

Constant Optimization Driven Database System Testing

CHI ZHANG, Nanjing University, China

MANUEL RIGGER, National University of Singapore, Singapore

Logic bugs are bugs that can cause database management systems (DBMSs) to silently produce incorrect results for given queries. Such bugs are severe, because they can easily be overlooked by both developers and users, and can cause applications that rely on the DBMSs to malfunction. In this work, we propose *Constant-Optimization-Driven Database Testing* (CDDTest) as a novel approach for detecting logic bugs in DBMSs. This method draws inspiration from two well-known optimizations in compilers: constant folding and constant propagation. Our key insight is that for a certain database state and query containing a predicate, we can apply constant folding on the predicate by replacing an expression in the predicate with a constant, anticipating that the results of this predicate remain unchanged; any discrepancy indicates a bug in the DBMS. We evaluated CDDTest on five mature and extensively-tested DBMSs—SQLite, MySQL, CockroachDB, DuckDB, and TiDB—and found 45 unique, previously unknown bugs in them. Out of these, 24 are unique logic bugs. Our manual analysis of the state-of-the-art approaches indicates that 11 logic bugs are detectable only by CDDTest. We believe that CDDTest is easy to implement, and can be widely adopted in practice.

CCS Concepts: • **Information systems** → **Database query processing**; • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: DBMSs testing, logic bugs, test oracle

ACM Reference Format:

Chi Zhang and Manuel Rigger. 2025. Constant Optimization Driven Database System Testing. *Proc. ACM Manag. Data* 3, 1 (SIGMOD), Article 24 (February 2025), 24 pages. <https://doi.org/10.1145/3709674>

1 Introduction

Database management systems (DBMSs) store data for numerous essential software systems. Similar to other software systems, DBMSs are susceptible to a range of bugs, including logic bugs. Logic bugs can be critical as they can cause the DBMSs to silently produce incorrect results for a given query and might be overlooked by both developers and users.

In recent years, various approaches have been proposed to find logic bugs in DBMSs. The state-of-the-art approaches are NoREC [30], TLP [31], PQS [32], DQE [35], and TQS [37]. NoREC and DQE assume that the same predicate should, for a given row, consistently evaluate to the same value, regardless of which clause it is used in. NoREC applies this to **WHERE** clauses in **SELECT**, while DQE applies this to **WHERE** clauses in **UPDATE** and **DELETE**. TLP decomposes a query into three partitioning queries, each retrieving rows based on predicates p , **NOT** p , and p **IS NULL**, respectively. All these approaches are black-box techniques, as they apply on the SQL level, not on the DBMSs' source code. Because of this, and since automated testing approaches typically lack guarantees, it is difficult to precisely characterize what bugs they find and which ones they miss. For one, none of

Authors' Contact Information: Chi Zhang, State Key Laboratory for Novel Software Technology, School of Computer Science, Nanjing University, Nanjing, China, zhangchi_seg@smail.nju.edu.cn; Manuel Rigger, National University of Singapore, Singapore, rigger@nus.edu.sg.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 2836-6573/2025/2-ART24
<https://doi.org/10.1145/3709674>

these approaches support testing subqueries, an important feature that allows query nesting, and it is not obvious how these approaches could be extended to support it. For another, as the results of this paper show, these approaches overlook various bugs due to their inherent limitations.

To tackle logic bugs in DBMSs, we propose a novel testing methodology that tests DBMSs through the lens of constant folding and constant propagation, two optimizations that were originally proposed for compilers. Specifically, we propose *Constant-Optimization-Driven Database Testing (CDDTest)*,¹ which is a black-box approach that, for a given database state and query, applies constant folding and propagation to expressions in the query. The *test oracle*, that is, the mechanism for validating the behavior of the DBMS for a given input, relies on validating that the transformed query produces the same result as the original query.

Constant folding is a well-known compiler optimization that evaluates constant expressions at compile time, rather than computing them at run time [27]. For example, it evaluates the statement $i = 1 + 2 + 3$; to $i = 6$; at compile time. Constant folding extends beyond numeric literal expressions by leveraging constant propagation. Constant propagation, through reachability analysis [11], determines constant values for variables by assessing their reachability at specific program points. By employing constant propagation in conjunction with constant folding, more intricate programs can be optimized. For example, the statement sequence $a = 1$; $i = a + 2 + 3$; can be optimized to $i = 6$;

In the context of DBMSs, the result of an expression in a query is deterministic for any given row under a specific database state. While, as an optimization, applying constant folding and propagation assuming a constant database state would be ineffective, since the state typically frequently changes, we can leverage these optimizations for a test oracle. Listing 1 illustrates this by using a bug-inducing test case that enabled us to find a subquery-related bug in SQLite.² We start from an initial database state, which includes a table t_0 with index i_0 , as well as a view v_0 . The query ①, which corresponds to the *original query*, uses a subquery in the **WHERE** clause, which we want to constant-fold. We can do so by extracting the subquery, using an *auxiliary query*, as shown in query ②, and executing it using the DBMS under test, which returns \emptyset . We derive the *folded query* as shown in query ③, by, in query ①, replacing the subquery with its result obtained in query ②. Unexpectedly, SQLite returned a different result for the original and folded queries; for query ①, SQLite returned 1, while it returned \emptyset for query ③. This discrepancy indicates a bug in the DBMS. We received feedback from the developers of SQLite, stating that this bug was caused by a query planner optimization, and an aggregate subquery is a necessary condition to trigger this bug. We found this bug, since our folded query (*i.e.*, query ③) no longer met this requirement. Our approach is akin to a combination of constant folding and constant propagation. In this example, using the DBMS to evaluate the expression under a certain database state can be seen as constant folding (*i.e.*, in query ②), while the substitution of the entire expression with its result can be viewed as constant propagation (*i.e.*, replacing the subquery in query ① with the constant \emptyset).

Our approach is effective in detecting logic bugs, because it can directly transform a predicate, resulting in the execution of different code within DBMSs, while the folded query should produce the same results as the original one. Any discrepancy between the queries indicates a bug in DBMSs. Our approach is capable of identifying bugs that occur simultaneously in all three **SELECT**, **UPDATE**, and **DELETE** statements, which will be missed by NoREC and DQE. Our approach can detect bugs where boolean predicates consistently evaluate to an opposite value, which would be missed by TLP. Our transformation-based oracle can handle complex queries, whereas PQS is limited by the ability of the self-implemented interpreter. Overall, we believe that CDDTest complements

¹This name pays homage to Edgar Frank Codd [12], who proposed the relational model for database management.

²<https://sqlite.org/forum/info/a68313d0545273c8>

Listing 1. An illustrative example, which triggered a bug in SQLite. The expression in the original query subject to constant folding and propagation, along with the corresponding constant in the folded query, are highlighted in red.

```
CREATE TABLE t0 (c0);
INSERT INTO t0(c0) VALUES (1);
CREATE INDEX i0 ON t0(c0 > 0);
CREATE VIEW v0(c0) AS SELECT AVG(t0.c0) FROM t0 GROUP BY 1 > t0.c0;
ⓐ SELECT COUNT(*) FROM t0 INDEXED BY i0 WHERE (SELECT COUNT(*) FROM v0 WHERE v0.c0
    BETWEEN 0 AND 0); -- 1 ✖
ⓑ SELECT COUNT(*) FROM v0 WHERE v0.c0 BETWEEN 0 AND 0; -- 0
ⓒ SELECT COUNT(*) FROM t0 INDEXED BY i0 WHERE 0; -- 0 ✓
```

Listing 2. A correlated subquery example.

```
CREATE TABLE t0(ID INT, score INT, classID INT);
INSERT INTO t0 VALUES (0, 90, 1), (1, 80, 1), (2, 83, 2);
ⓐ SELECT x.ID FROM t0 AS x WHERE x.score >
    (SELECT AVG(y.score) FROM t0 AS y WHERE x.classID = y.classID);
ⓑ SELECT x.classID,
    (SELECT AVG(y.score) FROM t0 AS y WHERE x.classID = y.classID) FROM t0 AS x;
ⓒ SELECT x.ID FROM t0 AS x WHERE x.score >
    (CASE WHEN x.classID = 1 THEN 85
        WHEN x.classID = 1 THEN 85
        WHEN x.classID = 2 THEN 83 END);
```

existing testing approaches for DBMSs, and provides a new conceptual angle on using compiler optimizations for DBMS testing.

We integrated CODDTest into SQLancer,³ a widely-used tool for testing DBMSs, in which also other state-of-the-art oracles have been integrated, and tested five mature DBMSs with it, all of which have been extensively tested by a number of state-of-the-art testing approaches. The results were surprisingly positive. We found 45 previously unknown, unique bugs. Of these, 12 were confirmed, 33 were fixed. 24 out of these were previously unknown unique logic bugs, which we aimed to find. Of these logic bugs, 5 were confirmed, and 19 were fixed. In summary, we make the following contributions:

- We propose the novel idea of testing DBMSs through the lens of constant propagation and constant folding.
- We provide a realization of this idea by, given an *original query*, deriving a *folded query*, in which an expression has been constant-folded based on the result of an *auxiliary query* to validate whether the DBMS computes a consistent result for the original and folded queries.
- We implemented the approach and evaluated it on five well-tested and mature DBMSs, uncovering 24 unique, previously unknown logic bugs. Furthermore, we provide a comparative analysis with the state-of-the-art approaches, as well as a performance comparison.

2 Background

Predicates. In SQL, a predicate is a boolean expression that evaluates to **TRUE**, **FALSE**, or **NULL** when applied to given values or rows. Predicates are used in various clauses of SQL, such as the **WHERE** clauses of **SELECT**, **UPDATE**, and **DELETE**, as well as the **JOIN ON**, **HAVING**, **GROUP BY**, and **ORDER BY** clauses of **SELECT**.

³<https://github.com/sqlancer/sqlancer/pull/1054>

Listing 3. A CASE expression example.

```

CREATE TABLE grade (score INT);
INSERT INTO grade(score) VALUES (100), (80), (60);
SELECT score, CASE WHEN score = 100 THEN "A"
                  WHEN score >= 80 AND score < 100 THEN "B"
                  ELSE "C" END FROM grade;

```

Subqueries. Subqueries are an important language feature in SQL, and can be used in predicates. Subqueries can be classified as correlated or non-correlated based on whether the subquery references columns from the outer query. Correlated subqueries are **SELECT** queries nested within outer queries, referencing columns from the outer queries to construct their predicate. A correlated subquery will be executed once for each row passed from the outer query. Consider the query ② depicted in Listing 2, where a correlated subquery is used to identify students (*i.e.*, represented by **ID**) whose scores exceed the average score in the class. The subquery calculates the average score of the class, whose **classID** is passed from the outer query. Therefore, for each student in t_0 , the subquery will be executed once to derive the average score in the class of that particular student. In contrast, non-correlated subqueries do not reference columns in the outer query and are evaluated once in execution. Various optimizations have been developed to enhance the performance of subqueries, such as subquery unnesting [6, 9] or correlated subquery decorrelation [33]. Like other language features, subquery support in DBMSs is susceptible to logic bugs.

SQL CASE expression. The **CASE** expression is a feature of SQL to process **if/then** logic. Listing 3 illustrates an example using the **CASE** expression. The predicate following the **WHEN** keyword determines whether the value after the **THEN** keyword should be returned. If all predicates evaluate to false, the result of this **CASE** expression is the value specified in the **ELSE** clause, or **NULL** if there is no **ELSE** clause. In our approach, we use the **CASE** expression to map each row of a table to constant values, corresponding to constant folding on a per-row basis.

Metamorphic testing. Metamorphic testing [10] is an automated testing methodology, which generates a new input for the program under test, based on an existing input and its corresponding output, with the expectation that the outcome of this new input can be predicted. Formally, given an input I and $P(I) = O$, where P is the program under test, a follow-up input I' is derived, so that a known relationship between O and $P(I') = O'$ is validated. Metamorphic testing has been applied in testing numerous critical systems, such as compilers [22], SMT solvers [38], and DBMSs. NoREC [30], TLP [31], and DQE [35] are all metamorphic testing approaches for detecting logic bugs in DBMSs. The core challenge of realizing a metamorphic testing approach is to identify a so-called metamorphic relation that derives the follow-up input I' and relates the outputs O and O' so that bugs in the system under test can be revealed.

3 Approach

We propose CODDTest as an approach to finding logic bugs in DBMSs through the lens of constant folding and propagation. Our key insight is that within an SQL query, by assuming a constant database state and given query, we can apply constant folding and constant propagation to a specific expression in a predicate, assuming that the query's result remains unchanged. More precisely, our approach for testing DBMSs generates two equivalent queries: the *original query*, which includes the randomly generated predicate, and the *folded query*, which is derived by substituting the predicate in the original query with the corresponding constant-folded predicate. Any discrepancy in the results of these two queries indicates a potential bug.

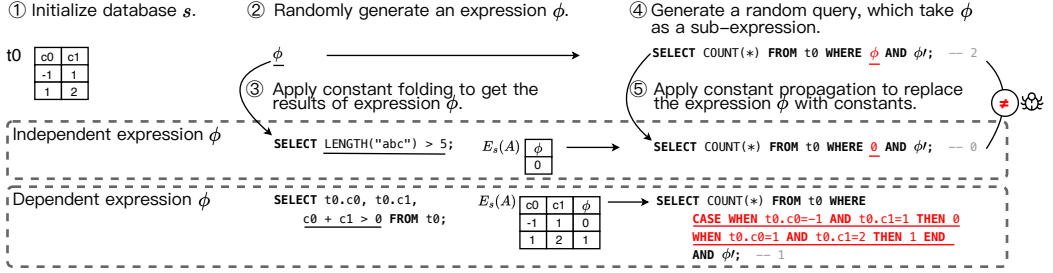


Fig. 1. Overview of approach. CDDTest generates pairs of equivalent queries by applying constant folding and propagation to the expression ϕ . The application of constant propagation and folding differs for independent and dependent expressions.

Metamorphic relation. CDDTest is a metamorphic testing approach, where the folded query F is derived from the original query O , and the results of these two queries with DBMS engine E are expected to be identical. More formally, for the potential database state space denoted by \mathcal{S} , we can establish the metamorphic relationship below. Any violation indicates a bug in E .

$$\forall s \in \mathcal{S} \quad E_s(O) = E_s(F)$$

where $E_s(O)$ denotes the outcome of executing a query O with the DBMS engine E under the state s , while $E_s(F)$ represents the result of F .

The generation of F involves applying constant folding and constant propagation to a randomly selected expression ϕ within the predicate of O , represented as $F = O[\phi/R_\phi]$, which includes substituting ϕ with the expression R_ϕ within the same query location in the query O . R_ϕ is the evaluated result of ϕ , and this process referred to as constant propagation. R_ϕ can be obtained through constant folding, which entails evaluating ϕ under the database s , denoted by $R_\phi = E_s(A[\phi])$. Here, $A[\phi]$ represents an auxiliary query A designed to obtain the evaluation result of ϕ . Finally, the above metamorphic relation can be elaborated as follows:

$$\begin{aligned} \forall s \in \mathcal{S} \quad R_\phi &= E_s(A[\phi]) \\ E_s(O) &= E_s(F) = E_s(O[\phi/R_\phi]) \end{aligned}$$

Approach overview. Figure 1 illustrates our approach. In step ①, we initialize the database and create non-empty tables. We do this randomly by using rule-based generators, such as also done by other testing works [1, 30]. Non-empty tables ensure that at least one row is available for us to apply constant folding. Here, we generate a table t_0 with two rows. Subsequently, we enter a loop from step ② to step ⑤ to thoroughly test the generated database state. Each iteration is designed to test the DBMS once. Step ② generates a random expression ϕ , which will undergo constant folding. In step ③, we introduce the *auxiliary query* to retrieve the corresponding constants of ϕ with respect to the database state—we see this step as constant folding. In step ④, we generate the original query, using ϕ as part of a predicate. We execute the original query using the DBMS under test and obtain its results (i.e., 2). In step ⑤, we apply constant propagation to ϕ by replacing it with the corresponding constants to obtain the folded query. If the folded query produces a different result than the original query, we have found a bug-inducing test case.

Applying constant folding (i.e., step ③) and constant propagation (i.e., step ⑤) differs based on whether the expression ϕ is a so-called *independent expression* or *dependent expression*. Independent expressions yield constant results irrespective of the outer context, allowing us to perform constant folding and propagation on the expression ϕ by replacing it with a constant or a constant list. For

Algorithm 1: Algorithm of CODDTest

```

1 function TestOracleGen(DatabaseState s)
    // Randomly generate an expression  $\phi$ , which will undergo constant folding and constant
    // propagation. We extract the set of the referenced columns  $\{c_i\}$  in  $\phi$ , which come from
    // outer context, along with the tables set  $\{t_i\}$  to which these columns  $c_i$  are associated.
2    $\phi, \{c_i\}, \{t_i\} \leftarrow \text{GenExpr}(s)$ 
    // Constant folding of  $\phi$  under  $s$ , this step differs based on  $\phi$  is dependent or
    // independent expressions. Specifically,  $\phi$  is considered an independent expression when
    //  $\{c_i\}$  is empty; otherwise,  $\phi$  is classified as a dependent expression.
3   if Size( $\{c_i\}$ ) == 0 then
    // Construct the auxiliary query for independent expression
4      $A \leftarrow \text{"SELECT"} + \phi$ 
5      $E_s(A) \leftarrow \text{ExecQuery}(A, s)$ 
    // Transform the constant result of  $\phi$  as a constant expression  $R_\phi$ 
6      $R_\phi \leftarrow E_s(A)$ 
7   else
    // Dependent expression has different result for each row of  $\{c_i\}$ 
8      $A \leftarrow \text{"SELECT"} + (\{c_i\}, \phi) + \text{"FROM"} + \{t_i\}$ 
9      $E_s(A) \leftarrow \text{ExecQuery}(A, s)$ 
    // Map the results of  $\phi$  to each row of  $\{c_i\}$  as an expression  $R_\phi$ 
10     $R_\phi \leftarrow \text{Map}(E_s(A))$ 
    // Generate the original query based on the current database state, using  $\phi$  as a
    // sub-expression in predicate
11     $O \leftarrow \text{QueryGenerate}(s, \phi, \{t_i\})$ 
12     $E_s(O) \leftarrow \text{ExecQuery}(O, s)$ 
    // Generate the folded query by replacing  $\phi$  with  $R_\phi$ 
13     $F \leftarrow \text{ReplaceExpr}(\phi, R_\phi, O)$ 
14     $E_s(F) \leftarrow \text{ExecQuery}(F, s)$ 
    // A bug is identified if there is a discrepancy between the results of the original
    // query and the folded query
15    if  $E_s(O) \neq E_s(F)$  then
16      ReportBug( $O, F, A$ )

```

example, the independent expression ϕ shown in Figure 1, $\text{LENGTH}(\text{"abc"}) > 5$, does not reference any columns and evaluates to **FALSE**, regardless of the database state and the clauses it is used in. Similarly, the constant-folded expression (*i.e.*, query \textcircled{A}) in Listing 1 is an independent expression, as the non-correlated subquery can be executed independently and yields constant results regardless of the outer query. We categorize expressions as dependent expressions if they reference columns from an outer context and cannot be executed independently. For example, the dependent expression $c_0 + c_1 > 0$ shown in Figure 1 references columns c_0 and c_1 from the **FROM** clause; the original query shown in Listing 2 is also a dependent expression, as the correlated subquery depends on the value of $x.\text{classID}$ from the outer query. For dependent expressions, we must consider that ϕ might evaluate to a different value for every row, which is why we represent the folded constant as a mapping from a row value to the constant.

Algorithmic sketch. We use Algorithm 1 to demonstrate the process of generating a test oracle. The function TestOracleGen corresponds to step ② to step ⑤ in Figure 1, taking the randomly generated database state s as its input.

For step ②, we randomly generate the expression ϕ with function `GenExpr`, based on the current database state s , as outlined in line 2. ϕ then undergoes constant folding and constant propagation. `GenExpr` returns two sets: $\{c_i\}$, which consists of referenced columns from the outer context and is utilized to determine if ϕ is an independent expression, and $\{t_i\}$, which denotes the tables to which c_i belong. For example, for the independent expression `LENGTH("abc") > 5`, `GenExpr` yields two empty sets for $\{c_i\}$ and $\{t_i\}$ as results, due to the absence of any referenced columns in this expression. Conversely, for the dependent expression `c0 + c1 > 0` shown in Figure 1, `GenExpr` returns $\{c0, c1\}$, and $\{t0\}$.

The logic for constant folding (*i.e.*, step ③) spans from line 3 to line 10. When the size of $\{c_i\}$ is zero, ϕ does not reference any columns from the outer context, leading to ϕ producing constant results regardless of the outer context. Therefore, ϕ is considered an independent expression. The value of an independent expression can be obtained through an auxiliary query consisting solely of the `SELECT` keyword followed by ϕ , as shown in line 4. Moreover, this `SELECT` keyword can be omitted when ϕ is a non-correlated subquery. The outcome of an independent expression is either a constant (the empty result can be considered as `NULL`) or a constant list. Once we ascertain the result of the auxiliary query (*i.e.*, line 5), we can directly convert this result into a constant expression in line 6. For a dependent expression, which produces different outcomes for various rows of $\{c_i\}$, we must construct R_ϕ to reflect this variability. In line 8, by placing $\{c_i\}$ and ϕ together following the `SELECT` keyword, we obtain the result of ϕ , as well as the values of $\{c_i\}$ that lead to the result of ϕ . To ensure syntactic correctness, we also need to append the list of tables $\{t_i\}$ after the `FROM` keyword. Finally, we construct R_ϕ as a mapping from the value of each row in $\{c_i\}$ to the corresponding result of ϕ , as shown in line 10. In the subsequent two subsections, we elaborate and further illustrate how we handle both independent and dependent expressions.

Step ④ is accomplished in line 11 through the `QueryGenerate` function, a random query generator that takes the current database state s , the expression ϕ , and the table set $\{t_i\}$ as input. ϕ is then randomly incorporated into a predicate of the generated query. It is imperative to include the table set $\{t_i\}$ as an input. This requirement arises from the need for auxiliary queries to replicate the original query's `JOIN` clauses, with the sole exception being instances where ϕ functions as the predicate within a `JOIN` clause. The set $\{t_i\}$ records the information of the `JOIN` clauses associated with each table. For a deeper understanding of the reason, we provide the explanation in Section 3.2, in which we discuss how to perform constant folding on dependent expressions.

Constant propagation (*i.e.*, step ⑤) is applied in line 13 using the `ReplaceExpr` function, which replaces the expression ϕ with the expression R_ϕ —representing the result of ϕ , in the corresponding location within query O . After obtaining the result of the original query O in line 12 and the result of the folded query in line 14, we can identify the presence of a bug if any discrepancies exist between these two results (*i.e.*, line 15 to line 16).

3.1 Folding Independent Expressions

Independent expressions can be executed without surrounding contexts and can have two possible shapes. First, if the expression ϕ has no column references, it is a constant expression that always yields a constant value. In such cases, in the auxiliary query, we can evaluate the expression using a `SELECT` statement. For example, the independent expression `LENGTH("abc") > 5` shown in Figure 1 is used in a `SELECT` statement to derive its results. Second, the expression ϕ can be a non-correlated subquery, which computes a constant result assuming a fixed database state. For example, in Listing 1, the subquery of query ① returns the same result regardless of the outer query's result. Therefore, in the auxiliary query ②, we can directly execute the extracted subquery. For both shapes, we create the folded query by replacing the expression with its evaluated result. For example, in Listing 1, ② evaluates to \emptyset , which is why we replace the subquery with \emptyset in ③.

Listing 4. **JOIN** can affect the values of ϕ for a given row.

```

CREATE TABLE t0 (c0 INT);
CREATE TABLE t1 (c0 INT);
INSERT INTO t0 VALUES (0);
INSERT INTO t1 VALUES (1);
① SELECT * FROM t0 LEFT JOIN t1 ON t0.c0 = t1.c0 WHERE t1.c0 IS NULL; -- 0|NULL
② SELECT t1.c0, t1.c0 IS NULL FROM t0 LEFT JOIN t1 ON t0.c0 = t1.c0; --NULL|1
③ SELECT * FROM t0 LEFT JOIN t1 ON t0.c0 = t1.c0 WHERE
CASE WHEN t1.c0 IS NULL THEN 1 END; --0|NULL

```

3.2 Folding Dependent Expressions

We conceptualize a dependent expression as a function $R_\phi = F(\theta)$, $\theta = \{c_1, c_2, \dots, c_i\}$, where the dependent expression can be seen as a mapping from its arguments c_1, c_2, \dots, c_i to its results R_ϕ . Each argument c_i is referenced in ϕ and represents a column from the outer context, and its domain consists of the rows within that column of the table. During constant folding and propagation, we must consider that ϕ might evaluate to a different value for every table row. We first obtain the results of the expression on each row (i.e., step ③ for dependent expression) and then represent them using a mapping (i.e., step ⑤ for dependent expression).

Constant folding. To construct the auxiliary query in step ③ for dependent expressions, we first identify θ , which serves as the keys of the mapping, by collecting all the columns that appear in ϕ and reference the outer context. We include these columns in the fetch clause. Additionally, we include ϕ itself in the fetch clause to obtain its results for each row, which corresponds to R_ϕ . Each row of the auxiliary query's result contains the values of θ , as well as the values of R_ϕ , which are the results of ϕ . Therefore, each row of the results represents a map entry, where the values of θ serve as keys and the values of R_ϕ as corresponding map values. For example, the dependent expression $c_0 + c_1 > 0$ shown in Figure 1 references two columns: c_0 and c_1 . These two columns as well as the expression itself are included in the fetch clause of the auxiliary query. The auxiliary query's results consist of two rows, each representing the result of ϕ for the corresponding row of the table. As another example, Listing 2 shows how we create the auxiliary query for correlated subqueries. To obtain the results of the subquery of query ① on all rows passed from the outer query, we include the subquery in the fetch clause as shown in query ②. While the subquery references two columns $x.classID$ and $y.classID$ in the predicate, only $x.classID$ is from the outer context. Thus, in query ②, θ contains only $x.classID$.

Supporting **JOINS** involves two considerations. First, the auxiliary queries must use the same **JOIN** clauses as the original query, except in cases where ϕ serves as the predicate within the **JOIN** clause. Listing 4 shows an example where we assign ϕ the concrete expression $t1.c0 \text{ IS NULL}$. The original query fetches a single row whose $t1.c0$ column holds a **NULL** value. This is because the predicate of the **LEFT JOIN** evaluates to false. Then, the expression $t1.c0 \text{ IS NULL}$ is evaluated to be true. As a result, the mapping assigns only one row in $t1.c0$ to the result of $t1.c0 \text{ IS NULL}$: from **NULL** to **TRUE**. To construct this mapping in the auxiliary query, the same **JOIN** as the original query must be used, ensuring that the only row in $t1.c0$ has a **NULL** value. Therefore, in the case of dependent expressions, we need to determine in advance whether to use the **JOIN** clause in the original query during the generation of ϕ . This corresponds to line 2 in Algorithm 1. Subsequently, this information will be stored within the set $\{t_i\}$. Conversely, if the folded expression ϕ is used as the predicate of **JOIN** clause, there is no need for a **JOIN** clause in the auxiliary query, as the expression ϕ would be evaluated with the row values before the **JOIN** operation. Second, although we generate non-empty tables, an empty result can still occur, for example, when using an **INNER JOIN** with a false predicate. In such scenarios, we discard the test.

Listing 5. The subquery, when used as the fetch keyword in a **SELECT**, must return only one column and one row.

```
CREATE TABLE t0(c0 INT);
CREATE TABLE t1(c0 INT);
INSERT INTO t0 VALUES (1);
INSERT INTO t0 VALUES (2), (3);
SELECT t0.c0, (SELECT t1.c0 FROM t1 WHERE t1.c0 > t0.c0) FROM t0;
-- Error: Subquery returns more than 1 row
SELECT t0.c0, (SELECT t1.c0, t1.c0 FROM t1 WHERE t1.c0=2) FROM t0;
-- Error: Operand should contain 1 column(s)
```

Listing 6. A bug found in TiDB caused by a predicate generating an incorrect result in **INSERT**.

```
CREATE TABLE t0(c0 BIGINT NOT NULL);
INSERT INTO t0(c0) VALUES (1);
CREATE TABLE ot0 (c0 BIGINT);
INSERT INTO ot0 SELECT t0.c0 AS c0 FROM t0 WHERE VERSION() >= t0.c0;
① SELECT * FROM ot0; -- empty result ✖
② SELECT t0.c0 AS c0 FROM t0 WHERE VERSION() >= t0.c0; -- 1
③ SELECT * FROM (SELECT 1) AS ft0; -- 1 ✓
```

Constant propagation. To represent the mapping obtained by the auxiliary query, we use a **CASE** expression in the folded query, as shown in red in step ⑤ of Figure 1. This mapping shows that ϕ generates 0 for the first row and 1 for the second row of the table. In Listing 2, the query ⑥ also uses the **CASE** expression to represent the mapping for the correlated subquery. The **CASE** pattern resembles a well-known optimization for dynamically-typed languages, known as a *polymorphic inline cache* [15].

3.3 Construction of the Original Query

The expression ϕ generated in step ② can be used in a predicate, which, in turn, can be used in any clause requiring a predicate to construct the original query in step ④. In Algorithm 1, the function **QueryGenerate**, called in line 11, generates the original query.

Predicate construction. We randomly generate predicates that contain or correspond to ϕ . While we use an existing random generation approach implemented in SQLancer [30, 32], subqueries require additional attention, as they were not supported by existing approaches. Subqueries can evaluate to three different result types: (1) a scalar value, which is a single value; (2) a row value, which is an ordered list of two or more scalar values; (3) multiple row values. When using a subquery in the fetch clause of a **SELECT** statement, DBMSs typically allow the subquery only to return a scalar, as illustrated for MySQL in Listing 5. This restriction applies to auxiliary queries for dependent expressions generated in step ③. To this end, for correlated subqueries, we either use an aggregate function without using a **GROUP BY** clause, or use a **LIMIT** clause. Both ensure that a scalar value is returned. For both non-correlated and correlated subqueries, we can use one of the multiple subquery operators, which apply for subqueries returning any of the result types. Such operators include **EXISTS**, **IN**, **ANY**, and **ALL**.

Query construction. The generated predicate can be used in any SQL statement where a predicate is required. Our approach supports placing these predicates not only in the **WHERE**, **JOIN**, **HAVING**, **GROUP BY**, and **ORDER BY** clauses of **SELECT**, but also in other statements that require predicates, such as **CREATE INDEX**, **CREATE VIEW**, **UPDATE**, **INSERT**, and **DELETE**.

Listing 7. A bug found in CockroachDB caused by a false predicate being always evaluated to true.

```

CREATE TABLE t1 (v VARBIT);
INSERT INTO t1 VALUES (B'11');
Ⓒ WITH t2 AS (SELECT NULL AS b) SELECT t1.v FROM t1, t2 WHERE t1.v NOT BETWEEN
    t1.v AND (CASE WHEN NULL THEN t2.b ELSE t1.v END); -- empty result ✓
Ⓐ SELECT NULL AS b; -- NULL
CREATE TABLE t2 (b BIT);
INSERT INTO t2 VALUES (NULL);
Ⓔ SELECT t1.v FROM t1, t2 WHERE t1.v NOT BETWEEN t1.v AND
    (CASE WHEN NULL THEN t2.b ELSE t1.v END); -- '11' ✗

```

Implementation details. Some DBMSs follow strict data type rules when using binary operators, for example, DuckDB and CockroachDB. Before constructing the predicate, it is necessary to query the return type of the expressions. DuckDB provides the `typeof()` function for this purpose, whereas CockroachDB provides `pg_typeof()`. Additionally, DBMSs such as CockroachDB lack automatic implicit casts for converting predicates from any type to boolean. Therefore, we must generate a boolean expression explicitly as the predicate. Some DBMSs, such as SQLite and MySQL, can automatically convert the data type of expressions used as operands of operators. This allows us the freedom to choose operators for the expression ϕ when testing these DBMSs.

The **ALL** and **ANY** operators are not supported in SQLite and DuckDB, and they disallow a value list as an argument in MySQL and TiDB. To overcome this limitation in MySQL and TiDB, we use the **UNION** operator. For example, we can represent a set of values `[1, 2, 3]` as `SELECT 1 UNION SELECT 2 UNION SELECT 3`.

3.4 Language Features Beyond Predicates

Our approach is applicable not only to constant-folding predicate expressions, but also to contexts that expect a relation (e.g., a table or a view reference). A subquery computing a non-empty result can be used as the source of values for tables, allowing us to test other language features such as **INSERT**, common table expressions, and derived tables. In step ④, we construct the original query by referencing an *original relation*, whose values are obtained from a subquery. In step ⑤, we change the reference of the original relation in the original query to a *folded relation* to construct the folded query, whose values are sourced from a table value constructor. A table value constructor in SQL is a constant list, typically represented by **VALUES** (...). Similar to the process of testing predicates, this approach also applies constant folding and constant propagation to subqueries. As before, the folded query should yield identical results to the original query.

CODDTest introduces three approaches to constructing original (i.e., step ④) and folded (i.e., step ⑤) relations. We randomly select one of them to construct the original relation and another one for the folded relation. The first approach creates a table based on the structure of the subquery results and then uses the subquery in an **INSERT** statement to add values to the table. For example, the query ③ shown in Listing 6 references the original relation $\text{ot}\theta$, whose value originates from the subquery in the **INSERT** statement. The second one is a derived table, which is used in the **FROM** clause of **SELECT**. The query ⑥ shown in Listing 6 references the folded relation $\text{ft}\theta$, which is sourced from a constant. We use a subquery as the operand of **AS**, as the left operand of **AS** can only be a subquery. The third one is common table expression (CTE) [8], which was introduced in the SQL Standard of 1999. It is a component of a **SELECT** query that allows creating a temporary table inside the query and retrieving values from subqueries or table value constructors. The query ⑦, shown in Listing 7, demonstrates the application of the CTE in our approach. The bug-inducing

test case shown in Listing 6 uses `INSERT` statement and derived table,⁴ while the bug-inducing test case shown in Listing 7 uses `INSERT` statement and CTE.⁵ We consider this as an extension of our original approach, because the folded table is not necessarily a constant.

4 EVALUATION

We implemented CODDTest, and evaluated three important aspects. First, we evaluated CODDTest's effectiveness by studying how many unique, unknown bugs our approach could find in widely-used DBMSs that have been extensively tested by the state-of-the-art approaches. Second, we sought to determine whether existing approaches would indeed be unable to find the bugs we reported. Third, we compared CODDTest's performance with the state-of-the-art approaches.

Implementation. We implemented CODDTest in SQLancer, a popular tool for DBMS testing, which supports multiple state-of-the-art test oracles for finding logic bugs, including NoREC, PQS, and TLP. SQLancer provides manually-written, rule-based generators specific to the DBMS under test, which can be used to generate statements and expressions. CODDTest initializes the database states using SQLancer's random generation method, which generates random statements to create tables, views, indexes, and insert values into tables. For generating the original query, we adopted the logic used in NoREC, and added support for generating subqueries. For generating auxiliary queries, we implemented an additional generator. For folded queries, we replaced the expression in the original query with the constant by replacing child nodes in the Abstract Syntax Tree (AST) that SQLancer provides. CODDTest is easy to understand and implement, so we could have chosen alternative database and query generators that we could similarly extend to realize our approach.

Tested DBMSs. We selected five DBMSs as our test targets: SQLite,⁶ MySQL,⁷ CockroachDB [36], DuckDB [29], and TiDB [16]. We chose these specific DBMSs for several reasons. Firstly, these DBMSs are widely used, popular, and considered mature. MySQL and SQLite are long-established DBMSs and thus rank highly in rankings such as the DB-Engines ranking.⁸ CockroachDB, DuckDB, and TiDB are relatively recent DBMSs, which are highly popular on GitHub, with star counts of 27.2K, 10.5K, and 34.2K respectively. Secondly, these DBMSs represent various types of DBMSs. SQLite and DuckDB are embedded DBMSs that run in the same process as the application that uses them. MySQL is a traditional, relational, and client-server DBMS. Both CockroachDB and TiDB are distributed relational DBMSs, designed to handle large-scale deployments and ensure high availability and scalability in a distributed environment [28]. Finally, all five DBMSs have been extensively tested by various methods such as NoREC [30], TLP [31], PQS [32], and DQE [35], suggesting that any newly found bugs might have been overlooked by these existing approaches.

DBMSs versions. We tested the latest development versions of the aforementioned DBMSs. We downloaded SQLite from its source code repository, and downloaded the other DBMSs from GitHub. For SQLite, we tested version c1f2a1d5 and later trunk versions. For MySQL, we tested commit version ea7087d. For CockroachDB, we tested commit version 07c7d4b and later versions. For DuckDB, we tested commit version b8cf6a9 and later versions. For TiDB, we tested version c233969 and later commit versions.

Baselines. We selected NoREC [30], TLP [31], DQE [35], and EET [17] as baselines for evaluation. NoREC, TLP, PQS [32], TQS [37], DQE, and EET are the state-of-the-art approaches for automatically

⁴<https://github.com/pingcap/tidb/issues/43373>

⁵<https://github.com/cockroachdb/cockroach/issues/102110>

⁶<https://sqlite.org/index.html>

⁷<https://www.mysql.com/>

⁸<https://db-engines.com/en/ranking>

Table 1. CODDTest found 45 unique bugs in five mature DBMSs.

DBMS	Bug type				Bug status	
	Logic bug	Internal error	Crash	Hang	Fixed	Verified
SQLite	6	1	0	0	7	0
MySQL	1	1	0	0	0	2
CockroachDB	7	4	0	2	11	2
DuckDB	5	2	2	3	12	0
TiDB	5	6	0	0	3	8
Total	24	14	2	5	33	12

testing DBMSs. NoREC and DQE assume that the same predicate always accesses the same row of a table, regardless of which clause it is used in. TLP decomposes a query into three partitioning queries, which retrieve rows based on predicates p , $\text{NOT } p$, and $p \text{ IS NULL}$, respectively. EET is a concurrent work with CODDTest, which has proposed the first expression-level manipulation approach, known as Equivalent Expression Transformation (EET), which introduces tautologies and contradictions to construct equivalent queries. PQS generates a `SELECT` query to retrieve a pivot row, and checks whether the DBMS fetches it as expected. PQS requires a high implementation effort, as operations of the DBMS must be implemented also in the testing tool. For this reason, it is no longer actively maintained in the SQLancer project, and currently triggers false alarms.⁹ TQS was proposed to find bugs in join optimizations; however, it is not publicly available. Thus, we omitted both PQS and TQS from the comparison.

4.1 Effectiveness

Methodology. We evaluated the effectiveness of CODDTest by testing the five aforementioned DBMSs. We intermittently ran CODDTest for a period of four months, during which we also implemented it. This is a standard methodology used to evaluate the effectiveness of automated testing tools [23, 26, 30]. Before reporting bugs, we manually reduced the bug-inducing test cases [39]. To avoid reporting duplicate issues, for multiple reports that we suspected to affect the same component, we reported the subsequent bug only after the previously reported one was fixed. For some DBMSs (e.g., TiDB), we refrained from reporting more bugs, due to the large number of unfixed bugs.

Results. Table 1 shows the number of bugs found by CODDTest, as well as the status of the bugs. CODDTest found a total of 45 previously unknown bugs. Out of these, 24 were logic bugs, 14 were internal errors, 2 were crashes, and 5 were hang-related issues. Out of all 45 bugs, 33 were fixed, and 12 were verified. Out of the 24 logic bugs, 19 were fixed, and 5 were verified. These results are highly encouraging, considering that these DBMSs have been the focus of many testing works as mentioned above.

Logic bugs causes. We analyzed the queries that triggered the bugs and identified that 12 bugs were triggered by folded queries—the queries generated and executed in step ⑤—and 12 by original queries. Out of the bug-inducing folded queries, 11 queries used folded constants that were derived from non-correlated subqueries. Of the bugs caused by them, 6 were triggered by replacing a constant, and the remaining 5 bugs were triggered by a value list, all of which were related to the `IN` operator. Only one query used a constant that was derived from the query with a simple

⁹<https://github.com/sqlancer/sqlancer/issues/527>

expression. Out of the remaining 12 bugs triggered by original queries, all of them used queries containing non-correlated subqueries. Feedback from the developers indicated that 3 bugs were bugs in subquery processing, while 5 bugs were unrelated to subqueries. For the 4 remaining bugs, we are unclear about the root cause due to limited developer feedback. Overall, we detected most of the bugs through non-correlated subqueries or folded queries derived from non-correlated subqueries. This is primarily because non-correlated subqueries were our initial test focus, as they can be executed independently and are straightforward to implement. Additionally, half of all the bugs were triggered by the folded queries. We believe that applying CODDTest to basic expressions or constant expressions could also have been used to detect some of these bugs. For example, for the bug-inducing test case shown in Listing 8, using a constant expression $1 - 1$ as ϕ could have detected this bug.

Other bugs. CODDTest found a total of 21 other bugs in five DBMSs, including crashes, internal errors, and hangs. Out of the 21 bugs, 14 have been fixed, and 7 have been verified. We found two crashes in DuckDB, both of which caused our tool to terminate unexpectedly, and the test cases resulted in segmentation faults when using DuckDB’s command line interface. These two crashes were introduced in the IEJoin optimization [20]. One crash was caused by an index out-of-bounds error, while the other was due to a type mismatch. These two crashes have been fixed. We found 14 internal errors in these five DBMSs, 6 have been fixed, and the others were verified. All 5 hang-related bugs found in DuckDB and CockroachDB have been fixed. We believe these crashes, internal errors, and hangs could be detected by other automated testing approaches [13, 40].

Bug importance. Although anecdotal, developer feedback is an important indicator of an approach’s effectiveness and the found bugs’ importance. Two DBMS companies reached out to us about our testing efforts; both were interested in how we found the bugs, and one invited us to present the approach to the development team. Furthermore, we received positive feedback on the public bug trackers. For example, a developer of CockroachDB commented in one of our reports “we really appreciate your work!”¹⁰ Three bugs in CockroachDB were assigned the “S-0” or “S-2” label, indicating high-impact bugs that were difficult to resolve. Three bug reports in TiDB were labeled as “Major”, which represents the highest bug severity.

False alarms. While realizing our approach, we identified corner cases that resulted in false alarms in our initial implementation. First, applying constant folding with floating-point numbers can result in false alarms. We avoid these in practice by eschewing test cases with small or large float-point values. Second, SQLite’s relaxed type system allows, in some context, values of different types to be returned. We avoid this through explicit casts to the data type we expect. We have not observed any false alarms after addressing these issues.

Discussion. Although we intermittently ran CODDTest over a period of four months, during which we also implemented it, most of the bugs were found at the beginning of CODDTest’s execution. As shown in the evaluation in Section 4.2, CODDTest identified 25 unique bugs in an older version of SQLite within 24 hours, with 13 of those bugs found in the first hour. Since CODDTest is a black-box method, it is an ongoing research challenge to determine whether a system has been sufficiently tested [25].

¹⁰<https://github.com/cockroachdb/cockroach/issues/104319>

Table 2. The number of detectable bugs by test oracles.

Oracles	NoREC	TLP	DQE	Only by CODDTest
Num	11	12	4	11

4.2 Test Oracle Comparison

A testing approach is valuable if it finds bugs that were overlooked by existing approaches. Thus, we sought to confirm whether the state-of-the-art test oracles are indeed ineffective in finding the bugs found by CODDTest, as well as to identify the types of bugs that CODDTest can not detect.

Methodology. To address these two questions, we designed three experiments. First, we ran CODDTest on the release version 3.30.0 of SQLite, in which NoREC and TLP detected 31 bugs, to determine how many of these bugs CODDTest could identify over 24 hours with 10 threads. The developers of SQLite were proactive in fixing bugs, which allowed us to determine whether two bug reports trigger the same bug by applying the bug-fixed commit to SQLite and checking if the bug could still be triggered. We selected NoREC, TLP, and EET for this comparison, because DQE is conceptually similar to NoREC. Additionally, the test cases generated by DQE are incomplete and cannot be automatically analyzed. Second, we analyzed the earliest versions of DBMSs where logic bugs identified by CODDTest can be triggered, to determine whether these bugs were introduced before the state-of-the-art approaches were published, suggesting that the existing approaches missed the bugs that CODDTest subsequently found. This methodology was also used in the EET work [17]. Third, we implemented a best-effort comparison by manually inspecting and analyzing the bug-inducing test cases and bugs found by CODDTest and analyzing whether the state-of-the-art test oracles could have found them. We selected NoREC, TLP, and DQE for this comparison, because EET’s transformations explore an extensive search space, making it difficult to conduct the transformations and check their results manually. For every analyzed bug, we include an analysis in the supplementary materials including the test cases constructed by the test oracles, allowing scrutinization of the results.

Results of oracles on an older version of SQLite. During a 24-hour period, NoREC, TLP, EET and CODDTest generated 97,003, 103,103, 16, and 6,990 bug reports, respectively, with 27, 27, 6, and 25 of the bugs being unique. Additionally, NoREC, TLP, EET, and CODDTest reported 3, 2, 3, and 4 bugs that were found by that oracle alone. Although the significant overlap might be surprising, we believe it is reasonable—a testing approach is useful if it can find new bugs overlooked by other approaches. Note that we conducted this experiment on a single and stable version of SQLite and that in other, especially older versions of SQLite, likely more unique bugs could be detected.

As all of these three oracles are black-box testing methods, it is challenging to determine which types of bugs will be overlooked in general. However, we noticed that CODDTest overlooked 14 bugs, 4 of which were related to indexing. Therefore, we speculate that CODDTest cannot effectively test indexing functionality.

Although it is an established methodology to run testing approaches for 24 hours [21], we noticed that it might be insufficient for some bugs to be consistently found. By running the 24-hour experiments, we found that 9, 13, and 16 bugs found by NoREC, TLP, and CODDTest, respectively, had fewer than 10 reports. Therefore, generating certain corner cases might require a substantial amount of time. This motivated us to conduct the manual comparison study, whose results are detailed next.

Listing 8. A bug found in SQLite related to **JOIN**.

```

CREATE TABLE vt0(c2);
CREATE TABLE t1(c0 TEXT);
INSERT INTO t1(c0) VALUES (1);
INSERT INTO vt0(c2) VALUES (-1);
CREATE VIEW v0(c0) AS SELECT 0 FROM t1;
① SELECT vt0.c2 AS c1 FROM t1 CROSS JOIN v0 ON (
    EXISTS (SELECT v0.c0 FROM v0 WHERE false)) FULL OUTER JOIN vt0 ON 1; -- -1 ✓
Ⓐ SELECT v0.c0 FROM v0 WHERE false; -- empty result
Ⓔ SELECT vt0.c2 AS c1 FROM t1 CROSS JOIN v0 ON (0) FULL OUTER JOIN vt0 ON 1;
    -- empty result ✗

```

Results on bugs introduction times. CODDTest can detect long-latent bugs in DBMSs that have been overlooked by state-of-the-art approaches. NoREC supports SQLite and MySQL. CODDTest found 7 logic bugs in these two systems, and 1 bug was introduced before NoREC was published. All DBMSs supported by CODDTest are also supported by TLP; 6 out of the 24 bugs found by CODDTest were introduced before TLP was published. DQE supports SQLite, MySQL, CockroachDB, and TiDB, and 15 out of 19 bugs found in these DBMSs were introduced before DQE was published. EET supports MySQL, SQLite, and TiDB, and the 12 bugs found in these DBMSs were all introduced prior to EET’s first bug report, indicating that EET may have missed these bugs. Out of the 24 logic bugs found by CODDTest, 6 were introduced before 2020, and 20 were introduced before 2023. The bug with the longest latency—14 years—was found in MySQL. This result demonstrates that CODDTest can effectively identify long-latent bugs.

Results of manual comparison. As shown in the manual-analysis results in Table 2, out of 24 bugs detected by CODDTest, 11 can be detected by NoREC, 12 can be detected by TLP, 4 can be detected by DQE, and 11 can only be detected by CODDTest. We found that these 11 bugs were in language features that were not supported by existing test oracles, including subqueries, **ON** clauses, **CASE** and **ANY** expressions, the **AVG** function, and **INSERT** statements.

Subqueries. Three bugs related to subqueries, which are out-of-scope for the three test oracles. Listing 1 shows such a bug-inducing test case for SQLite. The SQLite developers communicated that reproducing it requires five conditions: (1) the query must contain an aggregate subquery; (2) the aggregate subquery must have a **GROUP BY** clause; (3) the **GROUP BY** clause must reference terms that are not included in the result set of the query; (4) the query planner must choose to implement the **GROUP BY** clause by doing a sort operation; (5) the outer query that contains the aggregate subquery must make use of indexed expressions. Under these conditions, SQLite assigned an incorrect value to a variable in the SQL AST associated with the subquery. The transformations performed by existing test oracles are incapable of generating queries that violate these conditions, preventing them from detecting this bug. We found the second subquery-related bug in DuckDB, which was caused by incorrect processing of the return type of a subquery. CODDTest executed the auxiliary query to obtain the subquery’s results with the correct data type, which were then used in the folded query. The folded query subsequently generated the correct results. However, due to the incorrect handling of the subquery’s return type, the original query produced incorrect results. The third one was found in TiDB and caused by the incorrect recognition of columns with identical names, resulting in the misinterpretation of a non-correlated subquery as a correlated one.

ON clause. The **ON** clause of **JOIN** remained untested by NoREC, TLP, and DQE. Since these three test oracles rely on the direct mapping relationship between predicates and corresponding rows, which cannot be directly determined for the **JOIN ON** clause, they are unable to effectively test the clause. Two bugs in SQLite found by CODDTest were missed by these three test oracles for

Listing 9. A bug found in CockroachDB related to **IN** operator.

```
CREATE TABLE t (c INT4);
INSERT INTO t (c) VALUES (0);
❷ SELECT c FROM t WHERE c IN (0, 862827606027206657::INT8); -- empty result ❌
```

Listing 10. A bug found in TiDB related to **IN** operator.

```
CREATE TABLE t0(c0 BOOL);
INSERT INTO t0 VALUES (true);
❷ SELECT t0.c0 FROM t0 WHERE ((CASE ((CASE t0.c0 WHEN 6 THEN 0.03 ELSE t0.c0 END )
    LIKE (t0.c0)) WHEN t0.c0 THEN 1 END) AND (t0.c0 IN (1))); -- empty result ❌
```

this reason. Listing 8 shows one of the bugs that was triggered when replacing a predicate in the **JOIN ON** clause with a constant.¹¹ The original query that embedded the subquery generated the correct results and provided an opportunity to find this bug.

CASE and ANY expressions. Listing 7 shows a bug found in CockroachDB, where a bug in processing the **CASE** caused the expression to be incorrectly evaluated to **TRUE**, regardless of in which clause the predicate was placed. Consequently, NoREC, TLP, and DQE were unable to detect this bug. We observed two bugs related to **CASE** expressions and one bug related to **ANY** expressions, which consistently evaluated to an incorrect result.

AVG function. A bug in CockroachDB related to **AVG** function, which produced inconsistent results when the argument order was altered. The original and folded queries produced results with different orders, enabling us to find this bug.

Language features beyond queries. Listing 6 shows a bug found in TiDB, specifically in an **INSERT** statement, whose argument was a non-correlated subquery. The inner query may return a non-empty result, but it fails to produce a non-empty result within the **INSERT** statement. None of the existing test oracles support testing **INSERT** statements, as **INSERT** does not directly use predicates. The extension of our method described in Section 3.4, can find bugs in **INSERT** statements even when the bug cannot be reproduced using other statements.

Variations in expression behavior across clauses. NoREC and DQE were designed to identify logic bugs by assuming that a specific predicate evaluates to the same value irrespective of in which clause it is used. However, we discovered that either DQE or NoREC failed to detect some of these bugs. There are two bugs in this category. Listing 9 illustrates a bug in CockroachDB associated with the **IN** operator.¹² CockroachDB produced the correct results when the right operand was a subquery, but yielded incorrect results when the right operand was a value list. This predicate produced incorrect results in **SELECT** queries, while correctly functioning in **UPDATE** and **DELETE** statements. As a result, NoREC failed to detect this bug, whereas DQE successfully found it. Listing 10 illustrates a bug in TiDB that is also associated with the **IN** operator.¹³ Similar to the bug mentioned above, it generated correct results when the right operand was a subquery, but yielded incorrect results when the right operand was a value list. The difference, however, lies in the behavior exhibited in **WHERE** clauses, where it consistently produced incorrect results, while correctly functioning when used in the fetch clause of a **SELECT**. As a result, DQE failed to detect this bug, while NoREC successfully found it.

¹¹<https://sqlite.org/forum/forumpost/96cd4a7e9e>

¹²<https://github.com/cockroachdb/cockroach/issues/102864>

¹³<https://github.com/pingcap/tidb/issues/43624>

Listing 11. A bug found in DuckDB, which triggers an error when applying NoREC to it. The expressions highlighted in blue illustrate the application of the NoREC oracle.

```
CREATE TABLE t0(c1 INT8);
INSERT INTO t0(c1) VALUES ((1));
❶ SELECT t0.c1 FROM t0 WHERE (((-1314689763) + (-1947665992)) <=
    (EXISTS( SELECT t0.c1 FROM t0 WHERE false))); -- empty result ❗
❷ SELECT t0.c1 FROM t0 WHERE false; -- empty result
❸ SELECT t0.c1 FROM t0 WHERE (((-1314689763) + (-1947665992)) <= (false)); -- 1 ✓
SELECT (((-1314689763)+(-1947665992)) <= (EXISTS( SELECT t0.c1 FROM t0 WHERE false)))
FROM t0;
-- Out of Range Error: Overflow in addition of
-- INT32 (-1314689763 + -1947665992)!
```

Furthermore, four bug-inducing test cases triggered errors when other oracles were applied to them, indicating that these logic bugs could not be found by them. Two of them triggered internal errors when we applied NoREC to them, and two of them triggered semantic errors when we applied DQE to them. We encountered scenarios where a predicate triggered a logic bug within **WHERE** clauses, while simultaneously causing an internal error when used in a **SELECT** statement. Listing 11 shows one of the two such bugs found by CODDTest.¹⁴ When we used the predicate triggering the logic bug in the **SELECT**'s fetch clause, the query triggered an internal error. As a result, although NoREC cannot directly detect this logic bug, the internal error could have still been found. However, when applying DQE to the logic bug we discovered in TiDB, we observed that the predicate, which executed normally in **SELECT**, triggered a semantic error when used in **UPDATE** and **DELETE** statements. Although the authors of DQE referenced a similar case found in MySQL in their paper, they did not provide an explanation for the underlying reason. To investigate further, we searched the MySQL bug list and uncovered the root cause.¹⁵ In MySQL, the **SELECT** statement allows comparing values with different data types. However, it does not permit this in **UPDATE** and **DELETE** statements. Therefore, DQE triggers a semantic error in this particular test case, and is unable to detect this logic bug.

Summary. CODDTest can find bugs that are missed by the state-of-the-art approaches, such as in subqueries, **JOIN**, operators, and functions. We believe that the existing state-of-the-art test oracles cannot be extended to find these missed bugs in any obvious way. While query generators could still generate the same queries, the existing test oracles would miss the logic bugs associated with these language features. Additionally, CODDTest can test language features beyond predicates, as discussed in Section 3.4. This includes common table expressions (CTE) and derived tables.

4.3 Efficiency Comparison

In this section, we compare the performance of CODDTest to the state-of-the-art approaches.

Methodology. NoREC, TLP, and DQE are all implemented based on SQLancer, which enables a fair comparison of the approach, as we could use the same settings (e.g., the same number of threads). Besides CODDTest, we also evaluated two other configurations to explore the impact of testing different expression types on performance. Specifically, we explored the performance of CODDTest when using only expressions with no subqueries (i.e., CODDTest & Expression) and only subqueries (i.e., CODDTest & Subquery).

Metrics. We define six metrics for this evaluation. We measured test throughput by counting the number of successful test cases executed. We recorded the number of queries that were successfully

¹⁴<https://github.com/duckdb/duckdb/issues/7094>

¹⁵<https://bugs.mysql.com/bug.php?id=111483>

Table 3. The number of tests conducted by each approach.

Oracle	# of tests	# of successful queries	# of unsuccessful queries	QPT	# of unique query plans	branch coverage
NoREC	2,086,646k	4,207,286k	149,036k	2.05	172,808	63.18%
TLP	976,216k	2,180,736k	398,919k	2.23	137,743	63.63%
DQE	441,350k	7,502,402k	21,997k	17.00	486	46.71%
CODDTest	497,092k	1,655,518k	53,102k	3.33	2,577,603	63.06%
CODDTest & Expression	1,423,068k	4,411,510k	326,849k	3.10	7,399	63.23%
CODDTest & Subquery	423,310k	1,488,817k	47,141k	3.51	2,755,619	62.19%

executed (*i.e.*, successful queries) and those that encountered expected errors (*i.e.*, unsuccessful queries). Expected errors either refer to queries triggering unfixed internal errors in the DBMS, or cases where the query is semantically incorrect (*e.g.*, errors like unexpected integer overflows are difficult to avoid during construction). Oracles differ in how many queries they execute for each test case. Note that this number is not static as, for example, generating a test might fail due to an unsuccessful query. Thus, we compute the average number of queries for each successfully-executed test, that is, the queries per test (QPT). We also evaluated the number of unique query plans generated by each method. We collected only the query plan of the most complex query, which was the optimized query in NoREC, the partitioning query in TLP, the `SELECT` query in DQE, and the original query in CODDTest. Lastly, we compare the branch coverage of each oracle at the end of execution, as this provides a more rigorous assessment than statement coverage.¹⁶

Experimental setup. We conducted the performance evaluation on a server with a 64-Core AMD EPYC 7763 Processor at 2.45GHz and 512GB of memory running Ubuntu 22.04. We selected SQLite as our test target, and conducted this experiment with release version 3.42.0.0. We executed each approach with 10 threads for a duration of 24 hours and recorded the number of tests conducted. SQLancer provides a *MaxDepth* option, which controls the maximum depth of an expression. We used the default configuration of SQLancer, where this option is set to 3. We conducted separate experiments with the same configuration and times to collect statistics on unique query plans and branch coverage, as querying these metrics requires additional time.

Results. Table 3 shows the results. On average, CODDTest has a lower test throughput compared to NoREC and TLP, but a higher test throughput than DQE. Specifically, CODDTest was approximately $4.20\times$ slower than NoREC, $1.96\times$ slower than TLP, and $1.13\times$ faster than DQE. Upon further investigation, we have found two reasons why CODDTest is slower than NoREC and TLP. Firstly, for each test, CODDTest required executing at least three queries to derive the test oracle, including original query, folded query, and auxiliary query. When using subqueries to create relations, an additional query is required to retrieve the types of the subquery's result, which is used for table creation. This explains why the QPT exceeded 3 for CODDTest and its configurations. Additionally, for CODDTest, when applying a subquery in an `INSERT` statement, additional statements are needed to create and drop tables to maintain the database state. In contrast, NoREC executed only about two queries per test. For TLP, the average number of queries per test was 2.23; the number exceeds 2 as the test oracle randomly either executes the partitioning queries as one query—using `UNION ALL`—or executes three queries. For DQE, a test requires not only the three statements `SELECT`, `UPDATE`, and

¹⁶<https://www.sqlite.org/testing.html>

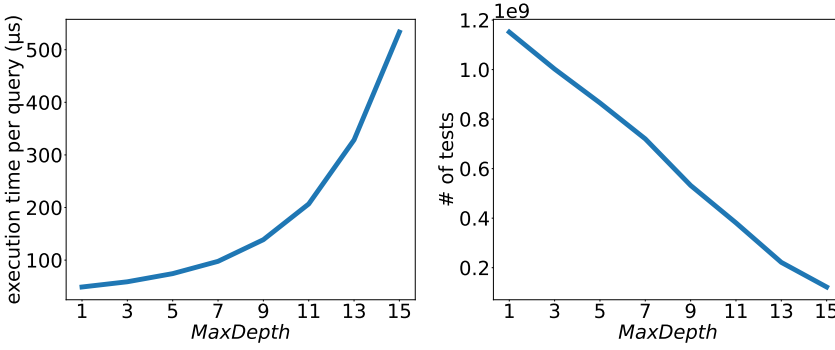


Fig. 2. The impact of expression complexity on query execution time and test throughput.

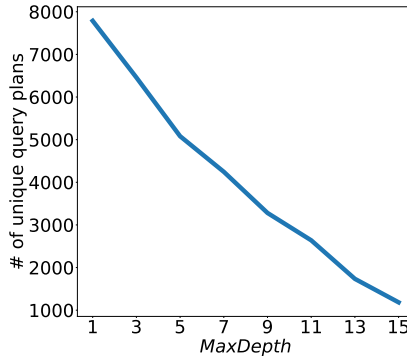


Fig. 3. The impact of expression complexity on unique query plans.

DELETE, but also additional statements for two extra columns. These extra columns are used to uniquely identify each row and to track whether a row has been modified. Secondly, the execution of queries with subqueries is significantly slower compared to queries with expressions (*i.e.*, 7.18 \times). We found that, on average, queries with expressions alone required only 44.73 microseconds for execution, whereas queries with subqueries required 321.19 microseconds. Although subqueries require more execution time, they cover a significantly greater number of unique query plans, resulting in CODDTest covering 14.92 \times (compared to NoREC) to 5303.71 \times (compared to DQE) more unique query plans than other oracles. NoREC, TLP, and CODDTest have a similar result of branch coverage, as they use the same statement generator. DQE, however, has lower branch coverage, because it cannot test certain language features, such as **JOIN**. However, branch coverage does not effectively illustrate the oracle's effectiveness, as SQLite already achieves 100% branch coverage with its own test suites, yet still contains logic bugs. Recent research suggests that exercising more unique query plans corresponds to uncovering more interesting and potentially erroneous behaviors in the DBMS under test [2].

Expression complexity. We also evaluated the impact of expression complexity on the efficiency of CODDTest. To isolate the effect of expression complexity, we conducted our evaluation on CODDTest & Expression, excluding the influence of subqueries. We define the complexity of an expression using *MaxDepth*. This experiment focuses on examining how complexity influences

the average execution time of each query (*i.e.*, execution time per query), as well as the overall throughput of CODDTest (# of tests). We conducted this evaluation using the same setup as before (*e.g.*, executing the experiment with 10 threads for 24 hours). The results, as depicted in Figure 2, reveal how varying of *MaxDepth* influence these performance metrics. We can observe that, as *MaxDepth* increases from 1 to 15, the average running time for each query increases by 9.91 \times , and the throughput of our method decreases by 89.4%. Therefore, the complexity of expressions significantly impacts the efficiency of our approach, as the DBMSs require more time to execute each query. However, in our evaluation, we found that most of the bugs we found by the bug-inducing test cases use only shallow expressions after reduction. Therefore, during testing, there is no need to generate expressions with great depth.

We also investigated the impact of expression complexity on the diversity of unique query plans (# of unique query plans) generated when varying *MaxDepth*. As shown in Figure 3, we find that the complexity of the expression has a significant impact on the number of unique query plans, which shows a decreasing trend similar to that of throughput. We can conclude that increasing expression depth with language features other than subqueries does not significantly exercise additional logic in DBMSs; the use of subqueries in expressions not only increases expression complexity, but also effectively triggers more logic within DBMSs.

Summary. CODDTest generates more unique query plans than existing oracles, suggesting that it exercises interesting functionality in the DBMSs. However, this results in slightly lower throughput. Prior research [5] suggests that the resources required to find bugs increase exponentially—a constant factor will not significantly decrease the bug-finding effectiveness.

5 Discussion

CODDTest scope. One potential concern is how CODDTest can be applied to features not considered in this work. First, we believe that any kind of expression can be supported by CODDTest assuming that it is deterministically computing its result. One advanced feature we did not consider are window functions, which can be supported by being used in subqueries. In addition, our approach is not specifically designed to test functionality unrelated to expressions, such as clauses (*e.g.*, `LIMIT`). While we can apply constant folding and propagation to the predicates of clauses, it is unclear how to replace the clauses themselves. Furthermore, like other logic bug detection methods, our approach lacks support for expressions with non-deterministic functions or ambiguous queries.

Our approach primarily applies to the statement types Data Query Language (DQL) and Data Manipulation Language (DML). It is inapplicable to most statements of Data Definition Language (DDL), Transaction Control Language (TCL), Miscellaneous Language (ML), and Data Control Language (DCL). DDL expressions produce results based on data that will be inserted into tables in the future, making CODDTest inapplicable to constant folding and propagation based on the current database state. DCL, TCL, and ML typically do not support the use of expressions as conditions.

The support of closed-source commercial DBMSs. In our evaluation, we considered only DBMSs whose source code is publicly available. One concern could thus be whether our approach could potentially also find bugs in closed-source commercial DBMSs. First, multiple of these open-source DBMSs are developed by companies, such as Cockroach Labs for CockroachDB and PingCAP for TiDB. Second, we attempted engaging with multiple closed-source commercial DBMS vendors in the past, in whose systems we also found bugs. However, we did not receive a response to our bug reports and any potential fixes would be reflected only in subsequent versions, making it difficult to identify their root causes and avoid reporting duplicate bugs.

6 Related Work

Detecting logic bugs in DBMSs. Several approaches have been proposed to detect logic bugs in DBMSs. RAGS [34] applies differential testing to find logic bugs by comparing the results of a query on different DBMSs or different versions of the same system. A key challenge that limits its applicability is that SQL dialects differ widely across DBMSs. Non-optimizing Reference Engine (NoREC) [30] and Differential Query Execution (DQE) [35] place the same predicate in different clauses, with the expectation that this predicate will retrieve the same rows. NoREC primarily focuses on the **WHERE** clause of **SELECT** statements, while DQE expands its scope to the **WHERE** clauses of **SELECT**, **UPDATE**, and **DELETE** statements. Ternary Logic Partitioning (TLP) [31] decomposes a query into three partitioning queries, each of which retrieves rows based on the predicates p , **NOT** p , and p **IS NULL**, respectively. TLP can test the predicate in **WHERE**, **GROUP BY**, **HAVING** clauses, as well as aggregate functions, and **DISTINCT** queries. NoREC and DQE rely on predicates by assuming that their results remain consistent regardless of the clause in which they are placed, and TLP leverages that for any given row, exactly one of p , **NOT** p , and p **IS NULL** evaluates to true. These three approaches are unable to detect logic bugs that affect the evaluation of expressions irrespective of how they are used, making them overlook, for example, the bug shown in Listing 1. Pivoted Query Synthesis (PQS) [32] generates queries that are guaranteed to retrieve a selected row, based on a naive implementation of operators and functions to be tested, Transformed Query Synthesis (TQS) [37] generates queries by decomposing a table into multiple sub-tables, to derive a test case and ground truth for queries that join these tables. The synthesis methods used by PQS and TQS ensure that the generated query matches a specific value or relationship. This is achieved using their self-built evaluation engine. Therefore, their synthesis techniques have limited support for various language features (*i.e.*, TQS only supports equi-join as predicates), as it is challenging to compute the expected results across all of them. This is also why we chose not to compare with PQS and TQS (the latter which is also not publicly available); many language features we test would be difficult to support by these approaches. A concurrent work called Equivalent Expression Transformation (EET) [17] closely relates to our approach. Both CODDTest and EET operate on expressions and aim to derive queries whose results are the same as the original queries, but the conceptual angle on how they derive them is different. EET introduces tautologies and contradictions while ensuring that the result remains equivalent to the original query. In contrast, CODDTest applies constant folding and propagation on expressions to replace certain language features with constants, thereby creating an equivalent, but simpler query. Both are black-box approaches, making it difficult to conceptualize which bugs are overlooked by one of the two approaches, but not the other. However, assuming a perfect query optimizer that can simplify any expressions, EET would be ineffective as the unnecessarily complex expressions could be simplified. For CODDTest, this is not the case, as the query optimizer cannot assume that the database's contents remains the same across queries.

Random and targeted queries. Multiple works have improved query generation for DBMSs, which is complementary to our contribution. Query Plan Guidance (QPG) [2] mutates the database state to cause the DBMS to guide test case generation towards potentially unseen query plans for subsequent queries. SQLRight [24] mutates the SQL statements based on code coverage feedback to cover more code in DBMSs. Squirrel [40] mutates SQL queries based on an intermediate representation to ensure syntax validity and uses coverage feedback for guidance. Griffin [13] proposes a grammar-free mutation approach for testing DBMSs, using a metadata graph to ensure semantic correctness.

Random and targeted databases. Many approaches have been proposed to automate the generation of databases. Gray et al. proposed approaches for generating billions of database records using multiple techniques [14]. Data Generation Language (DGL) [7] is a domain-specific language

designed for generating data that exhibits complex intra- and inter-table correlations. QAGen [4] is a query-aware database generator designed to produce database states by generating them based on a specified parametric query and set of user-defined constraints, ensuring that the outcomes of the query meet the user requirements. ADUSA [19] leverages a constraint solver to generate database data and the corresponding expected results for a given query and database schema.

Testing other aspects of DBMSs. Besides logic bugs, testing approaches were proposed to find performance issues. Cardinality Estimation Restriction Testing (CERT) [3] is designed to detect performance issues through the lens of cardinality estimation. For a given query, CERT derives a more restrictive query, for which the cardinality estimator is expected to predict that the query fetches fewer rows than the original query. APOLLO [18] detects performance regression issues in DBMSs by analyzing multiple versions of a given DBMS, and employs a suite of validation checks to minimize false positives.

7 Conclusion

In this paper, we have presented a black-box approach for detecting logic bugs in DBMSs, named *Constant-Optimization-Driven Database Testing (CDDTest)*. Our key insight is that, for a given query and fixed database state, constant propagation and folding can be applied to specific expressions of the query, assuming that the result remains unchanged. We believe that this idea is non-obvious, as constant propagation and constant folding were originally proposed as compiler optimizations, not for testing DBMSs. We have evaluated CDDTest on five mature and well-tested DBMSs, and found a total of 45 bugs. Of these, 24 were logic bugs, and the remaining were internal errors, crashes, and hang-related bugs. As indicated by our manual analysis, 11 logic bugs were missed by the state-of-the-art approaches. CDDTest generates test cases that exercise more unique query plans, suggesting that it can explore interesting functionality in the DBMS under test. Overall, we believe that CDDTest is a practical, widely applicable DBMS testing approach that complements existing test oracles for logic bugs.

Acknowledgments

This work was partially supported by the National Natural Science Foundation of China under Grant No. 62032010, 62202220. We would like to thank for the financial support from the program of China Scholarships Council (No.202106190065).

References

- [1] Andreas Seltenreich, Bo Tang, Sjoerd Mullender. 2018. SQLsmith: A random SQL query generator. <https://github.com/anse1/sqlsmith>.
- [2] Jinsheng Ba and Manuel Rigger. 2023. Testing Database Engines via Query Plan Guidance. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2060–2071. doi:10.1109/ICSE48619.2023.00174
- [3] Jinsheng Ba and Manuel Rigger. 2024. CERT: Finding Performance Issues in Database Systems Through the Lens of Cardinality Estimation. In *The 46th International Conference on Software Engineering (ICSE'24)*.
- [4] Carsten Binnig, Donald Kossmann, Eric Lo, and M. Tamer Özsu. 2007. QAGen: generating query-aware test databases. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, Beijing, China, June 12-14, 2007*, Chee Yong Chan, Beng Chin Ooi, and Aoying Zhou (Eds.). ACM, 341–352. doi:10.1145/1247480.1247520
- [5] Marcel Böhme and Brandon Falk. 2020. Fuzzing: on the exponential cost of vulnerability discovery. In *ESEC/FSE '20: 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Virtual Event, USA, November 8-13, 2020*, Prem Devanbu, Myra B. Cohen, and Thomas Zimmermann (Eds.). ACM, 713–724. doi:10.1145/3368089.3409729
- [6] Matthias Brantner, Norman May, and Guido Moerkotte. 2007. Unnesting Scalar SQL Queries in the Presence of Disjunction. In *Proceedings of the 23rd International Conference on Data Engineering, ICDE 2007, The Marmara Hotel,*

- Istanbul, Turkey, April 15-20, 2007*, Rada Chirkova, Asuman Dogac, M. Tamer Özsu, and Timos K. Sellis (Eds.). IEEE Computer Society, 46–55. doi:10.1109/ICDE.2007.367850
- [7] Nicolas Bruno and Surajit Chaudhuri. 2005. Flexible Database Generators. In *Proceedings of the 31st International Conference on Very Large Data Bases, Trondheim, Norway, August 30 - September 2, 2005*, Klemens Böhm, Christian S. Jensen, Laura M. Haas, Martin L. Kersten, Per-Åke Larson, and Beng Chin Ooi (Eds.). ACM, 1097–1107. <http://www.vldb.org/archives/website/2005/program/paper/wed/p1097-bruno.pdf>
 - [8] Miguel Cebollero, Jay Natarajan, Michael Coles, Miguel Cebollero, Jay Natarajan, and Michael Coles. 2015. Common Table Expressions and Windowing Functions. *Pro T-SQL Programmer's Guide* (2015), 233–268.
 - [9] Pedro Celis and Hansjörg Zeller. 1997. Subquery Elimination: A Complete Unnesting Algorithm for an Extended Relational Algebra. In *Proceedings of the Thirteenth International Conference on Data Engineering, April 7-11, 1997, Birmingham, UK*, W. A. Gray and Per-Åke Larson (Eds.). IEEE Computer Society, 321. doi:10.1109/ICDE.1997.583381
 - [10] Tsong Yueh Chen, S. C. Cheung, and Siu-Ming Yiu. 1998. Metamorphic Testing: A New Approach for Generating Next Test Cases. *Technical Report* Technical Report HKUST-CS98-01, Department of Computer Science, Hong Kong (1998).
 - [11] Yih-Fam Chen, Emden R Gansner, and Eleftherios Koutsosofos. 1998. A C++ data model supporting reachability analysis and dead code detection. *IEEE Transactions on Software Engineering* 24, 9 (1998), 682–694.
 - [12] E. F. Codd. 1970. A Relational Model of Data for Large Shared Data Banks. *Commun. ACM* 13, 6 (1970), 377–387. doi:10.1145/362384.362685
 - [13] Jingzhou Fu, Jie Liang, Zhiyong Wu, Mingzhe Wang, and Yu Jiang. 2022. Griffin : Grammar-Free DBMS Fuzzing. In *37th IEEE/ACM International Conference on Automated Software Engineering, ASE 2022, Rochester, MI, USA, October 10-14, 2022*. ACM, 49:1–49:12. doi:10.1145/3551349.3560431
 - [14] Jim Gray, Prakash Sundaresan, Susanne Englert, Kenneth Baclawski, and Peter J. Weinberger. 1994. Quickly Generating Billion-Record Synthetic Databases. In *Proceedings of the 1994 ACM SIGMOD International Conference on Management of Data, Minneapolis, Minnesota, USA, May 24-27, 1994*, Richard T. Snodgrass and Marianne Winslett (Eds.). ACM Press, 243–252. doi:10.1145/191839.191886
 - [15] Urs Hölzle, Craig Chambers, and David Ungar. 1991. Optimizing dynamically-typed object-oriented languages with polymorphic inline caches. In *ECOOP'91 European Conference on Object-Oriented Programming: Geneva, Switzerland, July 15–19, 1991 Proceedings* 5. Springer, 21–38.
 - [16] Dongxu Huang, Qi Liu, Qiu Cui, Zhuhe Fang, Xiaoyu Ma, Fei Xu, Li Shen, Liu Tang, Yuxing Zhou, Menglong Huang, Wan Wei, Cong Liu, Jian Zhang, Jianjun Li, Xuelian Wu, Lingyu Song, Ruoxi Sun, Shuaipeng Yu, Lei Zhao, Nicholas Cameron, Liqian Pei, and Xin Tang. 2020. TiDB: A Raft-based HTAP Database. *Proc. VLDB Endow.* 13, 12 (2020), 3072–3084. doi:10.14778/3415478.3415535
 - [17] Zu-Ming Jiang and Zhendong Su. 2024. Detecting Logic Bugs in Database Engines via Equivalent Expression Transformation. In *18th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2024, Santa Clara, CA, USA, July 10-12, 2024*, Ada Gavrilovska and Douglas B. Terry (Eds.). USENIX Association, 821–835. <https://www.usenix.org/conference/osdi24/presentation/jiang>
 - [18] Jinho Jung, Hong Hu, Joy Arulraj, Taesoo Kim, and Woon-Hak Kang. 2019. APOLLO: Automatic Detection and Diagnosis of Performance Regressions in Database Systems. *Proc. VLDB Endow.* 13, 1 (2019), 57–70. doi:10.14778/3357377.3357382
 - [19] Shadi Abdul Khalek, Bassem Elkarablieh, Yai O. Laleye, and Sarfraz Khurshid. 2008. Query-Aware Test Generation Using a Relational Constraint Solver. In *23rd IEEE/ACM International Conference on Automated Software Engineering (ASE 2008), 15-19 September 2008, L'Aquila, Italy*. IEEE Computer Society, 238–247. doi:10.1109/ASE.2008.34
 - [20] Zuhair Khayyat, William Lucia, Meghna Singh, Mourad Ouzzani, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Nan Tang, and Panos Kalnis. 2015. Lightning Fast and Space Efficient Inequality Joins. *Proc. VLDB Endow.* 8, 13 (2015), 2074–2085. doi:10.14778/2831360.2831362
 - [21] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. 2018. Evaluating Fuzz Testing. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (Toronto, Canada) (CCS '18)*. Association for Computing Machinery, New York, NY, USA, 2123–2138. doi:10.1145/3243734.3243804
 - [22] Vu Le, Mehrdad Afshari, and Zhendong Su. 2014. Compiler validation via equivalence modulo inputs. In *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14, Edinburgh, United Kingdom - June 09 - 11, 2014*, Michael F. P. O'Boyle and Keshav Pingali (Eds.). ACM, 216–226. doi:10.1145/2594291.2594334
 - [23] Vu Le, Chengnian Sun, and Zhendong Su. 2015. Finding deep compiler bugs via guided stochastic program mutation. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015, part of SPLASH 2015, Pittsburgh, PA, USA, October 25-30, 2015*, Jonathan Aldrich and Patrick Eugster (Eds.). ACM, 386–399. doi:10.1145/2814270.2814319
 - [24] Yu Liang, Song Liu, and Hong Hu. 2022. Detecting Logical Bugs of DBMS with Coverage-based Guidance. In *31st USENIX Security Symposium, USENIX Security 2022, Boston, MA, USA, August 10-12, 2022*, Kevin R. B. Butler and Kurt Thomas (Eds.). USENIX Association, 4309–4326. <https://www.usenix.org/conference/usenixsecurity22/presentation/liang>

- [25] Danushka Liyanage, Marcel Böhme, Chakkrit Tantithamthavorn, and Stephan Lipp. 2023. Reachable Coverage: Estimating Saturation in Fuzzing. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 371–383. doi:10.1109/ICSE48619.2023.00042
- [26] Muhammad Numair Mansur, Maria Christakis, and Valentin Wüstholtz. 2021. Metamorphic testing of Datalog engines. In *ESEC/FSE '21: 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Athens, Greece, August 23-28, 2021*, Diomidis Spinellis, Georgios Gousios, Marsha Chechik, and Massimiliano Di Penta (Eds.). ACM, 639–650. doi:10.1145/3468264.3468573
- [27] Steven S. Muchnick. 1997. *Advanced Compiler Design and Implementation*. Morgan Kaufmann.
- [28] Andrew Pavlo and Matthew Aslett. 2016. What's Really New with NewSQL? *SIGMOD Rec.* 45, 2 (2016), 45–55. doi:10.1145/3003665.3003674
- [29] Mark Raasveldt and Hannes Mühleisen. 2020. Data Management for Data Science-Towards Embedded Analytics.. In *CIDR*.
- [30] Manuel Rigger and Zhendong Su. 2020. Detecting optimization bugs in database engines via non-optimizing reference engine construction. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1140–1152.
- [31] Manuel Rigger and Zhendong Su. 2020. Finding bugs in database systems via query partitioning. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 1–30.
- [32] Manuel Rigger and Zhendong Su. 2020. Testing Database Engines via Pivoted Query Synthesis. In *14th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2020, Virtual Event, November 4-6, 2020*. USENIX Association, 667–682. <https://www.usenix.org/conference/osdi20/presentation/rigger>
- [33] Praveen Seshadri, Hamid Pirahesh, and T. Y. Cliff Leung. 1996. Complex Query Decorrelation. In *Proceedings of the Twelfth International Conference on Data Engineering, February 26 - March 1, 1996, New Orleans, Louisiana, USA*, Stanley Y. W. Su (Ed.). IEEE Computer Society, 450–458. doi:10.1109/ICDE.1996.492194
- [34] Donald R. Slutz. 1998. Massive Stochastic Testing of SQL. In *Vldb'98, Proceedings of 24rd International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA*, Ashish Gupta, Oded Shmueli, and Jennifer Widom (Eds.). Morgan Kaufmann, 618–622. <http://www.vldb.org/conf/1998/p618.pdf>
- [35] Jiansen Song, Wensheng Dou, Ziyu Cui, Qianwang Dai, Wei Wang, Jun Wei, Hua Zhong, and Tao Huang. 2023. Testing database systems via differential query execution. In *Proceedings of IEEE/ACM International Conference on Software Engineering (ICSE)*.
- [36] Rebecca Taft, Irfan Sharif, Andrei Matei, Nathan VanBenschoten, Jordan Lewis, Tobias Grieger, Kai Niemi, Andy Woods, Anne Birzin, Raphael Poss, Paul Bardea, Amruta Ranade, Ben Darnell, Bram Gruneir, Justin Jaffray, Lucy Zhang, and Peter Mattis. 2020. CockroachDB: The Resilient Geo-Distributed SQL Database. In *Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020*, David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo (Eds.). ACM, 1493–1509. doi:10.1145/3318464.3386134
- [37] Xiu Tang, Sai Wu, Dongxiang Zhang, Feifei Li, and Gang Chen. 2023. Detecting Logic Bugs of Join Optimizations in DBMS. *Proc. ACM Manag. Data* 1, 1 (2023), 55:1–55:26. doi:10.1145/3588909
- [38] Dominik Winterer, Chengyu Zhang, and Zhendong Su. 2020. Validating SMT solvers via semantic fusion. In *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020*, Alastair F. Donaldson and Emina Torlak (Eds.). ACM, 718–730. doi:10.1145/3385412.3385985
- [39] Andreas Zeller and Ralf Hildebrandt. 2002. Simplifying and Isolating Failure-Inducing Input. *IEEE Trans. Software Eng.* 28, 2 (2002), 183–200. doi:10.1109/32.988498
- [40] Rui Zhong, Yongheng Chen, Hong Hu, Hangfan Zhang, Wenke Lee, and Dinghao Wu. 2020. SQUIRREL: Testing Database Management Systems with Language Validity and Coverage Feedback. In *CCS '20: 2020 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, USA, November 9-13, 2020*, Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna (Eds.). ACM, 955–970. doi:10.1145/3372297.3417260

Received July 2024; revised September 2024; accepted November 2024