

SPOTIT : EVALUATING TEXT-TO-SQL EVALUATION WITH FORMAL VERIFICATION

Rocky Klopfenstein^{1*}, Yang He^{2*}, Andrew Tremante^{1*}, Yuepeng Wang²,
Nina Narodytska³, Haoze Wu^{1,3}

¹Amherst College ²Simon Fraser University ³VMware Research by Broadcom
{rklopfenstein27, atremante26, hwu}@amherst.edu
{yha244, yuepeng}@sfu.ca, nina.narodytska@broadcom.com

ABSTRACT

Community-driven Text-to-SQL evaluation platforms play a pivotal role in tracking the state of the art of Text-to-SQL performance. The reliability of the evaluation process is critical for driving progress in the field. Current evaluation methods are largely test-based, which involves comparing the execution results of a generated SQL query and a human-labeled ground-truth on a static test database. Such an evaluation is optimistic, as two queries can coincidentally produce the same output on the test database while actually being different. In this work, we propose a new alternative evaluation pipeline, called SPOTIT, where a formal bounded equivalence verification engine actively searches for a database that differentiates the generated and ground-truth SQL queries. We develop techniques to extend existing verifiers to support a richer SQL subset relevant to Text-to-SQL. A performance evaluation of ten Text-to-SQL methods on the high-profile BIRD dataset suggests that test-based methods can often overlook differences between the generated query and the ground-truth. Further analysis of the verification results reveals a more complex picture of the current Text-to-SQL evaluation.

1 INTRODUCTION

Text-to-SQL is one of the fundamental building blocks for designing natural language (NL) interfaces that enable users to access and analyze structured data sources. Translating human questions into executable database queries bridges the gap between non-technical users and complex data systems. This functionality underpins modern chatbots and smart assistants across a wide range of industrial applications, such as observability platforms for monitoring system health (BitsAI, 2025; Splunk 2025), critical business processes (Amazon, 2025), and healthcare (Amazon Web Services, 2024).

Due to its practical relevance for commercial products, Text-to-SQL has recently attracted significant attention, leading to the development of a wide range of solutions (Shi et al., 2024). New Text-to-SQL frameworks are announced regularly, and thanks to community-driven evaluation platforms such as BIRD (Li et al., 2024) and Spider (Lei et al., 2024), their performance can be benchmarked and compared in near real time. Given the pivotal role these platforms play in tracking the state of the art, the reliability of their evaluation processes is crucial for driving progress in the field.

In this paper, we take a close look at the evaluation process for the accuracy of Text-to-SQL methods. Currently, the process usually involves checking whether the SQL queries generated by a method produce results equivalent to those of the *gold SQLs* (i.e., human-written ground-truth SQLs), under a pre-defined notion of equivalence. Most state-of-the-art evaluation frameworks (Li et al., 2024; Lei et al., 2024) perform this equivalence check through *testing*: executing both queries on a static test database and comparing the results. If the results match, the generated SQL is labeled as correct. Although widely used in practice, the testing-based approach has clear limitations. Because the check is performed on a single database, two different SQL queries may appear equivalent by chance, purely due to the specific data contained in that database. This raises an important question: when the test-based approach marks a generated SQL as correct, how often does it truly produce the same

*Equal contribution.

results as the gold SQL in general? The next broader question is: to what extent can the current evaluation process accurately measure the performance of Text-to-SQL methods?

We investigate these questions by exploring an alternative correctness evaluation methodology. Instead of relying on test databases to assess equivalence, we propose to actively *search* for databases that can differentiate the generated SQL from the gold SQL. The search-based evaluation naturally provides stronger correctness guarantees and enables a more rigorous measurement of accuracy. Since providing complete equivalence guarantee is in general undecidable, we perform SMT-based bounded verification (He et al., 2024), which searches for differentiating databases with specified sizes. We develop a new Text-to-SQL evaluation workflow, SPOTIT, on top of those verification techniques. We significantly extend these techniques to support a new set of SQL operators over strings and dates which are commonly used for Text-to-SQL benchmarks.

Experiments on ten state-of-the-art Text-to-SQL methods on the popular BIRD dataset (Li et al., 2024) suggest that the reported accuracy of these methods drops by 11.3%–14.2% when switching from the official test-based evaluation to SPOTIT. The varying levels of decrease in absolute precision also lead to substantial changes in the order of ranking of the Text-to-SQL methods. Moreover, SPOTIT produces minimal differentiating databases, which enables us to pinpoint the sources of inconsistencies between the generated and gold SQLs. Analysis of these databases uncovers several shortcomings of the current Text-to-SQL evaluation process. Most surprisingly, we find that when the predicted SQL disagrees with the ground truth, it is often the gold SQL that is incorrect.

To summarize, our contributions include:

- SPOTIT, a new evaluation pipeline for Text-to-SQL powered by formal equivalence verification;
- novel SMT-encoding for a set of SQL operators over strings and dates, and proof of its correctness;
- practical strategies for the efficient deployment of SPOTIT;
- a large-scale evaluation of ten state-of-the-art Text-to-SQL methods on the BIRD dataset, which reveals several potential shortcomings of current Text-to-SQL evaluation.

2 PRELIMINARIES

We provide background on Text-to-SQL and formal equivalence checking. Due to space limitation, an overview of related work is present in App. A.

Text-to-SQL problem statement. Given a natural language query N and a database D with schema \mathcal{S} , the goal of Text-to-SQL is to map (N, D) to an SQL query Q , such that executing Q on D , denoted $Q(D)$, produces an output relation (table) that answers N .

Text-to-SQL evaluation. The main evaluation mechanism for a Text-to-SQL framework relies on a gold SQL query produced by a human annotator. Hence, for each natural language query N over a database, there exists a gold SQL query Q that represents the human-labelled ground truth of translating N into SQL. Given a generated SQL P and the corresponding gold query Q , current evaluation performs the following check:

$$\text{EX-TEST}(P, Q, D_{\text{test}}) = \begin{cases} 1, & \text{if } \forall r. r \in P(D_{\text{test}}) \leftrightarrow r \in Q(D_{\text{test}}) \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

where D_{test} is a test database provided by the benchmark set, and r denotes a row in the result table. In words, EX-TEST compares whether the two tables, $P(D_{\text{test}})$ and $Q(D_{\text{test}})$, contain the same set of rows. In order to more rigorously analyze the equivalence between P and Q , we use formal verification to search for a differentiating database D_{cex} such that $\text{EX-TEST}(P, Q, D_{\text{cex}}) = 0$.

Bounded SQL equivalence checking. Given two SQL queries Q_1 and Q_2 over a schema \mathcal{S} and an upper bound K on the relation size, the problem of bounded equivalence checking is to decide whether Q_1 and Q_2 are equivalent, denoted $Q_1 \simeq_{\mathcal{S}, K} Q_2$, for all databases D conforming to \mathcal{S} such that each relation in D has at most K tuples. Formally,

$$Q_1 \simeq_{\mathcal{S}, K} Q_2 \stackrel{\text{def}}{=} \forall D \in \text{Instances}(\mathcal{S}). \forall R \in \text{Relations}(D). |R| \leq K \Rightarrow Q_1(D) = Q_2(D),$$

where $\text{Instances}(\mathcal{S})$ represents all database instances conforming to \mathcal{S} , and $\text{Relations}(D)$ represents all relations in D . In general, the goal is either to prove the bounded equivalence holds, or to find a counterexample database D_{cex} that disproves the equivalence. Compared with unbounded equivalence

<p>N_1: "Which is the youngest patient with an abnormal anti-ribonuclear protein level? Please list his or her date of birth."</p> <hr/> <pre> /*Gold SQL Q*/: SELECT T1.birthday FROM patient AS T1 INNER JOIN laboratory AS T2 ON T1.ID = T2.ID WHERE T2.rnp != '-' OR '+-' ORDER BY T1.birthday DESC LIMIT 1 /*Generated SQL P*/: SELECT patient.birthday FROM patient INNER JOIN laboratory ON patient.ID = laboratory.ID WHERE NOT laboratory.rnp IN ('-', '+-') ORDER BY patient.birthday DESC LIMIT 1 </pre>	<p>N_2: "How many male patients who underwent testing between 1995 and 1997 and were subsequently diagnosed with Behcet disease did not stay in the hospital for treatment?"</p> <hr/> <pre> /*Gold SQL Q*/: SELECT COUNT(T1.id) FROM patient AS T1 INNER JOIN examination AS T2 ON T1.id = T2.id WHERE T2.diagnosis = 'Behcet' AND T1.sex = 'M' AND STRFTIME('%Y', T2.examination_date) BETWEEN '1995' AND '1997' AND T1.admission = '-'; /*Generated SQL P*/: SELECT COUNT(DISTINCT patient.id) FROM patient INNER JOIN examination ON patient.id = examination.id WHERE patient.sex = 'M' AND examination.examination_date BETWEEN '1995-01-01' AND '1997-12-31' AND examination.diagnosis = 'Behcet' AND patient.admission = '-'; </pre>
--	---

Figure 1: Examples of cases where the generated SQL produces the same output as the gold SQL on the BIRD’s official test database, but SPOTIT finds a database that differentiates the the queries. The parts that explain the mismatch are highlighted. For N_1 , the gold SQL is incorrect. And for N_2 , both SQL queries can be right depending on the interpretation of the NL question.

checking, which is generally undecidable (Mohamed et al., 2024), bounded equivalence checking can handle a more expressive SQL subset and is guaranteed to uncover small counterexamples (if they exist). These features make bounded verification suitable for large-scale Text-to-SQL evaluation.

VERIEQL. VERIEQL (He et al., 2024) is a recently proposed bounded equivalence checker for SQL queries and, to the best of our knowledge, supports the most expressive subset of SQL among existing tools. It reduces the verification task to a satisfiability problem by encoding the symbolic execution of the two SQL queries and the *non-equivalence* of the execution results as a satisfiability modulo theories (SMT) formula (Barrett & Tinelli, 2018), which can be solved by an off-the-shelf SMT solver (De Moura & Bjørner, 2008). The bounded equivalence property holds if and only if the formula is unsatisfiable, which means it is not possible to find a database that result in different execution results. Otherwise, a satisfying interpretation of the formula can be decoded to a counterexample database. We significantly extend VERIEQL to support our verification use cases.

3 MOTIVATING EXAMPLES

Before we describe our new verification-based evaluation pipeline, we first discuss main sources of mismatches between the gold SQL and the generated SQL in Text-to-SQL evaluation. There are three main such sources: (1) NL query N is ambiguous, so both the gold and generated SQL queries are justifiable interpretations; (2) N is unambiguous, but the gold SQL query is incorrect (gold SQLs are created manually and thus prone to human errors); (3) N is unambiguous, the gold SQL query is correct, but the generated SQL query is incorrect. Our framework focuses on checking equivalence between the gold SQL and the generated SQL, treating the latter as the best-effort, semantically correct formalization of N . We show that SPOTIT can successfully detect incorrect generated SQLs that are overlooked by existing test-based evaluation. Perhaps more surprisingly, SPOTIT also allows us to spot the first and second sources of mismatch. Fig. 1 shows two illustrative examples.

Example 3.1. Consider the query N_1 : "Which is the youngest patient with an abnormal anti-ribonuclear protein level? Please list his or her date of birth." together with the gold and generated SQL queries. On the development database that BIRD provides, both queries return "1989-08-28". However, SPOTIT found a database on which these two queries are not equivalent (Appendix D.1). In fact, we observe that all ten frameworks that we tested generated SQLs that are not equivalent to the gold query. Upon closer inspection, we find that the gold query is incorrect: its WHERE clause is equivalent to $T2.rnp \neq '-' \text{ OR } FALSE$, as a string literal like '+' is interpreted as FALSE in a boolean context, which is not the intended behavior. □

Example 3.2. Consider another query N_2 : "How many male patients who underwent testing between 1995 and 1997 and were subsequently diagnosed with Behcet disease did not stay in the hospital for treatment?" together with the gold and generated SQL queries. These two queries both return "2" on the BIRD test database. However, the two queries are clearly not equivalent (*id* is not a primary key of the *examination* table therefore duplicates are allowed): the generated query counts all

Query Q_r	::=	Q OrderBy (Q, \vec{E}, b)
Subquery Q	::=	R $\Pi_L(Q)$ $\sigma_\phi(Q)$ $\rho_R(Q)$ $Q \oplus Q$ Distinct (Q) $Q \otimes Q$ GroupBy (Q, \vec{E}, L, ϕ) With (\vec{Q}, \vec{R}, Q)
Attr List L	::=	$id(A)$ $\rho_a(A)$ L, L
Attr A	::=	Cast (ϕ) E $\mathcal{G}(E)$ $A \diamond A$
Pred ϕ	::=	b Null $A \odot A$ IsNull (E) $\vec{E} \in \vec{v}$ $\vec{E} \in Q$ $\phi \wedge \phi$ $\phi \vee \phi$ $\neg\phi$ PrefixOf (s, E) SuffixOf (s, E) Like (s, E) Contain (s, E)
Expr E	::=	a v $E \diamond E$ ITE (ϕ, E, E) Case ($\vec{\phi}, \vec{E}, E$) SubStr (E_1, E_2, E_3) Concat (E_1, E_2) Strftime (κ, E) JulianDay (E) DateShift (E, i, δ) ToInt (E) ToDate (E) ToStr (E)
Join Op \otimes	::=	\times \bowtie_ϕ $\bowtie_{\neg\phi}$ $\bowtie_{\phi, \psi}$ $\bowtie_{\neg\phi, \psi}$
Collection Op \oplus	::=	\cup \cap \setminus \sqcup \sqcap $-$
Arith Op \diamond	::=	$+$ $-$ \times $/$ $\%$
Logic Op \odot	::=	\leq $<$ $=$ \neq $>$ \geq

$R \in$ Relation Names $a \in$ Attribute Names $v \in \{\text{Null}\} \cup \text{Integers} \cup \text{Dates} \cup \text{Strings}$ $b \in$ Bools $i \in$ Integers
 $s \in$ Strings $\mathcal{G} \in \{\text{Count, Min, Max, Sum, Avg}\}$ $\kappa \in \{\text{"\%Y"}, \text{"\%M"}, \text{"\%d"}\}$ $\delta \in \{\text{"Year"}, \text{"Month"}, \text{"Day"}\}$

Figure 2: Extended syntax of SQL Queries. New features are in bold.

examinations per patient, whereas the gold query counts only distinct patients. SPOTIT easily found a database that differentiate the two queries (Appendix D.2). Note that depending on the interpretation of the question, both SQL queries can be correct: the gold SQL can be reasonable if the goal is to understand the hospital workload, while the generated SQL can be reasonable if the goal is to understand the number of unique patients. Hence, we conclude that N_2 is ambiguous. \square

Note that these examples were overlooked by existing test-based evaluations. On the other hand, using SPOTIT, we found that undetected cases like those are quite common in the BIRD dataset.

4 METHODOLOGY

In this section, we introduce new SMT-encodings for a number of SQL operators over string and date types that were not supported by existing bounded equivalence verification methods but frequently appear in Text-to-SQL benchmarks. Then we present our verification-based evaluation pipeline SPOTIT and discuss practical implementation strategies.

4.1 EQUIVALENCE CHECKING FOR SQL QUERIES

To understand our extension, let us first walk through Example 4.1 to understand how equivalence checking can be encoded as an SMT formula in a verifier like VERIEQL (He et al., 2024).

Example 4.1. Consider a schema $\mathcal{S} = \{R \mapsto \{id:\text{int}, dob:\text{date}\}\}$ and the following two queries:

$$Q_1 = \text{SELECT } id \text{ FROM } R \text{ WHERE } id > 1 \quad Q_2 = \text{SELECT } id \text{ FROM } R \text{ WHERE } id > 2$$

We describe how to encode equivalence checking for a bound (K) of 1 as an SMT formula. First, variables are introduced to represent the database and the execution results. This includes a symbolic database $D = \{R \mapsto [t_1]\}$, where $t_1 = [x_1, x_2]$ is a tuple in R , and x_1, x_2 are integer variables. In addition, tuples $t_2 = [x_3]$ and $t_3 = [x_4]$, are introduced to encode query results: $Q_1(D) = [t_2]$ and $Q_2(D) = [t_3]$, where x_3, x_4 are both integer variables. Note that the number of tuples in R is equal to the bound K . Also note that a date (x_2) is represented as an integer, which is sufficient here but not in general. We later introduce precise encoding of date to support richer operations.

We now describe the constraints over the variables. The first set of constraints ensures that t_2 and t_3 correctly capture the semantics of Q_1 and Q_2 . In this case, t_2 tuple is constrained by $\Phi_{Q_1} = (x_1 > 1 \rightarrow (x_3 = x_1 \wedge \neg \text{Del}(t_2))) \wedge (x_1 \leq 1 \rightarrow \text{Del}(t_2))$, where Del is an uninterpreted function denoting the non-existence of a symbolic tuple. The formula Φ_{Q_1} ensures that only interpretations satisfying $x > 1$ can populate a concrete tuple; otherwise, Q_1 's result is empty. Similarly, t_3 is constrained by $\Phi_{Q_2} = (x_1 > 2 \rightarrow (x_4 = x_1 \wedge \neg \text{Del}(t_3))) \wedge (x_1 \leq 2 \rightarrow \text{Del}(t_3))$.

The second set of constraints encodes that $Q_1(D)$ and $Q_2(D)$ returns different results. In this case, it is simply $t_2 \neq t_3$. The full encoding is a conjunction of all constraints: $\Phi_{Q_1} \wedge \Phi_{Q_2} \wedge (t_2 \neq t_3)$, whose satisfiability can be checked by an SMT solver. A satisfying interpretation to this conjunction corresponds to a database instance that differentiates Q_1 and Q_2 . For example, the queries are not equivalent under the interpretation $\mathcal{I} = \{x_1 \mapsto 2\}$. \square

Extension in SQL encoding. Existing bounded SQL equivalence checker still lacks support for several important features, including precise encoding of dates and strings, which are highly relevant

in Text-to-SQL applications. Furthermore, SQL supports computations across many different data types with implicit type casting (e.g., $1 + \text{“a”}$ and $\text{date}(\text{“2000-01-01”}) + \text{“1”}$), which poses significant challenges to establish precise semantics and encodings. To address these limitations and challenges, we introduce techniques to support dates and strings, along with their manipulations, in the SQL equivalence checker VERIEQL. We also introduce type conversions across Null, integers, dates, and strings for implicit type casting. For example, in the gold SQL for N_2 (Fig. 1), the output of the STRFTIME function is implicitly converted from a date to an integer.

Fig. 2 presents our supported SQL grammar. Specifically, the query language introduces type conversions among various data types (e.g., $\text{ToInt}(E)$, $\text{ToDate}(E)$, and $\text{ToString}(E)$), which allows us to precisely establish the semantics of dates and strings and enhances the expressiveness of our SQL subset. We also incorporate additional expressions and predicates for data and string manipulations, such as date formatting $\text{Strftime}(\kappa, E)$, Julian day $\text{JulianDay}(E)$, string pattern matching $\text{PrefixOf}(s, E)$, $\text{SuffixOf}(s, E)$, $\text{Like}(s, E)$, and string truncation $\text{SubStr}(E_1, E_2, E_3)$. The symbolic encoding for these extended expressions and predicates is formally presented in Appendix F

As an example, we describe how to precisely encode a date variable, which is very common in Text-to-SQL. For instance, the date of birth and the time of a transaction are naturally modeled with the date type. Previously, date was encoded as a single integer variable (see Example 4.1). Although this coarse representation still enables the encoding of certain date operations (e.g., comparison), it does not necessarily support all date operations, such as date-formatting, which is used in the gold SQL query for N_2 in Fig. 1. As a date can be viewed as a triplet (year, month, day), we introduce three integer variables y , m , and d , and constrain their values with the following formula Φ :

$$\begin{aligned} \Phi &= \Phi_1 \wedge \Phi_2 \wedge \Phi_3, \text{ where } \Phi_1 = \text{MIN_YEAR} \leq y \leq \text{MAX_YEAR}, \Phi_2 = 1 \leq m \leq 12, \\ \Phi_3 &= 1 \leq d \wedge (\bigvee_{c \in \{1,3,5,7,8,10,12\}} m = c \rightarrow d \leq 31) \\ &\quad \wedge (m = 2 \rightarrow d \leq 28 + \text{ite}(\text{leap}(y), 1, 0)) \wedge (\bigvee_{c \in \{4,6,9,11\}} m = c \rightarrow d \leq 30) \end{aligned}$$

The term $\text{leap}(y)$ encodes the leap year condition: $y \% 4 = 0 \wedge (y \% 100 \neq 0 \vee y \% 400 = 0)$. Constraints Φ_1 , Φ_2 , and Φ_3 restrict the possible values of the year, the month, and the day, respectively. For example Φ_1 specifies the valid range of the year, which is specific to the database engine. For example, SQLite only accepts dates between “0000-01-01” and “9999-12-31”; in which case MIN_YEAR is 0 and MAX_YEAR is 9999. This refined representation allows us to precisely encode a rich set of date operations and analyze more SQL queries compared to the previous encoding.

Equivalence under set semantics. SQL equivalence checkers typically support equivalence under bag semantics and list semantics. However, some Text-to-SQL evaluation platforms, such as BIRD (Li et al., 2024), by default adopt equivalence under set semantics (see equation 1). This can be expressed as an SMT constraint. Given two query results with symbolic tables $R_1 = [t_1, \dots, t_n]$ and $R_2 = [r_1, \dots, r_m]$, the condition that R_1 and R_2 are equivalent under set semantics is as follows:

$$\bigwedge_{i=1}^n (\neg \text{Del}(t_i) \rightarrow \bigvee_{j=1}^m (\neg \text{Del}(r_j) \wedge t_i = r_j)) \wedge \bigwedge_{j=1}^m (\neg \text{Del}(r_j) \rightarrow \bigvee_{i=1}^n (\neg \text{Del}(t_i) \wedge r_j = t_i)) \quad (2)$$

On a high level, equivalence is defined by mutual set containment: $R_1 = R_2$ iff $R_1 \subseteq R_2$ and $R_2 \subseteq R_1$. But since some tuples might be deleted due to WHERE clauses, we restrict set containment to non-deleted tuples, i.e., those satisfying $\neg \text{Del}(t)$.

Correctness of the encodings. We now state the correctness of our symbolic encoding for the extended expressions and predicates, as well as the equivalence under set semantics. Proof of these theorems is in Appendix G. As we encode the symbolic execution of queries, to prove the correctness of our approach, we need to show that our symbolic execution coincides with the concrete execution. This involves showing that given an expression E , the satisfying interpretation of E ’s symbolic execution result is identical to the concrete execution result of E . Thm. 1 states that formally.

Theorem 1 (Correctness of expression encoding). *Let D be a database over schema \mathcal{S} , xs be a tuple list, and E be an expression. Consider a symbolic database Γ over \mathcal{S} , a list of symbolic tuples \mathcal{T} , and E ’s symbolic encoding $\llbracket E \rrbracket_{\mathcal{S}, \Gamma, \mathcal{T}}$. For any satisfying interpretation \mathcal{I} with $\mathcal{I}(\Gamma) = D \wedge \mathcal{I}(\mathcal{T}) = xs$, evaluating the expression E over the database D and the tuple list xs yields the interpretation of E ’s symbolic encoding $\mathcal{I}(\llbracket E \rrbracket_{\mathcal{S}, \Gamma, \mathcal{T}})$, i.e., $\mathcal{I}(\Gamma) = D \wedge \mathcal{I}(\mathcal{T}) = xs \Rightarrow \llbracket E \rrbracket_{D, xs} = \mathcal{I}(\llbracket E \rrbracket_{\mathcal{S}, \Gamma, \mathcal{T}})$.*

Similarly, given a predicate ϕ , the satisfying interpretation of ϕ ’s symbolic execution result is also identical to the concrete execution result of ϕ . This is formally stated in Appendix G. Lastly, we state the correctness of our encoding for equivalence under set semantics.

Algorithm 1 Bounded equivalence checking

Require: Database schema \mathcal{S} , gold SQL query Q , generated SQL query P , time limit T , bound K

Ensure: A counterexample D_{cex}

```

1: function EQUIVCHECK( $\mathcal{S}, Q, P, T, K$ )
2:   for  $k \in [1, K]$  do
3:     res,  $D_{\text{cex}} \leftarrow$  CHECKBOUND( $\mathcal{S}, P, Q, k, T$ )
4:     if res = EQUIVALENT then continue
5:     ▷ Bounded equivalence under  $k$ 
6:     else if res = NON-EQUIVALENT then
7:       ▷ Find a counterexample
8:       ▷ Validate the counterexample on the backend DBMS
9:       if  $\neg$ EX-TEST( $P, Q, D_{\text{cex}}$ ) then
10:        return  $\{D_{\text{cex}}\}$ 
11:     else break ▷ Timeout, unsupported, undecidable queries
12:   return  $\emptyset$ 

```

Algorithm 2 SPOTIT⁺

Require: Database \mathcal{S} , user query N , gold SQL query Q , Text-to-SQL frameworks \mathcal{M} , time limit T and bound K

Ensure: Counterexamples D_{cexs}

```

1: function SPOTIT+( $\mathcal{S}, N, \mathcal{M}, T, K$ )
2:    $D_{\text{cexs}} \leftarrow \emptyset$ 
3:   for  $m \in \mathcal{M}$  do
4:      $P \leftarrow m(\mathcal{S}, N)$  ▷ Generate SQL query  $P$  using  $m$ 
5:      $D_{\text{cexs}}[m] \leftarrow$  EQUIVCHECK( $\mathcal{S}, Q, P, T, K$ )
6:   ▷ Performing cross-referencing counterexamples
7:    $D_{\text{cexs}}^* \leftarrow \cup_{m \in \mathcal{M}} D_{\text{cexs}}[m]$ 
8:   for  $m \in \mathcal{M}$  do
9:     for  $D \in D_{\text{cexs}}^* \setminus D_{\text{cexs}}[m]$  do
10:      if  $\neg$ EX-TEST( $P, Q, D$ ) then
11:         $D_{\text{cexs}}[m] \leftarrow D_{\text{cexs}}[m] \cup \{D\}$ 
12:   return  $D_{\text{cexs}}$ 

```

Theorem 2 (Equivalence under set semantics). *Given two relations $R_1 = [t_1, \dots, t_n]$ and $R_2 = [r_1, \dots, r_m]$, if formula (2) is valid, then R_1 and R_2 are equivalent under set semantics.*

4.2 SPOTIT: A SEARCH-BASED TEXT-TO-SQL EVALUATION PIPELINE

Fig. 3 presents a high-level workflow of our approach that consists of three conceptual phases.

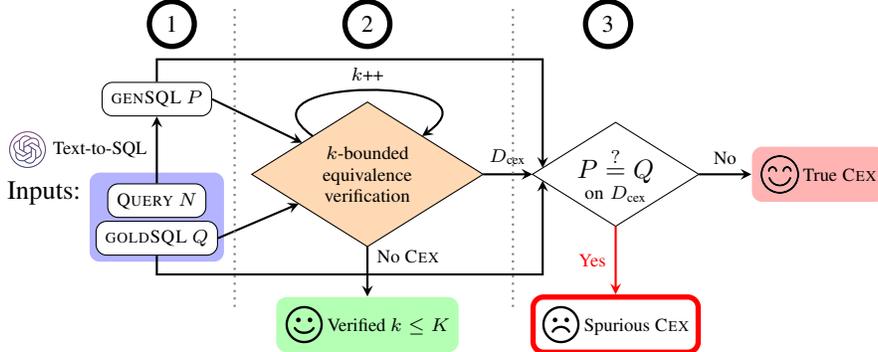


Figure 3: Three main phases of SPOTIT.

① **Input phase.** Given a NL question N and its corresponding gold SQL query Q , a Text-to-SQL framework takes as input N and generates a SQL query P . Both Q and P are passed to phase ②.

② **Verification phase.** The goal is to find a counterexample database instance on which the queries Q and P produce different outputs. For a given bound $k \leq K$, we perform bounded equivalence checking between Q and P . If the queries are proved equivalent, then we increase k by one for the next verification check. Furthermore, we cannot find any counterexample under all bounds and conclude that they are verified up to the bound k . On the other hand, if the queries are proved to be non-equivalent under some bound, we proceed to phase ③ for a further validation of D_{cex} .

③ **Validation phase.** Given the queries Q and P and a counterexample D_{cex} returned by verification algorithm, we must verify that this counterexample is non-spurious. There are two main reasons spurious counterexamples can arise in the verification engine. Either because some operators are over-approximated in the SMT encoding or the SQL query admits non-deterministic behaviors that cannot be modeled. Therefore, we execute the queries on the counterexample database (e.g., in SQLite) and check whether the results actually differ. D_{cex} is viewed valid if the results remain different; otherwise, we report this spurious case to the developers.

Alg 1 implements the second and third phases. For a given bound $k \leq K$, it first checks bounded equivalence between Q and P (line 3). If the queries are proved to be non-equivalent (line 6) under some bound, we validate that the counterexample database is indeed a true counterexample (line 9) and return it if this is the case. If the queries are proved to be equivalent in line 3, then we increase k by one for the next verification step. If the verifier cannot find any counterexample under all bounds,

Alg. 1 returns an empty set. Finally, if the verifier times out on a bound k , or the query is unsupported or undecidable, it also returns an empty set.

Cross-checking counter-examples. One observation we make is that as we progress through the frameworks, we collect a set of counterexamples that separate the gold query from the generated queries. Hence, we realized that these counterexamples can be reused as checks across all frameworks, as they might generalize across frameworks. Alg. 2 implements this idea. First, it obtains counterexample databases, if they exist, for all frameworks by calling Alg. 1 (lines 3-5). Then, it iterates over all frameworks again and tests equivalence between Q and P on these counterexample databases (lines 7-11). Empirically, this improves the effectiveness of our approach.

5 EXPERIMENTAL EVALUATION

In this section, we investigate the effect of using SPOTIT as the evaluation methodology for Text-to-SQL tasks. We are interested in the following questions:

- How much more SQL queries does our extension of VERIEQL support?
- Can SPOTIT provide more rigorous accuracy evaluation than test-based approaches?
- Can SPOTIT reveal shortcomings in existing Text-to-SQL evaluations?

Experimental Setup. We consider all 1,533 question-SQL pairs from the development set of BIRD (Li et al., 2023b), a state-of-the-art dataset for evaluating Text-to-SQL methods. The questions span 11 different databases from different professional domains, such as education, healthcare, and sports. The official BIRD leaderboard [1] contains over 80 Text-to-SQL methods and are updated frequently. Not all methods are open-source or have predictions publicly available. Therefore, we reached out to the developers of top-performing Text-to-SQL frameworks on the BIRD leaderboard and obtained the generated SQL queries for 10 of them, which constitutes a representative subset of state-of-the-art Text-to-SQL methods. The methods are listed in Tab. 1.

We first evaluate the predictions of each method using BIRD’s official test-based execution accuracy metric (EX-TEST), which, as described in Eq. 1, compares the results of executing the generated and gold queries on a given test database. For predictions that are deemed correct by EX-TEST, we apply SPOTIT to perform a more rigorous analysis. We implemented SPOTIT on top of VERIEQL (He et al., 2024), which we extended using the methods described in Sec. 4.1. To generate practically relevant counterexamples, we also extend the verification condition to exclude degenerate counterexamples that result in empty for one SQL and NULL for the other SQL.

We consider three variants of SPOTIT: (i) SPOTIT: Alg. 2 instantiated with the extended verification engine but without cross-checking (lines 7-11); (ii) SPOTIT⁻: Alg. 2 instantiated with vanilla VERIEQL and without cross-checking; (iii) SPOTIT⁺: Alg. 2 with cross-checking. We verify each generated-gold SQL pair up to a bound (K) of 5. Each verifier call is given one physical core, 8GB memory, and a CPU timeout of 600 seconds. In practice, a counterexample can typically be found within seconds, as reported below. Experiments were performed on a cluster equipped with Dell PowerEdge R6525 CPU servers featuring 2.6-GHz AMD CPU cores.

Performance of Verification Engine. We first evaluate the effect of our extensions to the original VERIEQL engine (He et al., 2024) in terms of *coverage*, defined as the fraction of generated-gold-SQL pairs that can be encoded into an SMT query. In addition, we measure the average runtime of SPOTIT on questions where a valid differentiating database is found. The results are shown in Tab. 2. Our extensions significantly increase the coverage of the verification engine on relevant questions (i.e., ones deemed correct by EX-TEST) for each method, allowing us to formally analyze a larger number of generated SQL queries. For example, for CSC-32B, the coverage increases from 84.83% to 94.88%, which corresponds to 110 additional supported questions $((94.88\% - 84.83\%) * 1094)$.

<https://bird-bench.github.io/>

Table 1: The Text-to-SQL methods we evaluated

Entry	Acronym
Alpha-SQL + Qwen2.5-Coder-32B (Li et al., 2025a)	ALPHA
CSC-SQL + Qwen2.5-Coder-7B (Sheng & Xu, 2025a)	CSC-7B
CSC-SQL + XiYanSQL (Sheng & Xu, 2025a)	CSC-32B
GenaSQL-1 (Dönder et al., 2025)	GENA-1
GenaSQL-2 (Dönder et al., 2025)	GENA-2
RSL-SQL + GPT-4o (Cao et al., 2024)	RSL
OmniSQL-32B (Li et al., 2025b)	OMNI-MAJ
GSR (anonymous authors)	GSR
CHESS _{IR+CG+UT} (Talaie et al., 2024a)	CHESS
SLM-SQL + Qwen2.5-Coder-1.5B (Sheng & Xu, 2025c)	SLM

Table 2: % of SQL pairs supported by SPOTIT⁻ and SPOTIT. For SPOTIT, also the average time in seconds on pairs where CEXs are found by the verifier and how many of them are non-spurious.

Method (# quest.)	SPOTIT ⁻ (%)	SPOTIT (%)	Avg. Time Valid. (s)	Valid. (%)
ALPHA (1064)	84.87	93.89	3.10	96.15
CHES (976)	87.40	97.13	1.40	93.34
CSC-32B (1094)	84.83	94.88	3.24	94.46
CSC-7B (1061)	85.77	96.14	3.93	95.10
GENA-1 (1062)	84.56	94.92	1.01	94.52
GENA-2 (1082)	84.47	94.55	0.93	95.42
GSR (1020)	84.51	93.63	1.12	94.86
OMNI-MAJ (1026)	86.65	95.61	1.36	95.83
RSL (1038)	86.03	95.18	1.64	95.62
SLM (973)	85.92	94.24	1.36	95.05

Table 3: Performance of Text-to-SQL methods using EX-TEST, EX-SPOTIT, and EX-SPOTIT⁺ on the 1533 BIRD-dev benchmarks.

	EX-TEST		EX-SPOTIT		EX-SPOTIT ⁺	
	Acc. (%)	Rank	Acc. (%)	Rank	Acc. (%)	Rank
CSC-32B	71.32	1	58.80	3	57.82	4
GENA-2	70.53	2	59.84	1	59.13	1
ALPHA	69.36	3	55.87	6	55.02	6
GENA-1	69.23	4	59.45	2	59.00	2
CSC-7B	69.17	5	58.54	4	57.95	3
RSL	67.67	6	56.58	5	55.80	5
OMNI-MAJ	66.88	7	54.69	7	54.04	7
GSR	66.49	8	54.56	8	53.72	8
CHES	63.62	9	52.87	9	52.35	9
SLM	63.43	10	51.37	10	50.98	10

The average time taken by SPOTIT to find a counterexample is under 4 seconds for all methods, which, combined with the fact that the analysis for each question can be done in parallel, confirms that SPOTIT is already a practical method for formally comparing generated SQLs with gold SQLs. Moreover, a high percentage (up to 96.15%) of the counterexamples found by the verifier are successfully validated, which suggests that our SMT encoding is sufficiently precise in practice.

Comparing test-based evaluation with SPOTIT. We now evaluate the accuracy of each Text-to-SQL method based on EX-TEST, EX-SPOTIT, and EX-SPOTIT⁺. As shown in Tab. 3, the accuracy of each method drops significantly when SPOTIT is used to check query equivalence. For example, the accuracy of CSC-32B drops from 71.32% to 58.80% with SPOTIT, and further to 57.82% when cross-checking is enabled. This means that there are 207 generated SQLs ($1533 * (71.32\% - 57.82\%)$) that passed the test on the official test databases, but were differentiated from the gold SQL by SPOTIT. Overall, SPOTIT resulted in a decrease in accuracy ranging from 9.8% to 13.5%, and cross-checking results in a small further decrease, by up to 1%. Interestingly, the ranking of the Text-to-SQL methods also changes substantially when evaluated under the more stringent, verification-based metrics, particularly in the top half of the table. For example, CSC-32B, which is ranked 1st by the official test-based metric, drops to 4th place when evaluated by SPOTIT⁺. And the 3rd place method ALPHA drops to the 6th place. These results indicate that test-based methods can in many cases overlook differences between the generated SQL and the gold SQL, which might lead to misrepresentation of the actual performance (both *absolute* and *relative*) of existing Text-to-SQL methods.

The effect of K. To study the effect of the choice of the verification bound K , we vary its value from 1 to 7 and run SPOTIT on the predictions of CSC-32B, the best model according to EX-TEST. As shown in Fig. 4, SPOTIT was able to find significantly more differentiating databases when K increases from 1 to 2 and the gain is marginal pass $K = 3$. This justifies our choice of $K = 5$.

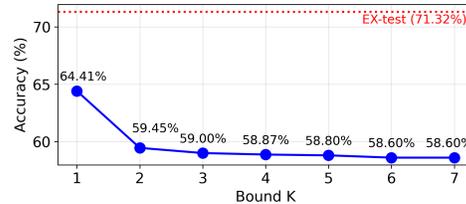


Figure 4: The effect of bound K .

Manual inspection of SPOTIT counterexamples. As SPOTIT performs bounded verification, the differentiating databases it finds are guaranteed to be minimal, which makes it easy to analyze them and understand the source of difference between the generated and gold SQLs. We manually examined the counterexamples for a random sample of 50 queries generated by CSC-32B and found that the difference between a generated SQL and a gold SQL can be primarily attributed to the three reasons that we described in Section 3: ambiguous question, incorrect gold SQL, and incorrect generated SQL. Fig. 5 shows a breakdown of the primary attributed reasons for those sampled questions. Surprisingly, while incorrect predictions do constitute a significant portion (26%), more often than not, the gold SQL itself is problematic. There are also a small fraction of cases (10%) where the question itself can be interpreted in multiple ways and therefore admits different answers. We discussed two

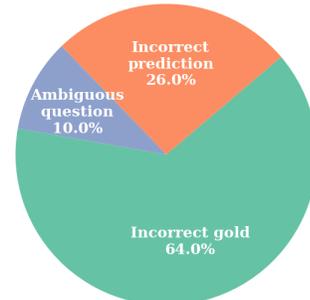


Figure 5: A breakdown of the primary reason for the difference between generated and gold SQLs.

examples of incorrect gold SQLs and ambiguous questions in Fig. I. Additional examples of each type of issues, along with the databases found by SPOTIT, are provided and discussed in App. D.

Additional analysis on Spider 2.0 benchmarks.

To assess the generalizability of SPOTIT to more complex Text-to-SQL tasks, we evaluate it on the recently introduced Spider 2.0 benchmark (Lei et al., 2024). We consider OMNISQL (Li et al., 2025b), a state-of-the-art Text-to-SQL method and GPT-5² on the 135 SQLite questions. These methods pass EX-TEST on 46 and 57 queries, respectively, which is competitive with the top entries on the Spider 2.0 leaderboard.³ As shown in Tab. 4, SPOTIT finds differentiating databases for 16 $((34.1\% - 22.2\%) * 135)$ and 8 query pairs deemed correct by test-based evaluation respectively for OMNISQL and GPT-5. SpotIt’s runtime for finding counterexamples remains low. We believe this is due to the fact that the schemas in Spider 2.0 are only moderately larger than those in BIRD (97.6 vs. 78.6 columns) and counter-examples can usually be detected with small values of K . The main verification challenge when it comes to Spider 2.0 lies in the number of SQL operators required for the queries but currently unsupported by our verification engine, which resulted in a smaller percentage of supported query pairs. Upon closer examination, 52.6% of the unsupported Spider 2.0 queries involve the window function (i.e., OVER clauses). Overall, these results indicate that SPOTIT is also useful for uncovering query discrepancies overlooked by test-based methods on the challenging Spider 2.0 benchmarks.

Summary of findings and implications. We now summarize the findings of our evaluation of a state-of-the-art Text-to-SQL evaluation dataset BIRD using SPOTIT and discuss their implications.

Finding 1: Existing test-based correctness metrics that involve executing the generated SQL and the gold SQL on static test databases can overlook significant variations in output data returned by the generated and gold SQLs. A search-based evaluation metric, such as SPOTIT, can serve as a practical alternative that provides additional perspectives on the performance of Text-to-SQL methods.

Finding 2: there is a significant number of problematic gold SQLs in existing Text-to-SQL benchmark sets. As shown by examples in Tab. I and App. D, in many cases, the issue can be hard to detect, yet can cause significantly different behaviors from the intended one. The presence of incorrect gold SQLs makes it hard to determine the true optimal performance on a benchmark set, as even a perfect Text-to-SQL method cannot achieve 100% accuracy.

Based on our result analysis for CSC-32B, we speculate that when most Text-to-SQL methods disagree with the gold SQL, the gold SQL is likely problematic. To validate this, we count the number of times that a prediction for a question is deemed correct by EX-TEST but incorrect by SPOTIT⁺ across all 10 Text-to-SQL methods. As shown in Fig. 6, there are 36 questions on which all methods generated queries that differ from the gold SQL. Manual inspection suggests that 31 of those 36 cases have problematic gold SQLs, 3 have ambiguous questions, and only 2 represent genuine errors in the generated SQLs.

While so far we have focused on incorrect gold SQLs overlooked by EX-TEST, our investigation begs the question: *when the generated query differs from the gold SQL, how often in general is the gold SQL problematic?* Fig. 7 shows the number of times the prediction for a question is deemed incorrect by EX-TEST across the 10 Text-to-SQL methods, for questions where CSC-32B’s predictions failed EX-TEST. There are 294 questions where at least 8 of the other 9 methods also failed EX-TEST. If shared disagreement with gold SQL is also a good indicator for problematic gold SQL in this case, then even a perfect Text-to-SQL method might

Table 4: Evaluation of SPOTIT on Spider 2.0.

	EX-TEST	EX-SPOTIT		
	Acc.(%)	Acc.(%)	Supported (%)	Avg. Time
OMNISQL	34.1	22.2	60.9%	3.4
GPT-5	42.2	36.3	50.9%	1.1

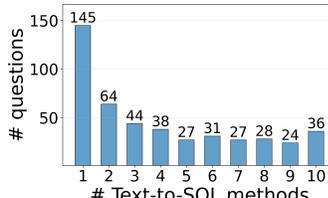


Figure 6: A breakdown of questions that passed EX-TEST but failed SPOTIT⁺.

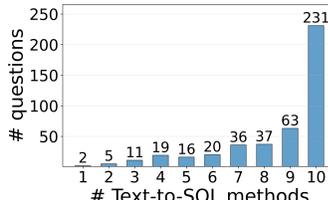


Figure 7: A breakdown of questions for which CSC-32B’s predictions failed EX-TEST.

²We use the same prompt as used in OMNISQL.

³<https://spider2-sql.github.io/>. We were not able to obtain predictions of Text-to-SQL methods on the leaderboard because they are predominantly closed source.

not be able to achieve an EX-TEST score much higher than 80% on BIRD-dev. As the time of completing this manuscript, the best EX-TEST score for BIRD-dev achieved by any method on the official leaderboard is 76.14%.

Large-scale benchmark sets inevitably contain problematic gold SQLs. Indeed, multiple sources have found examples of problematic gold SQLs in the BIRD dataset (Hui, 2024; Wretblad et al., 2024), and some of them have already been addressed by the maintainers. SPOTIT is the first approach that can provide minimal, easily analyzable databases to differentiate generated and gold SQLs, and can help to systematically uncover problematic gold SQLs.

Finding 3: A substantial number of questions in the Text-to-SQL dataset can be interpreted in different ways, thus admitting different SQL queries. While ambiguity is inherent in natural language, judging the correctness of a generated SQL query based on a single gold SQL query when the natural language question admits multiple interpretations might result in unfair penalization of Text-to-SQL methods.

Finding 4 (for the verification community): SMT-based equivalence verification techniques can already support a large fraction of practical SQL queries. Our results demonstrate that verification can often be completed within seconds. Due to the practical relevance of Text-to-SQL, we believe there is motivation for the verification community to invest more resources to precisely cover a larger fragment of SQL. In App. C, we discuss further extensions that would be especially useful according to our evaluation. Another significant next step to incorporate user preferences (potentially from natural languages) to search for particular types of counterexample databases. One way to achieve this is to encode the preferences as additional constraints in the SMT formulation of the verification problem.

6 RELATED WORK

Popular evaluation platforms such as BIRD-SQL (Li et al., 2024) and Spider 2.0 (Lei et al., 2024) evaluate query correctness by testing on predefined database instances. Several additional evaluations have been proposed to take into account partially correct generated queries (Pinna et al., 2025), efficiency of query executions (Zhang et al., 2024), and ambiguity in the questions (Li et al., 2023a; BIRD, 2025). However, the final correctness check is still via testing on a static database. Formal SQL equivalence checking broadly falls into two categories, full-fledged verification (Chu et al., 2017c; 2018; Zhou et al., 2022; 2024; Wang et al., 2024) and bounded verification (Veanes et al., 2010; Chu et al., 2017a; b; He et al., 2024). To the best of our knowledge, VERIEQL (He et al., 2024) supports the most expressive SQL fragments, while also offering extensibility for new features. We significantly extend the VERIEQL framework to support date and string types as well as a number of common operators for the Text-to-SQL evaluation task. Test data generation methods can also be useful for detecting query non-equivalence (Chandra et al., 2015; Somwase et al., 2024; Zhong et al., 2020). However, when the counterexamples require very specific structures, random fuzzing/testing can become unreliable in refuting equivalence. In contrast, SPOTIT systematically searches over the space of possible differentiating databases, finds minimal counter-examples, and provides a formal guarantee: if SPOTIT deems two SQL queries equivalent, then no counterexample of size less than a fixed number exists. A more detailed review of related work can be found in App. A.

7 CONCLUSION

We presented SPOTIT, the first verification-based evaluation pipeline for Text-to-SQL. We introduced techniques to support a richer SQL grammar, which enabled us to efficiently analyze a large fragment of SQL queries commonly seen in Text-to-SQL tasks. Our initial motivation for developing SPOTIT was to examine the extent to which the accuracy of a Text-to-SQL method is overestimated by test-based evaluation, which is widely adopted as the default metric on high-profile Text-to-SQL evaluation platforms. However, a closer inspection of the verification results revealed a far more complex picture. While SPOTIT can indeed detect incorrect generated SQL queries that were overlooked by test-based methods, a significant portion of the inconsistency between the gold and generated SQLs can be explained by the benchmarks themselves—either due to problematic gold SQLs or due to ambiguous natural language questions. We discussed the implications of and the next steps from our findings, and hope that our work will motivate further work on evaluating and improving Text-to-SQL evaluation frameworks.

REFERENCES

- Amazon. Build a text-to-sql solution for data consistency in generative ai using amazon nova, September 2025. URL <https://aws.amazon.com/blogs/machine-learning/build-a-text-to-sql-solution-for-data-consistency-in-generative-ai-using-amazon-nova/>.
- Amazon Web Services. How msd uses amazon bedrock to translate natural language into sql for complex healthcare databases. <https://aws.amazon.com/blogs/machine-learning/how-merck-uses-amazon-bedrock-to-translate-natural-language-into-sql-for-complex-healthcare-databases/>, 2024. Accessed: 2025-09-21.
- Clark Barrett and Cesare Tinelli. Satisfiability modulo theories. In *Handbook of model checking*, pp. 305–343. Springer, 2018.
- Team BIRD. Bird-interact: Re-imagine text-to-sql evaluation via lens of dynamic interactions. <https://github.com/bird-bench/BIRD-Interact>, 2025. Accessed: 2025-06-01.
- Team BitsAI. Bits ai:scale your teams with autonomous ai agents across monitoring, development, and security, September 2025. URL <https://www.datadoghq.com/product/platform/bits-ai/>.
- Zhenbiao Cao, Yuanlei Zheng, Zhihao Fan, Xiaojin Zhang, Wei Chen, and Xiang Bai. Rsl-sql: Robust schema linking in text-to-sql generation. *arXiv preprint arXiv:2411.00073*, 2024.
- Bikash Chandra, Bhupesh Chawda, Biplab Kar, KV Maheshwara Reddy, Shetal Shah, and S Sudarshan. Data generation for testing and grading sql queries. *The VLDB Journal*, 24(6):731–755, 2015.
- Shuaichen Chang and Eric Fosler-Lussier. How to prompt llms for text-to-sql: A study in zero-shot, single-domain, and cross-domain settings. *arXiv preprint arXiv:2305.11853*, 2023.
- Shumo Chu, Daniel Li, Chenglong Wang, Alvin Cheung, and Dan Suciu. Demonstration of the cosette automated sql prover. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pp. 1591–1594, 2017a.
- Shumo Chu, Chenglong Wang, Konstantin Weitz, and Alvin Cheung. Cosette: An automated prover for sql. In *CIDR*, pp. 1–7, 2017b.
- Shumo Chu, Konstantin Weitz, Alvin Cheung, and Dan Suciu. Hottsql: Proving query rewrites with univalent sql semantics. *Acm sigplan notices*, 52(6):510–524, 2017c.
- Shumo Chu, Brendan Murphy, Jared Roesch, Alvin Cheung, and Dan Suciu. Axiomatic foundations and algorithms for deciding semantic equivalences of sql queries. *arXiv preprint arXiv:1802.02229*, 2018.
- Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 337–340. Springer, 2008.
- Yusuf Denizay Dönder, Derek Hommel, Andrea W Wen-Yi, David Mimno, and Unso Eun Seo Jo. Cheaper, better, faster, stronger: Robust text-to-sql without chain-of-thought or fine-tuning. *arXiv preprint arXiv:2505.14174*, 2025.
- Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou. C3: Zero-shot text-to-sql with chatgpt, 2023.
- Yingqi Gao, Yifu Liu, Xiaoxia Li, Xiaorong Shi, Yin Zhu, Yiming Wang, Shiqi Li, Wei Li, Yuntao Hong, Zhiling Luo, Jinyang Gao, Liyu Mou, and Yu Li. A preview of xiyan-sql: A multi-generator ensemble framework for text-to-sql, 2025. URL <https://arxiv.org/abs/2411.08599>.
- Yang He, Pinhan Zhao, Xinyu Wang, and Yuepeng Wang. Verieql: Bounded equivalence verification for complex sql queries with integrity constraints. *Proc. ACM Program. Lang.*, 8(OOPSLA1), April 2024. doi: 10.1145/3649849. URL <https://doi.org/10.1145/3649849>.

- Binyuan Hui. Panel of BIRD Annotation Issues. <https://github.com/AlibabaResearch/DAMO-ConvAI/issues/39>, 2024. Accessed: May 2024.
- Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin Su, Zhaoqing Suo, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, et al. Spider 2.0: Evaluating language models on real-world enterprise text-to-sql workflows. *arXiv preprint arXiv:2411.07763*, 2024.
- Boyan Li, Jiayi Zhang, Ju Fan, Yanwei Xu, Chong Chen, Nan Tang, and Yuyu Luo. Alpha-sql: Zero-shot text-to-sql using monte carlo tree search. *arXiv preprint arXiv:2502.17248*, 2025a.
- Haoyang Li, Shang Wu, Xiaokang Zhang, Xinmei Huang, Jing Zhang, Fuxin Jiang, Shuai Wang, Tieying Zhang, Jianjun Chen, Rui Shi, et al. Omnisql: Synthesizing high-quality text-to-sql data at scale. *arXiv preprint arXiv:2503.02240*, 2025b.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiayi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin Chen-Chuan Chang, Fei Huang, Reynold Cheng, and Yongbin Li. Can LLM already serve as A database interface? A big bench for large-scale database grounded text-to-sqls. *CoRR*, abs/2305.03111, 2023a. doi: 10.48550/ARXIV.2305.03111. URL <https://doi.org/10.48550/arXiv.2305.03111>.
- Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36:42330–42357, 2023b.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiayi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin Chen-Chuan Chang, Fei Huang, Reynold Cheng, and Yongbin Li. BIRD_SQ. <https://bird-bench.github.io/>, 2024. Accessed: Sep 2025.
- Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S. Yu. A comprehensive evaluation of chatgpt’s zero-shot text-to-sql capability, 2023.
- Yifu Liu, Yin Zhu, Yingqi Gao, Zhiling Luo, Xiaoxia Li, Xiaorong Shi, Yuntao Hong, Jinyang Gao, Yu Li, Bolin Ding, and Jingren Zhou. Xiyan-sql: A novel multi-generator framework for text-to-sql, 2025. URL <https://arxiv.org/abs/2507.04701>.
- Mudathir Mohamed, Andrew Reynolds, Cesare Tinelli, and Clark Barrett. Verifying sql queries using theories of tables and relations. In *Proceedings of 25th Conference on Logic for Pro*, volume 100, pp. 445–463, 2024.
- Giovanni Pinna, Yuriy Perekhoshin, Luca Manzoni, Mauro Castelli, and Andrea De Lorenzo. Redefining text-to-sql metrics by incorporating semantic and structural similarity. *Scientific Reports*, 15, July 2025. ISSN 2045-2322. doi: 10.1038/s41598-025-04890-9.
- Juan Sequeda, Dean Allemang, and Bryon Jacob. A benchmark to understand the role of knowledge graphs on large language model’s accuracy for question answering on enterprise sql databases, 2023.
- Lei Sheng and Shuai-Shuai Xu. Csc-sql: Corrective self-consistency in text-to-sql via reinforcement learning. *arXiv preprint arXiv:2505.13271*, 2025a.
- Lei Sheng and Shuai-Shuai Xu. Csc-sql: Corrective self-consistency in text-to-sql via reinforcement learning, 2025b. URL <https://arxiv.org/abs/2505.13271>.
- Lei Sheng and Shuai-Shuai Xu. Slm-sql: An exploration of small language models for text-to-sql. *arXiv preprint arXiv:2507.22478*, 2025c.
- Liang Shi, Zhengju Tang, Nan Zhang, Xiaotong Zhang, and Zhi Yang. A survey on employing large language models for text-to-sql tasks. *ACM Computing Surveys*, 2024.
- Vladislav Shkapenyuk, Divesh Srivastava, Theodore Johnson, and Parisa Ghane. Automatic metadata extraction for text-to-sql, 2025. URL <https://arxiv.org/abs/2505.19988>.

- Sunanda Somwase, Parismita Das, and S Sudarshan. Data generation for testing complex queries. *arXiv preprint arXiv:2409.18821*, 2024.
- Splunk. Ai assistant in observability cloud, September 2025. URL https://www.splunk.com/en_us/products/splunk-ai-assistant-in-observability-cloud.html
- Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. Chess: Contextual harnessing for efficient sql synthesis. *arXiv preprint arXiv:2405.16755*, 2024a.
- Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. Chess: Contextual harnessing for efficient sql synthesis, 2024b. URL <https://arxiv.org/abs/2405.16755>.
- TiDBCloud. Get Started with Chat2Query API. <https://docs.pingcap.com/tidbcloud/use-chat2query-api>, 2020. Accessed: May 2024.
- Margus Veanes, Nikolai Tillmann, and Jonathan De Halleux. Qex: Symbolic sql query explorer. In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*, pp. 425–446. Springer, 2010.
- Shuxian Wang, Sicheng Pan, and Alvin Cheung. Qed: A powerful query equivalence decider for sql. *Proceedings of the VLDB Endowment*, 17(11):3602–3614, 2024.
- Niklas Wretblad, Fredrik Gordh Riseby, Rahul Biswas, Amin Ahmadi, and Oskar Holmström. Understanding the effects of noise in text-to-sql: An examination of the bird-bench benchmark, 2024.
- Bohan Zhai, Canwen Xu, Yuxiong He, and Zhewei Yao. Excot: Optimizing reasoning for text-to-sql with execution feedback, 2025. URL <https://arxiv.org/abs/2503.19988>.
- Bin Zhang, Yuxiao Ye, Guoqing Du, Xiaoru Hu, Zhishuai Li, Sun Yang, Chi Harold Liu, Rui Zhao, Ziyue Li, and Hangyu Mao. Benchmarking the text-to-sql capability of large language models: A comprehensive evaluation, 2024. URL <https://arxiv.org/abs/2403.02951>.
- Ruiqi Zhong, Tao Yu, and Dan Klein. Semantic evaluation for text-to-sql with distilled test suites. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 396–411, 2020.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. Solving challenging math word problems using GPT-4 code interpreter with code-based self-verification. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=c8McWs4Av0>.
- Qi Zhou, Joy Arulraj, Shamkant B Navathe, William Harris, and Jinpeng Wu. Spes: A symbolic approach to proving query equivalence under bag semantics. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pp. 2735–2748. IEEE, 2022.