph over tables and
136 define each table’s *one-hop join neighborhood*; per-table guidance is conditioned on its neighborhood
137 so that join structure is reflected locally. Multi-hop structure is recovered indirectly via overlapping
138 neighborhoods at the summarization step.

139 **Stage 3: Per-table guidance generation.** For each table T_j , a per-table operator ϕ produces a
140 structured artifact g_j with slots for indexed columns, performance filters, common joins, field notes,
141 and example query patterns. The operator composes (i) a deterministic *heuristic skeleton* that fills
142 index/filter/join slots from column names, types, and declared relationships; (ii) an *LLM enrichment*
143 step conditioned on the table, its neighborhood, and any domain hints in $\theta^{(d)}$; and (iii) an optional
144 *curation merge* that overlays subject-matter-expert annotations when available.

145 **Per-database instantiation.** The complete policy is the composition

$$\pi_g(\mathcal{S}^{(d)}; \theta^{(d)}) = \text{Summarize}\left(\left\{\phi(T_j, \mathcal{N}(T_j); \theta^{(d)})\right\}_{j=1}^{m_d}; B\right),$$

146 applied independently to each database d . The procedure $\pi_g(\cdot; \cdot)$ is identical across all databases
147 studied; only the inputs $(\mathcal{S}^{(d)}, \theta^{(d)})$ change. Whether this construction invariance suffices for transfer
148 to unseen schemas is decided empirically by the held-out database evaluation.

149 3.4 Guidance as compressed sufficient structure and search-space reduction view

150 **Compressed sufficient structure.** We do not claim sufficiency in the information-theoretic sense.
151 Instead, we use *compressed sufficient structure* as a working notion of what guidance should encode.
152 Let $\mathcal{H}(\mathcal{S}, x)$ denote the SQL programs compatible with the full schema \mathcal{S} for query x , and let
153 $\mathcal{H}_g(\mathcal{S}, x) \subseteq \mathcal{H}(\mathcal{S}, x)$ denote the support effectively explored by the generator under guidance g . We
154 call g *task-sufficient* for (x, \mathcal{S}, y^*) when $y^* \in \mathcal{H}_g(\mathcal{S}, x)$, and *recall-preserving on a distribution \mathcal{D}*
155 when

$$\Pr_{(x, \mathcal{S}, y^*) \sim \mathcal{D}}[y^* \in \mathcal{H}_g(\mathcal{S}, x)] \geq 1 - \delta$$

156 for some small δ . Recall-preservation is a *necessary* condition for guidance to help: if $y^* \notin \mathcal{H}_g$,
157 decoding under g strictly loses regardless of what is pruned. Whether π_g achieves recall-preservation
158 in practice on real clinical schemas is an empirical question; the construction’s Stage 4 validator
159 certifies well-formedness of *what is kept*, not coverage of unseen queries (Appendix B.1).

160 **Search-space reduction.** Decoding under context c produces \hat{y} from a distribution supported on
161 $\mathcal{H}_c(\mathcal{S}, x)$. We adopt one working assumption about LLM decoders: that the probability of decoder
162 failure is upper bounded by a non-decreasing function $\varepsilon(\cdot)$ of the realized candidate-set size. Under
163 this assumption, recall-preserving and precision-improving guidance gives

$$\Pr[\hat{y} \neq y^* \mid g] \leq \varepsilon(|\mathcal{H}_g(\mathcal{S}, x)|) \leq \varepsilon(|\mathcal{H}(\mathcal{S}, x)|).$$

164 This is the formal content of the search-space-reduction view: gains scale with how much \mathcal{H}_g shrinks
165 \mathcal{H} *without losing the gold program* (Appendix B.1). The same logic predicts *where* gains should
166 concentrate. Spurious table selection and structurally invalid joins are exactly the failure modes
167 whose support is large in $\mathcal{H}(\mathcal{S}, x)$ when the schema is large; pruning them tightens the bound and
168 should manifest in *structural* metrics (Code F1, Table F1) before answer-level metrics.

169 **Connection to agentic context.** This perspective is useful beyond text-to-SQL. An agent does
170 not only retrieve context; it must choose which context to *maintain* as a stable representation
171 between steps. We argue that for SQL generation, the right stable representation encodes query-
172 construction knowledge—joins, anti-patterns, mandatory filters, do-not rules—rather than raw schema
173 text. Schema guidance is the operational instantiation: the construction pipeline (Section 3.3) produces
174 the representation.

175 **Empirical signature.** The mechanism makes a falsifiable prediction: if guidance primarily com-
176 presses *structural* hypothesis support while preserving gold reachability, gains in Code F1 and Table
177 F1 should exceed gains in execution-match. Our experiments (Section 4) are consistent with this
178 prediction.

179 4 Experimental Setup

180 4.1 SPIDER

181 To evaluate our method on widely studied text-to-SQL tasks, we adopt a unified test set constructed
182 from multiple classical benchmarks, following the preprocessing provided in prior work.[Yu et al.,
183 2019] Test examples are aggregated from eight text-to-SQL datasets across multiple domains; each
184 example consists of a natural language query, a corresponding ground-truth SQL query, and auxiliary
185 metadata. The SQL queries are rewritten to follow a standardized format consistent with the SPIDER
186 dataset, enabling uniform evaluation across datasets. Additional details in Appendix C

187 4.2 MIMIC-III

188 In the EHRSQL benchmark, MIMIC-III serves as one of the underlying relational databases over
189 which natural-language questions are paired with executable SQL queries. This makes it partic-
190 ularly suitable for evaluating text-to-SQL systems in the clinical domain, where questions often
191 require temporal reasoning, aggregation, and multi-table joins [Lee et al., 2024]. In our work, we
192 follow the benchmark setting provided by EHRSQL rather than reconstructing the database schema
193 independently.

194 4.3 eICU

195 Within EHRSQL, eICU complements MIMIC-III by providing a more heterogeneous clinical environ-
196 ment for evaluating text-to-SQL generalization [Lee et al., 2024]. Using both MIMIC-III and eICU is
197 important because it reduces the risk that performance reflects overfitting to a single hospital-specific
198 schema or documentation style. Therefore, the benchmark offers a more realistic test bed for clinical
199 question answering over relational EHR databases.

200 4.4 Closed Institutional Benchmark of Real World Queries

201 From a tertiary medical center 554 SQL queries tied to clinical requests to Information Technology
202 for data extraction were collated. Each SQL query was linked to a natural language requests. If
203 multiple queries were linked to the original request, natural language queries were derived from the
204 individual SQL query using zero shot gemini-2.5-pro.

205 5 Benchmark and Evaluation Metrics

206 5.1 Benchmark construction

207 We construct a SQL-centric benchmark from a corpus of clinically validated queries organized as
208 a hierarchy of task folders. Each folder corresponds to a high-level clinical task, described by a
209 human-authored *solution description*, and contains one or more SQL queries representing different
210 data extraction components (e.g., demographics, labs, medications).

211 Each SQL query is paired with a natural language *task description* derived from the query using a
212 preprocessing pipeline. This results in a dataset of approximately 554 query-level test cases, each

213 defined as:

$$(x_i, y_i^*),$$

214 where x_i is a natural language task description and y_i^* is the corresponding gold SQL query.

215 Multiple queries within a folder may share overlapping logic (e.g., cohort definitions and join
216 patterns), which is intentional and reflects real-world analytical workflows. This allows evaluation of
217 both independent query generation and implicit structural consistency.

218 5.2 Evaluation setup

219 For each test case, an agent receives a natural language task description x_i along with clinical concept
220 descriptions a_i (e.g., value-set aligned concepts). The agent generates SQL $\hat{y}_i = f(x_i, a_i)$, which is
221 executed against the database alongside the gold SQL y_i^* .

222 Evaluation is primarily execution-based, supplemented by structural metrics computed from generated
223 and gold SQL queries.

224 **Model and decoding configuration.** We use gemini-2.5-pro via the Vertex AI GenAI API
225 (vertexai=True) using the Google genai client, with the model identifier (e.g., gemini-2.5-pro)
226 and Google Cloud project and location supplied at initialization.

227 For comparison experiments, we additionally evaluate gemini-2.5-flash via the same Vertex AI API,
228 gpt-5.2 served through an enterprise API gateway (Apigee) using an Azure OpenAI-compatible
229 interface, and medgemma-1.5-4b-it deployed locally via vLLM.

230 Decoding uses temperature=0.0 for API-based models to encourage deterministic SQL generation;
231 other decoding hyperparameters (e.g., top- p , top- k , maximum output tokens) are not explicitly
232 set and therefore follow API defaults. For MedGemma, decoding uses temperature=0.2 and
233 max_tokens=2048, with other parameters following vLLM defaults.

234 For agent control, we allow up to five function-calling rounds per query. Transient API failures are
235 handled with up to two retries and a five-second delay between retry attempts.

236 **Compute setup.** API-based models are accessed through external inference APIs from Google
237 Cloud Workbench notebooks and are not executed locally.

238 Medgemma-1.5-4b-it is deployed locally on a Google Cloud VM with a g2-standard-16 configuration
239 (16 vCPUs, 64 GB RAM) and a single NVIDIA L4 GPU (24 GB VRAM). Each dataset evaluation
240 requires on the order of hours, depending on the number of queries.

241 5.3 Per-test-case metrics

242 **Code F1.** We compute F1 over extracted clinical codes (e.g., ICD-10) to measure correctness of
243 inclusion criteria.

244 **Table F1.** We compute F1 over referenced tables to evaluate schema selection.

245 **Column Recall.** Recall over expected output columns: $|\text{Col}_{\text{agent}} \cap \text{Col}_{\text{gold}}|/|\text{Col}_{\text{gold}}|$.

246 5.4 Composite score

247 We focus on the metrics with signal in the uploaded summary: Composite score, Code F1, Table F1,
248 and Column Recall. The exported patient-level, row-count, and null-quality metrics were uniformly
249 zero across all evaluated settings, and are therefore excluded from the analysis. Composite score
250 serves as the main end-to-end benchmark objective. We place weights to reflect the relative importance
251 of cohort identification (code) along with correctly selected tables and columns in the SQL query.
252 We do not normalize weights to sum to 1; the composite score should be interpreted as a weighted
253 utility rather than a probability-like metric. We plan to incorporate additional downstream evaluation
254 dimensions—such as patient counts returned, null-value handling, and column datatype validation—
255 that are currently excluded due to the inability to execute queries on the requisite databases. Metrics
256 are aggregated across test cases using mean, min, and max statistics.

$$\text{Score} = 0.20 \cdot \text{Code F1} + 0.15 \cdot \text{Table F1} + 0.15 \cdot \text{Column Recall}$$

257 **Comparison Protocol** We compare *with query guidance* against *without query guidance but access*
 258 *to schemas*. The CSV includes multiple runs/models per dataset; our tables report averages (across 3
 259 runs) and models per dataset. This comparison isolates the value of *guidance as managed context*.
 260 The baseline is not deprived of schema access entirely; instead, it has schema access without the
 261 curated guidance layer.

262 More precisely, the with-guidance condition pre-loads the global `query_guidance_index` and the
 263 `MASTER_DESCRIPTIONS.json` summary into the agent’s initial system prompt, while the without-
 264 guidance condition omits these. Both conditions retain access to the same `load_table_schemas`
 265 tool, which returns per-table `query_guidance` fields when present in the JSON artifacts. The
 266 comparison therefore isolates the value of *eager* guidance—guidance resident in the initial context
 267 window—relative to *lazy*, on-demand retrieval through tool calls (Appendix D.2). Improvements
 268 therefore support the claim that *how* and *when* schema information is represented matters, not only
 269 whether it is present.

270 6 Results

271 6.1 Guidance Improves Aggregate Performance

272 Table 1 reports aggregate results across all dataset–model pairs, averaged over three runs per configu-
 273 ration, with clinical datasets reported separately from the SPIDER dataset.

274 Across clinical datasets, schema guidance improves Composite score from 0.220 ± 0.017 to $0.296 \pm$
 275 0.018 , Code F1 from 0.691 ± 0.036 to 0.880 ± 0.016 , Table F1 from 0.290 ± 0.055 to 0.451 ± 0.064 ,
 276 and Column Recall from 0.253 ± 0.029 to 0.351 ± 0.039 . The largest relative improvement is
 277 observed in Table F1 (+55.5%), consistent with the search-space reduction view that guidance helps
 278 the model select the correct structural backbone before fine-grained decoding. We provide additional
 279 results about failure mode analysis in the appendixE

Setting	Composite	Code F1	Table F1	Column Recall
<i>Clinical datasets</i>				
With guidance	0.296 ± 0.018	0.880 ± 0.016	0.451 ± 0.064	0.351 ± 0.039
Without guidance	0.220 ± 0.017	0.691 ± 0.036	0.290 ± 0.055	0.253 ± 0.029
<i>SPIDER</i>				
With guidance	0.178 ± 0.012	0.651 ± 0.040	0.258 ± 0.022	0.062 ± 0.003
Without guidance	0.056 ± 0.008	0.215 ± 0.029	0.062 ± 0.014	0.027 ± 0.004

Table 1: Aggregate results (mean \pm standard error) reported separately for clinical datasets and the SPIDER dataset. Results are averaged over dataset–model pairs and three runs per configuration.

280 6.2 Consistent Gains Across Datasets

281 Table 2 reports results by dataset. The Composite score improves under guidance on every dataset:
 282 from 0.133 ± 0.007 to 0.180 ± 0.005 on the closed institutional benchmark, from 0.243 ± 0.044
 283 to 0.419 ± 0.005 on EHRSQL-eICU, from 0.289 ± 0.013 to 0.320 ± 0.012 on EHRSQL-MIMIC-
 284 III, and from 0.056 ± 0.008 to 0.178 ± 0.012 on SPIDER. Code F1 likewise improves across all
 285 datasets: $0.591 \pm 0.034 \rightarrow 0.799 \pm 0.018$ (closed), $0.584 \pm 0.088 \rightarrow 0.961 \pm 0.011$ (eICU),
 286 $0.870 \pm 0.012 \rightarrow 0.901 \pm 0.025$ (MIMIC-III), and $0.215 \pm 0.029 \rightarrow 0.651 \pm 0.040$ (SPIDER).

287 The largest relative gains are observed on SPIDER, where the Composite score increases more
 288 than $3\times$ and Table F1 rises from 0.062 ± 0.014 to 0.258 ± 0.022 , and on EHRSQL-eICU, where
 289 Table F1 improves from 0.453 ± 0.105 to 0.839 ± 0.016 and Column Recall increases substantially
 290 ($0.389 \pm 0.071 \rightarrow 0.677 \pm 0.014$). In contrast, EHRSQL-MIMIC-III shows the smallest absolute
 291 gain, consistent with its stronger unguided baseline.

292 6.3 Significance of Improvement in Structural Metrics

293 The concentration of gains in structural metrics (Code F1 and Table F1) supports the search-space re-
 294 duction hypothesis: guidance is most effective when it suppresses structurally implausible hypotheses.

Dataset	N	Comp. (G)	Comp. (No G)	Code F1 (G)	Code F1 (No G)	Table F1 (G)	Table F1 (No G)	Column Recall (G)	Column Recall (No G)
Closed institutional benchmark	554	0.180 ± 0.005	0.133 ± 0.007	0.799 ± 0.018	0.591 ± 0.034	0.007 ± 0.001	0.004 ± 0.000	0.129 ± 0.012	0.094 ± 0.006
EHRSQL-eICU	1204	0.419 ± 0.005	0.243 ± 0.044	0.961 ± 0.011	0.584 ± 0.088	0.839 ± 0.016	0.453 ± 0.105	0.677 ± 0.014	0.389 ± 0.071
EHRSQL-MIMIC-III	1198	0.320 ± 0.012	0.289 ± 0.013	0.901 ± 0.025	0.870 ± 0.012	0.606 ± 0.043	0.454 ± 0.081	0.328 ± 0.004	0.310 ± 0.006
SPIDER	3509	0.178 ± 0.012	0.056 ± 0.008	0.651 ± 0.040	0.215 ± 0.029	0.258 ± 0.022	0.062 ± 0.014	0.062 ± 0.003	0.027 ± 0.004

Table 2: Per-dataset results (mean \pm standard error) averaged over three runs and models available in both guidance conditions. Clinical datasets are reported alongside the SPIDER dataset benchmark for comparison.

Model	Comp. (G)	Comp. (No G)	Code F1 (G)	Code F1 (No G)	Table F1 (G)	Table F1 (No G)	Column Recall (G)	Column Recall (No G)
gemini-2.5-flash	0.280 ± 0.029	0.170 ± 0.032	0.838 ± 0.033	0.491 ± 0.076	0.443 ± 0.095	0.277 ± 0.076	0.306 ± 0.070	0.198 ± 0.045
gemini-2.5-pro	0.301 ± 0.031	0.244 ± 0.038	0.897 ± 0.026	0.712 ± 0.072	0.493 ± 0.107	0.403 ± 0.105	0.320 ± 0.077	0.276 ± 0.069
gpt-5.2	0.267 ± 0.033	0.140 ± 0.029	0.810 ± 0.056	0.470 ± 0.079	0.406 ± 0.098	0.163 ± 0.068	0.293 ± 0.066	0.144 ± 0.031
medgemma-1.5-4b-it	0.201 ± 0.023	0.185 ± 0.016	0.729 ± 0.013	0.776 ± 0.035	0.186 ± 0.081	0.001 ± 0.001	0.183 ± 0.055	0.197 ± 0.060

Table 3: Performance comparison across models (mean \pm standard error), averaged over datasets where both guidance conditions are available and over three runs per configuration.

295 In clinical SQL, errors in table selection are often fatal to downstream interpretation, and guidance
 296 directly targets this failure mode.

297 One caveat is the closed institutional benchmark, where Table F1 remains near zero in both conditions
 298 (0.007 ± 0.001 with guidance, 0.004 ± 0.000 without). This likely reflects the use of views, CTEs,
 299 and schema-qualified or aliased table names in the reference queries, which are not fully captured
 300 by the table-name matching metric. Furthermore, complex queries sometimes integrated various
 301 databases where the schemas are separate from the current schemas used in this work. We therefore
 302 report Table F1 for completeness but rely primarily on Code F1 and Composite score as more reliable
 303 indicators on this dataset.

304 6.4 Relevance to Clinical Applications

305 The closed institutional benchmark results are especially important because they suggest the method
 306 is not limited to public benchmark schemas. Real clinical warehouses often contain many tables,
 307 multiple near-synonymous dimensions, and institutional conventions that are not obvious from
 308 schema names alone. Guidance provides a place to encode those conventions explicitly. Precisely
 309 this type of knowledge is difficult to recover from one-shot prompting but well suited to a structured,
 310 updateable context object.

311 7 Related Work

312 **Large Language Model Text-to-SQL Systems** Modern text-to-SQL systems increasingly rely
 313 on large language models, with a range of enhancements improving performance across domains.
 314 Initial work focused on single LLM enhancements with prompt-level methods such as in-context
 315 learning Sun et al. [2023], Zhang et al. [2023], Pourreza and Rafiei [2023] and prompt engineering
 316 Pandey et al. [2025], as well as retrieval-augmented pipelines [Shen et al., 2024] and reinforcement-
 317 learning-enhanced retrieval [Toteja et al., 2025]. Other approaches directly modify model weights
 318 through fine-tuning Chafik et al. [2026], Roberson et al. [2024], augment capabilities by incorporating
 319 expert knowledge Hong et al. [2024] or synthetic data Caferoğlu et al. [2025], or seek to introduce
 320 alternative query representations Eyal et al. [2023].

321 With the increasing capability of agentic systems, system-level approaches are emerging that incorpo-
 322 rate multiple LLMs via multi-step planning, with or without agentic control or reinforcement learning
 323 Chen et al. [2024], Pham et al. [2026], Xu et al. [2025], Wang et al. [2025], Shao et al. [2025], or
 324 leverage ensemble methods Shen et al. [2023].

325 8 Conclusion

326 Our results suggest that agentic SQL systems benefit from treating context as something to be
 327 *constructed* rather than only retrieved. In the clinical setting, raw schemas exceed practical context
 328 windows and a structured guidance object that is explicit, inspectable, and budget-feasible offers a

329 usable middle ground between flat retrieval and opaque latent memory. The empirical evidence here
330 is limited to the clinical databases we evaluate and the specific construction in Section 3.3.

331 **Limitations** Our empirical evaluation is based on exported benchmark summaries rather than
332 raw per-example traces. This restricts the depth of error analysis and ablation we can provide
333 since we do not provide it at the query level. Future iterations should include run-level paired
334 comparisons and error analysis. While we frame schema guidance as the primary contribution,
335 the system is to be implemented within a broader agentic pipeline that includes planning, value
336 resolution, validation, and repair. Improvements may arise in part from system-level redundancy,
337 error correction, or architectural choices rather than the guidance representation alone. A key issue
338 is the analysis primarily on the SQL queries, and not the returned data; the system may produce
339 syntactically and structurally plausible SQL that is nonetheless clinically incorrect, incomplete, or
340 misleading. This emphasis on structural metrics (e.g. Code F1 and Table F1) without assessment
341 of clinically meaningful correctness of the data returned is primarily due issues with accessing
342 identifiable data for the closed institutional dataset. However, we have built in downstream metrics
343 (e.g., patient-level accuracy, row-count agreement, null-quality) with future plans to use these metrics
344 to assess real-world utility. Finally, the approach is evaluated primarily on structured relational
345 databases with relatively well-defined schemas. Its applicability to noisier real-world settings—such
346 as partially documented schemas, missing or implicit foreign keys, heavy reliance on free-text fields,
347 or hybrid data models (e.g., FHIR, event streams)—is not established. More broadly, the system
348 does not address key deployment concerns in clinical environments, including uncertainty estimation,
349 interpretability, auditability, and human oversight. These limitations are critical for real-world
350 adoption and remain open directions for future work.

351 **Broader Impact** This work lays a foundation for iterating from static guidance construction
352 toward a future of *interactive, clinician-in-the-loop context refinement*. Schema guidance is naturally
353 interpretable and editable, which creates an opportunity for domain experts to directly inspect and
354 revise the context that governs generation. Incorporating such feedback into the guidance optimization
355 loop would allow the system to iteratively refine its representation of clinically relevant structure,
356 transforming guidance from a fixed, dataset-level artifact into a continuously evolving knowledge
357 layer. We anticipate integration of this into an agentic harness that incorporates deterministic
358 validators that check structural properties, a SQL fixer that repairs failed outputs; and an execution-
359 aware gate using dry-run validation. These can be further extended beyond single-query generation
360 into time-varying *longitudinal phenotyping* across clinical databases. Finally, we envision applications
361 to downstream tasks in clinical research and care, particularly *deep phenotyping and outcome*
362 *ascertainment*. By enabling more reliable and interpretable cohort definitions, schema guidance
363 could support the construction of high-quality observational datasets and clinically meaningful
364 endpoints. More broadly, this work points toward agentic systems that do not merely generate
365 queries, but maintain and refine structured representations of domain knowledge that can be reused
366 across tasks, datasets, and institutions. Despite these benefits, there are important risks to consider.
367 Incorrect or incomplete SQL generation may lead to misleading clinical analyses if used without
368 expert validation, particularly in high-stakes settings. Misuse of such systems for automated cohort
369 construction without appropriate oversight could introduce bias or propagate erroneous assumptions.
370 In addition, applying schema guidance to sensitive clinical data raises privacy and governance
371 concerns. Mitigation strategies include maintaining human-in-the-loop review, enforcing validation
372 checks prior to execution, and restricting deployment to controlled environments with appropriate
373 data governance.

374 **References**

- 375 Hasan Alp Caferoğlu, Mehmet Serhat Çelik, and Özgür Ulusoy. Sing-sql: A synthetic data generation
376 framework for in-domain text-to-sql translation, 2025. URL [https://arxiv.org/abs/2509.](https://arxiv.org/abs/2509.25672)
377 25672.
- 378 Salmane Chafik, Saad Ezzini, and Ismail Berrada. Lego-code: Can modular curriculum learning
379 advance complex code generation? insights from text-to-sql, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2604.18254)
380 [abs/2604.18254](https://arxiv.org/abs/2604.18254).
- 381 Ziru Chen, Michael White, Raymond Mooney, Ali Payani, Yu Su, and Huan Sun. When is tree search
382 useful for llm planning? it depends on the discriminator, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.10890)
383 [2402.10890](https://arxiv.org/abs/2402.10890).
- 384 Ben Eyal, Amir Bachar, Ophir Haroche, and Michael Elhadad. Semantic parsing for complex data
385 retrieval: Targeting query plans vs. sql for no-code access to relational databases, 2023. URL
386 <https://arxiv.org/abs/2312.14798>.
- 387 Satya K Gundabathula and Sriram R Kolar. Promptmind team at ehsql-2024: Improving reliability
388 of sql generation using ensemble llms, 2024. URL <https://arxiv.org/abs/2405.08839>.
- 389 Zijin Hong, Zheng Yuan, Hao Chen, Qinggang Zhang, Feiran Huang, and Xiao Huang. Knowledge-
390 to-sql: Enhancing sql generation with data expert llm, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.11517)
391 [2402.11517](https://arxiv.org/abs/2402.11517).
- 392 Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li-Wei H. Lehman, Mengling Feng, Mohammad
393 Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. Mimic-iii, a
394 freely accessible critical care database. *Scientific Data*, 3:160035, 2016. doi: 10.1038/sdata.2016.
395 35.
- 396 B. Ross Katz, Abdul Khan, James York-Winegar, and Alexander J. Titus. Nhanes-gcp: Leveraging
397 the google cloud platform and bigquery ml for reproducible machine learning with data from the
398 national health and nutrition examination survey, 2024. URL [https://arxiv.org/abs/2401.](https://arxiv.org/abs/2401.06967)
399 [06967](https://arxiv.org/abs/2401.06967).
- 400 Gyubok Lee, Sunjun Kweon, Seongsu Bae, and Edward Choi. Overview of the EHRSQL 2024
401 shared task on reliable text-to-sql modeling on electronic health records. In *Proceedings of the*
402 *6th Clinical Natural Language Processing Workshop*, pages 644–654, Mexico City, Mexico, jun
403 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.clinicalnlp-1.62. URL
404 <https://aclanthology.org/2024.clinicalnlp-1.62/>.
- 405 Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
406 and Percy Liang. Lost in the middle: How language models use long contexts, 2023. URL
407 <https://arxiv.org/abs/2307.03172>.
- 408 Ali Modarressi, Hanieh Deilamsalehy, Franck Deroncourt, Trung Bui, Ryan A. Rossi, Seunghyun
409 Yoon, and Hinrich Schütze. Nolima: Long-context evaluation beyond literal matching, 2025. URL
410 <https://arxiv.org/abs/2502.05167>.
- 411 Hari Mohan Pandey, Anshul Gupta, Subham Sarkar, Minakshi Tomer, Schneider Johannes, and
412 Yan Gong. Gemma-sql: A novel text-to-sql model based on large language models, 2025. URL
413 <https://arxiv.org/abs/2511.04710>.
- 414 Minh Tam Pham, Trinh Pham, Tong Chen, Hongzhi Yin, Quoc Viet Hung Nguyen, and Thanh Tam
415 Nguyen. Av-sql: Decomposing complex text-to-sql queries with agentic views, 2026. URL
416 <https://arxiv.org/abs/2604.07041>.
- 417 Tom J. Pollard, Alistair E. W. Johnson, Joseph D. Raffa, Leo A. Celi, Roger G. Mark, and Omar
418 Badawi. The eicu collaborative research database, a freely available multi-center database for
419 critical care research. *Scientific Data*, 5:180178, 2018. doi: 10.1038/sdata.2018.178.
- 420 Mohammadreza Pourreza and Davood Rafiei. Din-sql: Decomposed in-context learning of text-to-sql
421 with self-correction, 2023. URL <https://arxiv.org/abs/2304.11015>.

- 422 Richard Roberson, Gowtham Kaki, and Ashutosh Trivedi. Analyzing the effectiveness of large
423 language models on text-to-sql synthesis, 2024. URL <https://arxiv.org/abs/2401.12379>.
- 424 Zhihui Shao, Shubin Cai, Rongsheng Lin, and Zhong Ming. Enhancing text-to-sql with question
425 classification and multi-agent collaboration. In *Findings of the Association for Computational*
426 *Linguistics: NAACL 2025*, pages 4340–4349, Albuquerque, New Mexico, apr 2025. Association
427 for Computational Linguistics. doi: 10.18653/v1/2025.findings-naacl.245. URL [https://](https://aclanthology.org/2025.findings-naacl.245/)
428 aclanthology.org/2025.findings-naacl.245/.
- 429 Ran Shen, Gang Sun, Hao Shen, Yiling Li, Liangfeng Jin, and Han Jiang. Spsql: Step-by-step parsing
430 based framework for text-to-sql generation, 2023. URL <https://arxiv.org/abs/2305.11061>.
- 431 Zhili Shen, Pavlos Vougiouklis, Chenxin Diao, Kaustubh Vyas, Yuanyi Ji, and Jeff Z. Pan. Improving
432 retrieval-augmented text-to-sql with ast-based ranking and schema pruning. In *Proceedings of*
433 *the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7865–7879,
434 Miami, Florida, USA, nov 2024. Association for Computational Linguistics. doi: 10.18653/v1/
435 2024.emnlp-main.449. URL <https://aclanthology.org/2024.emnlp-main.449/>.
- 436 Ruoxi Sun, Sercan Ö. Arik, Rajarishi Sinha, Hootan Nakhost, Hanjun Dai, Pengcheng Yin, and
437 Tomas Pfister. Sqlprompt: In-context text-to-sql with minimal labeled data, 2023. URL [https://](https://arxiv.org/abs/2311.02883)
438 arxiv.org/abs/2311.02883.
- 439 Samuel Thio, Matthew Lewis, Spiros Denaxas, and Richard J. B. Dobson. Unlocking electronic
440 health records: a hybrid graph rag approach to safe clinical ai for patient qa. *Frontiers in*
441 *Digital Health*, 8, March 2026. ISSN 2673-253X. doi: 10.3389/fdgth.2026.1780700. URL
442 <http://dx.doi.org/10.3389/fdgth.2026.1780700>.
- 443 Rishit Toteja, Arindam Sarkar, and Prakash Mandayam Comar. In-context reinforcement learning with
444 retrieval-augmented generation for text-to-sql. In *Proceedings of the 31st International Conference*
445 *on Computational Linguistics*, pages 10390–10397, Abu Dhabi, UAE, jan 2025. Association for
446 Computational Linguistics. URL <https://aclanthology.org/2025.coling-main.692/>.
- 447 Bing Wang, Changyu Ren, Jian Yang, Xinnian Liang, Jiaqi Bai, LinZheng Chai, Zhao Yan, Qian-Wen
448 Zhang, Di Yin, Xing Sun, and Zhoujun Li. Mac-sql: A multi-agent collaborative framework for
449 text-to-sql. In *Proceedings of the 31st International Conference on Computational Linguistics*,
450 pages 540–557, Abu Dhabi, UAE, jan 2025. Association for Computational Linguistics. URL
451 <https://aclanthology.org/2025.coling-main.36/>.
- 452 Ping Wang, Tian Shi, and Chandan K. Reddy. Text-to-sql generation for question answering on
453 electronic medical records, 2020. URL <https://arxiv.org/abs/1908.01839>.
- 454 Zekun Xu, Siyu Xia, Chuhuai Yue, Jiajun Chai, Mingxue Tian, Xiaohan Wang, Wei Lin, Haoxuan
455 Li, and Guojun Yin. Mtir-sql: Multi-turn tool-integrated reasoning reinforcement learning for
456 text-to-sql, 2025. URL <https://arxiv.org/abs/2510.25510>.
- 457 Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene
458 Li, Qingning Yao, Shanell Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale
459 human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task, 2019.
460 URL <https://arxiv.org/abs/1809.08887>.
- 461 Hanchong Zhang, Ruisheng Cao, Lu Chen, Hongshen Xu, and Kai Yu. Act-sql: In-context learning
462 for text-to-sql with automatically-generated chain-of-thought, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2310.17342)
463 [abs/2310.17342](https://arxiv.org/abs/2310.17342).
- 464 Angelo Ziletti and Leonardo D’Ambrosi. Retrieval augmented text-to-sql generation for epi-
465 demiological question answering using electronic health records. In *Proceedings of the 6th*
466 *Clinical Natural Language Processing Workshop*, pages 47–53, Mexico City, Mexico, jun
467 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.clinicalnlp-1.4. URL
468 <https://aclanthology.org/2024.clinicalnlp-1.4/>.
- 469 Angelo Ziletti and Leonardo D’Ambrosi. Generating patient cohorts from electronic health records
470 using two-step retrieval-augmented text-to-sql generation, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2502.21107)
471 [abs/2502.21107](https://arxiv.org/abs/2502.21107).

472 **A Formal problem setting**

473 This appendix gives the formal version of the problem setting summarized in Section 2. Notation is
 474 consistent with the body; we restate the objects here so that the construction in Appendix B and the
 475 search-space view in Appendix B.1 can be read self-containedly.

476 **Inputs and generator.** Let $x \in \mathcal{X}$ denote a natural-language clinical query, $\mathcal{S} \in \mathfrak{S}$ a relational
 477 schema drawn from a space of schemas \mathfrak{S} , and y^* a target SQL program. A text-to-SQL system is a
 478 stochastic mapping

$$p_f(\hat{y} \mid x, c)$$

479 that produces SQL hypotheses conditioned on a context c . Two choices for c are of interest: the raw
 480 schema $c = \mathcal{S}$ and a derived guidance object $c = g$. Hypotheses are sampled as

$$\hat{y} \sim \text{Decode}(p_f(\cdot \mid x, c)),$$

481 where $\text{Decode}(\cdot)$ is a fixed decoding rule (e.g., greedy, beam, or best-of- k with execution feedback)
 482 held constant across all conditions compared in the paper. The generator f is also held fixed: all
 483 comparisons reported below vary c only.

484 **Definition A.1** (Schema-guidance policy). *Let Θ be a space of database-specific configuration*
 485 *parameters (e.g., key/time naming heuristics, slowly-changing-dimension flag handling, sensitivity*
 486 *rules). A schema-guidance policy is a parameterized mapping*

$$\pi_g : \mathfrak{S} \times \Theta \longrightarrow \mathcal{G}_B, \quad g = \pi_g(\mathcal{S}; \theta),$$

487 *producing a structured context object containing key filters, commonly joined tables, field notes,*
 488 *example query patterns, and critical notes.*

489 We treat $\pi_g(\cdot; \theta)$ as a *procedure* that is invariant across databases: only the inputs \mathcal{S} and the
 490 configuration θ vary, while the construction itself is shared. The pipeline realizing π_g is specified
 491 in Appendix B, and its per-database instantiation is written $g^{(d)} = \pi_g(\mathcal{S}^{(d)}; \theta^{(d)})$. We refer to this
 492 property as *construction invariance* and treat it as our operational notion of *database-agnostic*; we
 493 make no formal categorical claim.

494 **Unguided baseline** The unguided baseline π_{g_0} is defined as raw schema text serialized and truncated
 495 to fit $|g| \leq B$:

$$\pi_{g_0}(\mathcal{S}; \theta) = \text{Trunc}_B(\mathcal{S}),$$

496 where Trunc_B is a deterministic, schema-independent serialization-and-truncation rule. This isolates
 497 the value of *managed* context: π_g and π_{g_0} consume the same budget B and the same generator f ,
 498 and differ only in how schema information is represented.

499 **Empirical objective** We evaluate guidance policies on a held-out validation set

$$\mathcal{D}_{\text{val}} = \{(x_i, \mathcal{S}_i, y_i^*)\}_{i=1}^n.$$

500 Let $u(x, y^*, \hat{y}) \in [0, 1]$ be a bounded utility (e.g., execution match, Code F1, Table F1) with
 501 $u(x, y^*, y^*) = 1$. The empirical utility of a policy π_g under fixed θ is

$$\hat{J}(\pi_g) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\hat{y}_i \sim \text{Decode}(p_f(\cdot \mid x_i, \pi_g(\mathcal{S}_i; \theta)))} [u(x_i, y_i^*, \hat{y}_i)].$$

502 For each database d we report $\hat{J}(\pi_g)$ against $\hat{J}(\pi_{g_0})$, and the improvement $\Delta(d) = \hat{J}(\pi_g^{(d)}) - \hat{J}(\pi_{g_0}^{(d)})$
 503 quantifies the value of guidance on that database. Construction invariance is then operationalized
 504 as the requirement that a single procedure π_g yields $\Delta(d) > 0$ across the held-out databases in
 505 Section 4.

506 **What is held fixed** For all comparisons in the body, the following are held constant across π_g and
 507 π_{g_0} : the generator f , the decoding rule $\text{Decode}(\cdot)$, the validation set \mathcal{D}_{val} , and the utility u . The only
 508 quantity that varies is the context $c \in \{\pi_g(\mathcal{S}; \theta), \pi_{g_0}(\mathcal{S}; \theta)\}$ delivered to f at inference time.

509 **B Detailed construction pipeline**

510 We construct $g^{(d)} = \pi_g(\mathcal{S}^{(d)}; \theta^{(d)})$ via a pipeline whose shape is fixed across datasets and whose
 511 only per-database inputs are the raw schema $\mathcal{S}^{(d)}$ and the configuration $\theta^{(d)}$. A single approximation
 512 guarantee is invoked for the summarization step; all other stages are evaluated empirically.

513 **Stage 1: Input normalization.** For each database d , we normalize raw schema metadata into a
 514 canonical representation

$$\mathcal{S}^{(d)} = \{T_1, \dots, T_m\}, \quad T_j = \{(C_{j\ell}, \tau_{j\ell}, \delta_{j\ell})\}_{\ell=1}^{n_j},$$

515 where $C_{j\ell}$ is a column name, $\tau_{j\ell}$ its data type, and $\delta_{j\ell}$ an optional textual description. Foreign-key
 516 relationships $\mathcal{R} \subseteq (\cup_j T_j)$ are extracted from declarations when available and inferred from naming
 517 conventions otherwise. The output is a unified per-database JSON artifact with identical shape across
 518 all d .

519 **Stage 2: Schema graph and join neighborhoods.** We construct a schema graph $G^{(d)} = (V, E)$
 520 with $V = \{T_1, \dots, T_m\}$ and edges induced by \mathcal{R} . Per-table guidance is generated over a one-hop
 521 join neighborhood

$$\mathcal{N}(T_j) = \{T_k \in V : (T_j, T_k) \in E\},$$

522 so that each table’s guidance is conditioned on the tables it can be joined with directly. Multi-hop
 523 join structure is recovered indirectly via overlapping neighborhoods at the summarization stage; we
 524 report the effect of this locality choice.

525 **Stage 3: Per-table guidance generation.** We apply a per-table operator

$$\phi : (T_j, \mathcal{N}(T_j); \theta^{(d)}) \mapsto g_j \in \mathcal{G}_{\text{tab}},$$

526 producing a structured artifact g_j with slots for indexed columns, performance filters, common joins,
 527 field notes, and example query patterns. The operator ϕ is the composition of three sub-operators
 528 applied in sequence:

- 529 1. **Heuristic skeleton.** Index, filter, and join slots are populated by deterministic rules over
 530 column names, types, and declared relationships (e.g., key/time naming patterns, shared-
 531 column join inference). This stage is content-free regarding clinical semantics and depends
 532 only on $\theta^{(d)}$.
- 533 2. **LLM enrichment.** A language model conditioned on $T_j, \mathcal{N}(T_j)$, a global relationship
 534 summary, and any domain hints in $\theta^{(d)}$ proposes natural-language descriptions, recommen-
 535 dations, and example query patterns for the qualitative slots left empty by Stage 3.1.
- 536 3. **Curation merge.** When subject-matter-expert annotations exist for $\mathcal{S}^{(d)}$, they are merged
 537 into g_j with precedence over the LLM-generated content.

538 **B.1 Search-space reduction view**

539 Let $\mathcal{H}(\mathcal{S})$ denote the hypothesis class of SQL programs compatible with the full schema, and let
 540 $\mathcal{H}_g(\mathcal{S}, x) \subseteq \mathcal{H}(\mathcal{S})$ denote the subset that the generator effectively explores at inference time when
 541 conditioned on the guidance object g . Guidance is intended to bias the generator away from irrelevant
 542 tables, implausible joins, missing filters, and known anti-patterns, so \mathcal{H}_g is typically strictly smaller
 543 than the unguided support.

544 We adopt a deliberately abstract model of decoding error: let $\varepsilon(m)$ be a non-decreasing function
 545 of the size m of the candidate program set. This captures the common situation in which decoding
 546 becomes harder as more spurious hypotheses remain available; we treat this monotonicity as a
 547 working assumption about LLM decoders, not as a theorem.

548 **Remark B.1** (Search-space reduction under recall-preserving guidance). *Suppose, for an example*
 549 *(x, \mathcal{S}, y^*) , that (A1) $y^* \in \mathcal{H}_g(\mathcal{S}, x)$ (recall-preserving), (A2) $|\mathcal{H}_g(\mathcal{S}, x)| < |\mathcal{H}(\mathcal{S}, x)|$ (precision-*
 550 *improving), and (A3) the decoder error is upper bounded by a non-decreasing function $\varepsilon(\cdot)$ of*
 551 *candidate-set size. Then*

$$\Pr[\hat{y} \neq y^* | g] \leq \varepsilon(|\mathcal{H}_g(\mathcal{S}, x)|) \leq \varepsilon(|\mathcal{H}(\mathcal{S}, x)|).$$

552 We treat Remark B.1 as a *decomposition* rather than a substantive result: the inequality is an immediate
553 consequence of (A1)–(A3), and the substantive content lies in whether those assumptions hold for
554 our construction. (A3) is a working assumption about LLM decoders that we adopt to organize
555 the decomposition. Together, (A1) and (A2) motivate two distinct empirical questions—*recall* of
556 guidance (does g contain the schema elements required by y^* ?) and *precision* of guidance (does
557 g suppress structurally implausible alternatives?)—and predict that gains should concentrate in
558 structural metrics (Code F1, Table F1) when both hold. Section 4 reports the structural-metric pattern
559 predicted by this decomposition; a direct empirical decomposition of guidance recall vs. precision is
560 left to future work with run-level traces.

561 C Additional Dataset details

562 **Licensing and data usage** MIMIC-III and eICU are publicly available de-identified clinical
563 datasets accessed under credentialed PhysioNet data use agreements and used in accordance with
564 their respective terms. The SPIDER and related text-to-SQL benchmark datasets are used under
565 their original research licenses. The EHRSQL dataset is released under the CC-BY-4.0 license. The
566 closed institutional dataset is not publicly released due to privacy and data governance constraints.
567 All datasets are used strictly for research purposes and in compliance with their respective licenses
568 and terms of use.

569 C.1 SPIDER

570 Each data point includes (i) a database identifier (`db_id`), (ii) a ground-truth SQL query, (iii) extracted
571 constants and their associated schema-level context, (iv) a set of executable test databases for
572 denotation-based evaluation, and (v) the corresponding natural language utterance. This unified
573 format facilitates consistent benchmarking across heterogeneous datasets and reduces variability
574 caused by dataset-specific SQL styles. The SPIDER dataset is released under the Apache-2.0 license
575 and is used in compliance with its terms.

576 C.2 MIMIC-III

577 MIMIC-III is a publicly available, de-identified critical care database developed by the MIT Labora-
578 tory for Computational Physiology and collaborators. It contains detailed electronic health records
579 associated with intensive care unit (ICU) admissions at Beth Israel Deaconess Medical Center, includ-
580 ing demographics, laboratory measurements, medications, charted observations, and clinical events
581 [Johnson et al., 2016]. The database covers adult and neonatal ICU stays collected between 2001 and
582 2012, and has become a standard resource for clinical machine learning and medical NLP research.

583 C.3 eICU Collaborative Database

584 The eICU Collaborative Research Database is a large-scale, de-identified multi-center ICU database
585 collected from hospitals across the United States [Pollard et al., 2018]. Unlike MIMIC-III, which
586 is derived from a single medical center, eICU provides broader institutional diversity, including
587 variations in documentation practices, patient populations, and treatment patterns. The database
588 contains high-granularity information such as vital signs, diagnoses, treatment records, and care
589 documentation.

590 C.4 Closed Institutional Benchmark Dataset

591 The closed institutional dataset is not publicly released due to privacy and data governance constraints
592 and is used in accordance with institutional data use and privacy policies. All data are de-identified or
593 handled under appropriate governance procedures, with Institutional Review Board (IRB) approval
594 number 26-002996.

595 **D Runtime agent and prompt**

596 **D.1 Runtime agent implementation**

597 At inference time, we instantiate an agentic SQL generator that performs structured text-to-SQL
598 generation under schema constraints. Unless stated otherwise, the configuration includes the full
599 `query_guidance_index` summarized in (ii) below; the without-query-guidance baseline uses the
600 same agent and pipeline but omits that pre-loaded index, as detailed when we describe the simplified
601 prompt template. The agent consumes schema guidance as its primary context and maintains a
602 dynamically constructed system prompt that integrates (i) compact schema descriptions, (ii) a query
603 guidance index summarizing indexed columns, key filters, and join patterns (omitted in the baseline),
604 and (iii) domain-specific hints for mapping user queries to relevant schema components.

605 Given a natural language query x , the language model is prompted to decide whether additional
606 schema detail is needed, using the task text together with high-level guidance (e.g., the query-
607 guidance summary and domain hints). When additional detail is required, it emits a structured
608 call to `load_table_schemas` with tables identified in `SCHEMA.TABLE` form. The tool returns table
609 definitions drawn from the curated JSON schema artifacts: column names, data types, and—when
610 present in the file—optional `query_guidance` fields (indexed columns, performance-oriented filters,
611 and common join patterns). That payload is serialized into the dialogue and supplied to the next
612 decoding pass together with the standing system instructions and prior turns.

613 SQL generation follows two design rules communicated to the model: (i) in executable SQL, only
614 simple table names are used (no dataset or schema prefixes), and (ii) column identifiers are restricted
615 to those present in the loaded schema JSON (including outputs returned by `load_table_schemas`),
616 rather than being invented or copied from unrelated examples.

617 These rules enforce schema-grounded generation through prompting, rather than via an external
618 syntactic validator. The model is instructed to return structured JSON with fields such as the SQL
619 string, a short explanation, tables used, assumptions, and optional notes.

620 In the minimal inference configuration evaluated here, `load_table_schemas` is the only callable
621 tool for retrieving detailed column definitions from the JSON artifacts; optional routines that sample
622 or summarize rows in a live database, when enabled, belong to an extended SQLite-connected variant
623 and are outside the core schema-only agent.

624 The agent supports multi-turn behavior with a fixed cap on function-calling rounds. Within a single
625 user query, previously loaded table definitions are cached so repeated `load_table_schemas` requests
626 for the same file avoid redundant I/O. Across independent benchmark instances, our evaluation
627 harness calls an explicit reset after each test case, clearing conversation history and the schema cache
628 so that no context carries over between examples. The same agent object can be reused without reset
629 for multi-turn interaction, in which case prior dialogue and cached schema files persist; this behavior
630 is determined by the caller (e.g., whether the benchmark adapter invokes reset), not by a separate
631 named configuration flag.

632 **D.2 Simplified runtime prompt template**

633 The following is a simplified version of the dynamically generated runtime prompt used in both con-
634 ditions. In the with-guidance condition, `<<<DYNAMIC_QUERY_GUIDANCE_SUMMARY>>>` is populated
635 from `query_guidance_index`. In the schema-only condition, the agent is initialized directly
636 from schema JSON files without `MASTER.DESCRPTIONS.json`, leaving `query_guidance_index` empty
637 and removing all pre-loaded guidance summaries from the system prompt. As a result, the model
638 receives only raw schema information in its initial context. Any per-table `query_guidance` fields, if
639 present in the schema artifacts, are only accessible after explicit calls to `load_table_schemas` and
640 are therefore not part of the initial prompt context.

641 You are an expert SQL query generator for a SQLite healthcare database.

642

643 [SQLite Configuration]

644 This database uses SQLite. All SQL queries must follow these rules:

- 645 - Use simple table names (e.g., `FROM FACT_DIAGNOSIS`, `FROM DIM_PATIENT`)
- 646 - Do not use schema prefixes or BigQuery-style paths

```

647 - Use double quotes for identifiers if necessary
648
649 [Naming Conventions]
650 There are two distinct conventions:
651 1. Schema loading (for load_table_schemas):
652     - Use SCHEMA.TABLE format
653     - Example: "FACT_CLINICAL_DOCS_MCC.FACT_CLINICAL_DOCUMENTS"
654 2. SQL generation:
655     - Use simple table names only (no schema prefixes)
656
657 [Schema Context]
658 <<<DYNAMIC_JSON_SCHEMA_DESCRIPTIONS>>>
659 <<<DYNAMIC_QUERY_GUIDANCE_SUMMARY>>> # empty in schema-only
660 <<<DYNAMIC_SCHEMA_NOTE>>>
661 [Instructions]
662
663
664 1. Query Guidance Usage:
665     - If query guidance is provided, review it before loading schemas
666     - If query guidance is absent, select tables using schema descriptions
667       and domain hints only
668
669 2. Schema Loading:
670     - Use load_table_schemas when detailed schema information is required
671     - Always use SCHEMA.TABLE format when loading schemas
672
673 3. SQL Generation:
674     - Use simple table names (e.g., FACT_DIAGNOSIS)
675     - Include appropriate JOIN conditions
676     - Add LIMIT clauses (e.g., LIMIT 100)
677     - Use ROW_CURRENT_INDICATOR = 'Y' where applicable
678
679 4. Column Validation:
680     - Only use columns that exist in the loaded schemas
681     - Do not infer or assume column names
682     - Verify all columns in SELECT, WHERE, JOIN, and ORDER BY clauses
683
684 5. Response Format:
685
686 For SQL queries:
687 {
688     "sql_query": "...",
689     "explanation": "...",
690     "tables_used": [...],
691     "assumptions": [...],
692     "notes": [...]
693 }
694
695 For conversational responses:
696 {
697     "response": "..."
698 }
699
700 6. Follow-up Handling:
701     - Maintain conversation context
702     - Reuse previously loaded schemas when possible
703
704 7. Example Workflow:
705     Step 1: Identify relevant tables

```

```
706     Step 2: Load schemas using SCHEMA.TABLE format
707     Step 3: Generate SQL query using SQLite syntax
708
709 Example:
710 SELECT *
711 FROM FACT_CLINICAL_DOCUMENTS f
712 JOIN FACT_DIAGNOSIS d ON f.patient_dk = d.patient_dk
```

713 **E Failure Mode Analysis and Demonstrated Benefit of Query Guidance**

714 Across public datasets without SME curation (eICU, MIMIC-III, SPIDER), guidance consistently
 715 reduces non-SQL errors and improves exact match—primarily by decreasing hallucinated tables and
 716 columns—demonstrating that gains arise from improved task framing rather than reliance on curated
 717 annotations, even as incorrect table selection remains the dominant residual failure mode.

Model	No SQL ↓	Wrong Table ↓	Exact Match ↑	Code F1 ↑
Closed Institutional Benchmark				
gpt-5.2 (No G)	~0.20	~0.60	~0.00	0.56
gpt-5.2 (G)	~0.05	~0.80–0.90	~0.01	0.82
gemini-2.5-flash (No G)	~0.40	~0.60	~0.00	0.42
gemini-2.5-flash (G)	~0.05	~0.95+	~0.006	0.84
gemini-2.5-pro (No G)	~0.10	~0.90	~0.00	0.68
gemini-2.5-pro (G)	~0.01	~0.95+	~0.007	0.85
medgemma-1.5-4b-it (No G)	~0.15	~0.84	~0.00	0.70
medgemma-1.5-4b-it (G)	~0.18	~0.81	~0.00	0.70
EHRSQL-eICU				
gpt-5.2 (No G)	0.92–0.96	~0.30	~0.02	0.28
gpt-5.2 (G)	0.20–0.30	~0.30	0.50+	0.98
gemini-2.5-flash (No G)	0.55–0.65	~0.25	~0.05	0.56
gemini-2.5-flash (G)	~0.05–0.10	~0.20	0.70+	0.91
gemini-2.5-pro (No G)	0.10–0.20	~0.30	0.60	0.89
gemini-2.5-pro (G)	< 0.01	~0.25	0.75+	0.98
EHRSQL-MIMIC-III				
gpt-5.2 (No G)	~0.90+	~0.10–0.20	~0.05	0.84–0.86
gpt-5.2 (G)	~0.20–0.30	~0.40+	0.50+	0.95+
gemini-2.5-flash (No G)	~0.40–0.50	~0.30	~0.20	0.84
gemini-2.5-flash (G)	~0.05–0.10	~0.30	0.70+	0.90+
gemini-2.5-pro (No G)	~0.20–0.30	~0.30	~0.30	0.93–0.94
gemini-2.5-pro (G)	< 0.05	~0.30	0.75+	0.98
medgemma-1.5-4b-it (No G)	~0.12	~0.88	~0.00	0.85
medgemma-1.5-4b-it (G)	~0.17	~0.21	~0.00	0.76
SPIDER				
gpt-5.2 (No G)	0.80–0.90	~0.15	~0.05–0.10	0.17
gpt-5.2 (G)	0.15–0.25	~0.25	0.55+	0.65+
gemini-2.5-flash (No G)	0.85–0.90	~0.10	~0.05	0.14
gemini-2.5-flash (G)	~0.05–0.10	~0.20	0.70+	0.85+
gemini-2.5-pro (No G)	0.70–0.80	~0.20	~0.10–0.15	0.33
gemini-2.5-pro (G)	< 0.05	~0.25	0.75+	0.90+

Table 4: Failure mode analysis across datasets comparing models with and without query guidance (G). Guidance reduces non-SQL generation failures, especially on eICU, MIMIC-III, and SPIDER, but residual errors remain dominated by incorrect table selection. This indicates that guidance improves task framing and generation reliability, while schema grounding remains a bottleneck.

718 **NeurIPS Paper Checklist**

719 **1. Claims**

720 Question: Do the main claims made in the abstract and introduction accurately reflect the
721 paper’s contributions and scope?

722 Answer: [Yes]

723 Justification: The abstract and introduction consistently describe schema guidance as a
724 structured representation for managing context in agentic SQL generation and as a mecha-
725 nism for search-space reduction. The stated contributions, including the schema-guidance
726 construction pipeline, its interpretation as a compressed context representation, and its role in
727 improving SQL generation, are aligned across sections. The empirical scope—covering SPI-
728 DER, open clinical benchmarks (MIMIC-III and eICU), and a closed institutional dataset—is
729 consistently reflected in both the abstract and introduction. Differences in emphasis do not
730 affect the accuracy or scope of the stated claims.

731 Guidelines:

- 732 • The answer [N/A] means that the abstract and introduction do not include the claims
733 made in the paper.
- 734 • The abstract and/or introduction should clearly state the claims made, including the
735 contributions made in the paper and important assumptions and limitations. A [No] or
736 [N/A] answer to this question will not be perceived well by the reviewers.
- 737 • The claims made should match theoretical and experimental results, and reflect how
738 much the results can be expected to generalize to other settings.
- 739 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
740 are not attained by the paper.

741 **2. Limitations**

742 Question: Does the paper discuss the limitations of the work performed by the authors?

743 Answer: [Yes]

744 Justification: The paper includes a dedicated Limitations paragraph that explicitly discusses
745 key assumptions and constraints of the proposed approach. These include the reliance on
746 benchmark summaries rather than query-level traces, which limits detailed error analysis;
747 the dependence on structural SQL metrics without full validation of clinically meaningful
748 outputs; and potential system-level confounds arising from the broader agentic pipeline. The
749 discussion also reflects on the scope of empirical evaluation, including limitations related to
750 structured relational schemas and the lack of validation in noisier real-world settings (e.g.,
751 incomplete schemas, free-text-heavy databases, or hybrid data models). Additionally, the
752 authors acknowledge important deployment challenges such as interpretability, auditability,
753 and human oversight. Overall, the paper provides a clear and transparent account of the
754 main factors that may affect performance and generalization.

755 Guidelines:

- 756 • The answer [N/A] means that the paper has no limitation while the answer [No] means
757 that the paper has limitations, but those are not discussed in the paper.
- 758 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 759 • The paper should point out any strong assumptions and how robust the results are to
760 violations of these assumptions (e.g., independence assumptions, noiseless settings,
761 model well-specification, asymptotic approximations only holding locally). The authors
762 should reflect on how these assumptions might be violated in practice and what the
763 implications would be.
- 764 • The authors should reflect on the scope of the claims made, e.g., if the approach was
765 only tested on a few datasets or with a few runs. In general, empirical results often
766 depend on implicit assumptions, which should be articulated.
- 767 • The authors should reflect on the factors that influence the performance of the approach.
768 For example, a facial recognition algorithm may perform poorly when image resolution
769 is low or images are taken in low lighting. Or a speech-to-text system might not be
770 used reliably to provide closed captions for online lectures because it fails to handle
771 technical jargon.

- 772 • The authors should discuss the computational efficiency of the proposed algorithms
773 and how they scale with dataset size.
- 774 • If applicable, the authors should discuss possible limitations of their approach to
775 address problems of privacy and fairness.
- 776 • While the authors might fear that complete honesty about limitations might be used by
777 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
778 limitations that aren't acknowledged in the paper. The authors should use their best
779 judgment and recognize that individual actions in favor of transparency play an impor-
780 tant role in developing norms that preserve the integrity of the community. Reviewers
781 will be specifically instructed to not penalize honesty concerning limitations.

782 3. Theory assumptions and proofs

783 Question: For each theoretical result, does the paper provide the full set of assumptions and
784 a complete (and correct) proof?

785 Answer: [Yes]

786 Justification: The paper states its theoretical claims with explicit assumptions and cross-
787 references them to the appendix. The submodular coverage guarantee is stated in Remark ??
788 with the required monotonicity, non-negativity, submodularity, and budget assumptions,
789 and references the standard greedy approximation result. The search-space reduction result
790 is stated as a proposition in Appendix B.1, with the recall-preservation and monotone
791 decoder-error assumptions made explicit.

792 Guidelines:

- 793 • The answer [N/A] means that the paper does not include theoretical results.
- 794 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
795 referenced.
- 796 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 797 • The proofs can either appear in the main paper or the supplemental material, but if
798 they appear in the supplemental material, the authors are encouraged to provide a short
799 proof sketch to provide intuition.
- 800 • Inversely, any informal proof provided in the core of the paper should be complemented
801 by formal proofs provided in appendix or supplemental material.
- 802 • Theorems and Lemmas that the proof relies upon should be properly referenced.

803 4. Experimental result reproducibility

804 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
805 perimental results of the paper to the extent that it affects the main claims and/or conclusions
806 of the paper (regardless of whether the code and data are provided or not)?

807 Answer: [Yes]

808 Justification: The paper describes the schema-guidance construction pipeline, runtime
809 agent behavior, datasets, evaluation metrics, comparison protocol, model configurations,
810 decoding settings, and compute setup. The appendix further details the runtime prompt
811 structure and agentic SQL generation process, including tool-calling behavior, schema-
812 loading mechanisms, caching, and reset behavior between benchmark examples. Additional
813 implementation details, including fully instantiated prompts and generated guidance artifacts,
814 are provided in the supplemental material and accompanying code.

815 While the institutional dataset cannot be released due to privacy and governance constraints,
816 the same evaluation protocol is applied to publicly available benchmarks (SPIDER, MIMIC-
817 III, and eICU), enabling independent verification of the main claims.

818 Guidelines:

- 819 • The answer [N/A] means that the paper does not include experiments.
- 820 • If the paper includes experiments, a [No] answer to this question will not be perceived
821 well by the reviewers: Making the paper reproducible is important, regardless of
822 whether the code and data are provided or not.
- 823 • If the contribution is a dataset and/or model, the authors should describe the steps taken
824 to make their results reproducible or verifiable.

- 825 • Depending on the contribution, reproducibility can be accomplished in various ways.
826 For example, if the contribution is a novel architecture, describing the architecture fully
827 might suffice, or if the contribution is a specific model and empirical evaluation, it may
828 be necessary to either make it possible for others to replicate the model with the same
829 dataset, or provide access to the model. In general, releasing code and data is often
830 one good way to accomplish this, but reproducibility can also be provided via detailed
831 instructions for how to replicate the results, access to a hosted model (e.g., in the case
832 of a large language model), releasing of a model checkpoint, or other means that are
833 appropriate to the research performed.
- 834 • While NeurIPS does not require releasing code, the conference does require all submis-
835 sions to provide some reasonable avenue for reproducibility, which may depend on the
836 nature of the contribution. For example
 - 837 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
838 to reproduce that algorithm.
 - 839 (b) If the contribution is primarily a new model architecture, the paper should describe
840 the architecture clearly and fully.
 - 841 (c) If the contribution is a new model (e.g., a large language model), then there should
842 either be a way to access this model for reproducing the results or a way to reproduce
843 the model (e.g., with an open-source dataset or instructions for how to construct
844 the dataset).
 - 845 (d) We recognize that reproducibility may be tricky in some cases, in which case
846 authors are welcome to describe the particular way they provide for reproducibility.
847 In the case of closed-source models, it may be that access to the model is limited in
848 some way (e.g., to registered users), but it should be possible for other researchers
849 to have some path to reproducing or verifying the results.

850 5. Open access to data and code

851 Question: Does the paper provide open access to the data and code, with sufficient instruc-
852 tions to faithfully reproduce the main experimental results, as described in supplemental
853 material?

854 Answer: [Yes]

855 Justification: The paper provides anonymized code, prompts, and configuration details in
856 the supplemental material to support reproducibility. The submission includes the schema
857 and query guidance construction pipeline, runtime agent implementation, and evaluation
858 scripts. All experiments are conducted on publicly available datasets (SPIDER, MIMIC-III,
859 and eICU), for which access instructions are provided. The institutional dataset cannot
860 be released due to privacy and governance constraints, but it is used only as an additional
861 evaluation setting.

862 Guidelines:

- 863 • The answer [N/A] means that paper does not include experiments requiring code.
- 864 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
865 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 866 • While we encourage the release of code and data, we understand that this might not
867 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
868 including code, unless this is central to the contribution (e.g., for a new open-source
869 benchmark).
- 870 • The instructions should contain the exact command and environment needed to run to
871 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 872 • The authors should provide instructions on data access and preparation, including how
873 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 874 • The authors should provide scripts to reproduce all experimental results for the new
875 proposed method and baselines. If only a subset of experiments are reproducible, they
876 should state which ones are omitted from the script and why.
- 877 • At submission time, to preserve anonymity, the authors should release anonymized
878 versions (if applicable).
- 879

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies datasets, benchmark construction, evaluation metrics, comparison protocols, model configurations, decoding settings, agent control parameters, retry policies, and compute setup. As the system is evaluated in an inference-only setting, no training procedure or optimizer is required. The appendix further details runtime agent behavior, schema-loading tools, and prompt structure across examples. Additional implementation details, including fully instantiated prompts and generated guidance artifacts, are provided in the supplemental material and accompanying code.

Some low-level implementation details, such as exact fully instantiated prompts and generated guidance artifacts, are summarized rather than exhaustively enumerated. These details do not prevent interpretation of the reported results, and additional implementation details will be provided with the released code.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The paper reports results as mean \pm standard error across aggregated evaluations over datasets and runs, as shown in the tables. Each evaluation corresponds to a repeated inference under identical experimental conditions, capturing variability arising from model inference and agent execution.

The reported error bars correspond to the standard error of the mean (SEM), computed as the sample standard deviation divided by \sqrt{n} , where n is the number of aggregated dataset-run evaluations.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).

- 933 • If error bars are reported in tables or plots, the authors should explain in the text how
934 they were calculated and reference the corresponding figures or tables in the text.

935 8. Experiments compute resources

936 Question: For each experiment, does the paper provide sufficient information on the com-
937 puter resources (type of compute workers, memory, time of execution) needed to reproduce
938 the experiments?

939 Answer: [Yes]

940 Justification: Experiments were conducted using API-based inference for GPT and Gemini
941 models via Google Cloud Workbench notebooks; these models were accessed through
942 external cloud services and were not executed locally. As a result, hardware specifications
943 for these models are abstracted by the provider, and compute is primarily determined by
944 API latency and throughput.

945 In contrast, MedGemma-1.5-4B-IT was executed locally on a Google Cloud VM instance
946 with a g2-standard-16 configuration (16 vCPUs, 64 GB RAM) and a single NVIDIA L4
947 GPU (24 GB VRAM).

948 Each dataset evaluation required on the order of hours, depending on the number of queries.
949 These specifications reflect the primary computational requirements needed to reproduce
950 the reported experiments.

951 Guidelines:

- 952 • The answer [N/A] means that the paper does not include experiments.
- 953 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
954 or cloud provider, including relevant memory and storage.
- 955 • The paper should provide the amount of compute required for each of the individual
956 experimental runs as well as estimate the total compute.
- 957 • The paper should disclose whether the full research project required more compute
958 than the experiments reported in the paper (e.g., preliminary or failed experiments that
959 didn't make it into the paper).

960 9. Code of ethics

961 Question: Does the research conducted in the paper conform, in every respect, with the
962 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

963 Answer: [Yes]

964 Justification: The research complies with the NeurIPS Code of Ethics. The study uses
965 de-identified or publicly available datasets where applicable, and the closed institutional
966 dataset is handled in accordance with appropriate privacy and data governance practices.
967 The work does not involve human subject experimentation or deployment in a real-world
968 clinical setting.

969 Guidelines:

- 970 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
971 Ethics.
- 972 • If the authors answer [No], they should explain the special circumstances that require a
973 deviation from the Code of Ethics.
- 974 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
975 eration due to laws or regulations in their jurisdiction).

976 10. Broader impacts

977 Question: Does the paper discuss both potential positive societal impacts and negative
978 societal impacts of the work performed?

979 Answer: [Yes]

980 Justification: The paper discusses positive societal impacts, including more reliable and
981 interpretable cohort construction, deep phenotyping, and outcome ascertainment in clinical
982 research. It also discusses potential negative impacts, including misleading clinical analyses
983 from incorrect SQL generation, misuse of automated cohort construction without oversight,
984 bias propagation, and privacy or governance risks when applied to sensitive clinical data. The

985 paper further identifies mitigation strategies such as human-in-the-loop review, validation
986 checks before execution, and deployment within controlled environments with appropriate
987 data governance.

988 Guidelines:

- 989 • The answer [N/A] means that there is no societal impact of the work performed.
- 990 • If the authors answer [N/A] or [No], they should explain why their work has no societal
991 impact or why the paper does not address societal impact.
- 992 • Examples of negative societal impacts include potential malicious or unintended uses
993 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
994 (e.g., deployment of technologies that could make decisions that unfairly impact specific
995 groups), privacy considerations, and security considerations.
- 996 • The conference expects that many papers will be foundational research and not tied
997 to particular applications, let alone deployments. However, if there is a direct path to
998 any negative applications, the authors should point it out. For example, it is legitimate
999 to point out that an improvement in the quality of generative models could be used to
1000 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1001 that a generic algorithm for optimizing neural networks could enable people to train
1002 models that generate Deepfakes faster.
- 1003 • The authors should consider possible harms that could arise when the technology is
1004 being used as intended and functioning correctly, harms that could arise when the
1005 technology is being used as intended but gives incorrect results, and harms following
1006 from (intentional or unintentional) misuse of the technology.
- 1007 • If there are negative societal impacts, the authors could also discuss possible mitigation
1008 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1009 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1010 feedback over time, improving the efficiency and accessibility of ML).

1011 11. Safeguards

1012 Question: Does the paper describe safeguards that have been put in place for responsible
1013 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1014 image generators, or scraped datasets)?

1015 Answer: [N/A]

1016 Justification: The paper does not release a high-risk model, pretrained language model, image
1017 generator, or scraped dataset requiring special safeguards. The closed institutional clinical
1018 dataset is not publicly released, and the work focuses on a schema-guidance construction
1019 method evaluated with existing models rather than releasing a deployable model artifact.

1020 Guidelines:

- 1021 • The answer [N/A] means that the paper poses no such risks.
- 1022 • Released models that have a high risk for misuse or dual-use should be released with
1023 necessary safeguards to allow for controlled use of the model, for example by requiring
1024 that users adhere to usage guidelines or restrictions to access the model or implementing
1025 safety filters.
- 1026 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1027 should describe how they avoided releasing unsafe images.
- 1028 • We recognize that providing effective safeguards is challenging, and many papers do
1029 not require this, but we encourage authors to take this into account and make a best
1030 faith effort.

1031 12. Licenses for existing assets

1032 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1033 the paper, properly credited and are the license and terms of use explicitly mentioned and
1034 properly respected?

1035 Answer: [Yes]

1036 Justification: The paper cites the original sources for the existing datasets used, including
1037 MIMIC-III, eICU, EHRSQL, and Spider-style text-to-SQL benchmarks. Appendix C

1038 explicitly describes licensing and data usage terms: Spider is released under Apache-
1039 2.0, EHRSQL under CC-BY-4.0, MIMIC-III and eICU are accessed under credentialed
1040 PhysioNet data use agreements, and the closed institutional dataset is not publicly released
1041 due to privacy and governance constraints. All datasets are used for research purposes in
1042 compliance with their respective licenses and terms of use.

1043 Guidelines:

- 1044 • The answer [N/A] means that the paper does not use existing assets.
- 1045 • The authors should cite the original paper that produced the code package or dataset.
- 1046 • The authors should state which version of the asset is used and, if possible, include a
1047 URL.
- 1048 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1049 • For scraped data from a particular source (e.g., website), the copyright and terms of
1050 service of that source should be provided.
- 1051 • If assets are released, the license, copyright information, and terms of use in the
1052 package should be provided. For popular datasets, paperswithcode.com/datasets
1053 has curated licenses for some datasets. Their licensing guide can help determine the
1054 license of a dataset.
- 1055 • For existing datasets that are re-packaged, both the original license and the license of
1056 the derived asset (if it has changed) should be provided.
- 1057 • If this information is not available online, the authors are encouraged to reach out to
1058 the asset's creators.

1059 13. New assets

1060 Question: Are new assets introduced in the paper well documented and is the documentation
1061 provided alongside the assets?

1062 Answer: [N/A]

1063 Justification: The paper does not release new assets such as a public dataset, codebase, or
1064 model checkpoint. The closed institutional benchmark is used for evaluation but is not
1065 publicly released due to privacy and data governance constraints, as described in Appendix C.

1066 Guidelines:

- 1067 • The answer [N/A] means that the paper does not release new assets.
- 1068 • Researchers should communicate the details of the dataset/code/model as part of their
1069 submissions via structured templates. This includes details about training, license,
1070 limitations, etc.
- 1071 • The paper should discuss whether and how consent was obtained from people whose
1072 asset is used.
- 1073 • At submission time, remember to anonymize your assets (if applicable). You can either
1074 create an anonymized URL or include an anonymized zip file.

1075 14. Crowdsourcing and research with human subjects

1076 Question: For crowdsourcing experiments and research with human subjects, does the paper
1077 include the full text of instructions given to participants and screenshots, if applicable, as
1078 well as details about compensation (if any)?

1079 Answer: [N/A]

1080 Justification: The study does not involve crowdsourcing or prospective experiments with
1081 human participants (e.g., no participant recruitment, instructions, or compensation). It is
1082 based on retrospective analysis of de-identified clinical and/or publicly available datasets.
1083 The use of institutional data was reviewed and approved by an Institutional Review Board
1084 (IRB), and all analyses were conducted in compliance with applicable data governance and
1085 privacy regulations.

1086 Guidelines:

- 1087 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1088 with human subjects.

- 1089
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- 1090
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 1091
- 1092
- 1093
- 1094

1095 **15. Institutional review board (IRB) approvals or equivalent for research with human**
1096 **subjects**

1097 Question: Does the paper describe potential risks incurred by study participants, whether
1098 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1099 approvals (or an equivalent approval/review based on the requirements of your country or
1100 institution) were obtained?

1101 Answer: [Yes]

1102 Justification: The study involves retrospective analysis of clinical data. Institutional Review
1103 Board (IRB) approval was obtained (IRB No. 26-002996). All data were handled in
1104 compliance with relevant privacy and data protection regulations.

1105 Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1116 **16. Declaration of LLM usage**

1117 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1118 non-standard component of the core methods in this research? Note that if the LLM is used
1119 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1120 scientific rigor, or originality of the research, declaration is not required.

1121 Answer: [Yes]

1122 Justification: Large language models (LLMs) are a core component of both the schema
1123 guidance construction pipeline and the runtime SQL generation process. In the construction
1124 pipeline, LLMs are used for enrichment during per-table guidance generation. At inference
1125 time, LLMs are used within the agent to generate SQL queries conditioned on schema
1126 guidance. In addition, LLMs are used for zero-shot derivation of natural language task
1127 descriptions in preprocessing of the closed institutional dataset. These uses are described in
1128 the Method and Experimental Setup sections.

1129 Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.

1130

1131

1132

1133