# Scaling Automated Database System Testing

Suyang Zhong
suyang@u.nus.edu
National University of Singapore
Singapore, Singapore

Manuel Rigger
rigger@nus.edu.sg
National University of Singapore
Singapore, Singapore

## Abstract

Recently, various automated testing approaches have been proposed that use specialized test oracles to find hundreds of logic bugs in mature, widely-used Database Management Systems (DBMSs). These test oracles require database and query generators, which must account for the often significant differences between the SQL dialects of these systems. Since it can take weeks to implement such generators, many DBMS developers are unlikely to invest the time to adopt such automated testing approaches. In short, existing approaches fail to scale to the plethora of DBMSs. In this work, we present both a vision and a platform, *SQLancer++*, to apply test oracles to any SQL-based DBMS that supports a subset of common SQL features. Our technical core contribution is a novel architecture for an *adaptive SQL statement generator*. This adaptive SQL generator generates SQL statements with various features, some of which might not be supported by the given DBMS, and then learns through interaction with the DBMS, which of these are understood by the DBMS. Thus, over time, the generator will generate mostly valid SQL statements. We evaluated *SQLancer++* across 18 DBMSs and discovered a total of 196 unique, previously unknown bugs, of which 180 were fixed after we reported them. While *SQLancer++* is the first major step towards scaling automated DBMS testing, various follow-up challenges remain.

***CCS Concepts:*** • **Software and its engineering** → **Software testing and debugging**; • **Information systems** → **Data management systems**.

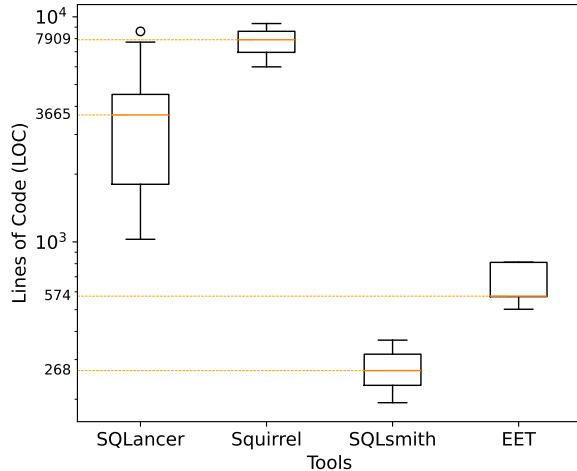***Keywords:*** Database systems; automated testing; test case generation

## 1 Introduction

Database Management Systems (DBMSs) are large, complex software systems. For example, MySQL consists of more than 5.5 million lines of code (LOC), and PostgreSQL has more than 1.7 million LOC. Unsurprisingly, such DBMSs can be affected by bugs. Automated testing approaches for DBMSs have been proposed to find so-called *logic bugs* [2, 15, 29, 35–37, 43, 47], which are bugs that cause a system to silently compute an incorrect result, making such bugs notoriously difficult to find. Many existing works aiming to find logic bugs proposed so-called *test oracles* that can validate whether a query computes the correct result by transforming it in a semantic-preserving way and checking both queries' results' equivalence. Overall, these approaches have found hundreds of bugs in widely-known DBMSs such as SQLite, MySQL, and PostgreSQL.

It would be ideal to apply automated DBMS testing approaches to the thousands of existing DBMSs.[1] The market for DBMSs is significant, currently being 162.25 billion USD and growing at a compound annual growth rate of 15.2% [7], fueling the development of new DBMSs, as well as further development of existing ones. With the end of Moore's law, various trends have set in posing new reliability challenges, such as the development of new, increasingly specialized DBMSs, often based on SQL and the relational model. In addition, existing DBMSs are becoming increasingly complex, by using accelerators [17, 23, 42], or incorporating learned components [16, 24, 32]. Overall, ensuring the correctness of DBMSs will become increasingly difficult.

Despite the significant, increasing need to apply automated testing approaches for DBMSs at scale, much effort is needed to implement and operate them. A key challenge for current automated testing approaches is that test-case generators, which generate SQL statements that create a database, populate it with data, and execute queries, must be manually implemented to account for the differences in SQL dialects [35, 38]. Doing so typically requires implementing thousands of lines of code. For example, SQLancer [2, 3, 35–37], a state-of-the-art tool for DBMS testing, currently supports generators for 22 DBMSs, which, on average, are implemented in 3,729 LOC (see Figure 1), with some DBMS-specific components being contributed by major companies. However, most DBMS development teams are unlikely to invest this effort. For example, a Vitess blog post describes that they considered using SQLancer, but realized that *"It would take a*

---

[1]See https://dbdb.io/

**Figure 1.** The lines of code needed in popular DBMS testing tools [22, 37, 38, 54] for DBMS-specific components, where the orange line shows the median. Note that we only considered the core part of the generator of each DBMS.

*lot of work to properly integrate Vitess with SQLancer, due to each DBMS tester in SQLancer essentially being written completely separately with similar logic."* [48] While Vitess aims to be MySQL-compatible, and SQLancer provides a MySQL implementation, in practice, various differences still exist that make it difficult to reuse generators in such a scenario.

In this work, we introduce our vision of applying testing DBMSs *at scale.* We propose a new automated testing platform, *SQLancer++,* as part of our ongoing effort to reduce the effort of implementing and operating DBMS testing approaches. The main technical contribution of this work is an adaptive query generator, which infers to generate statements that are understood by the DBMS under test. This generator repeatedly sends SQL test cases to the DBMS; if a SQL feature is not understood, the generator will prevent its generation and thus, over time, increase the validity rate of generated statements. Such a generator is unlikely to perform as well as one that is implemented for a specific system; however, it can be directly applied without additional implementation effort. In addition, any added features might also be immediately applicable to other DBMSs. Other key design decisions include an internal model of the database schema, to eschew avoiding querying schema information from DBMSs, which often requires using DBMS-specific interfaces. Finally, we propose a synergistic and pragmatic bug-prioritization approach, which reports bug-inducing tests only if its features are not present in previously reported bugs. Overall, we believe that these techniques are both simple and practical.

Our results are promising, as we have found and reported 196 unique bugs across 18 DBMSs, demonstrating that the approach is effective and scalable. As our bug-finding efforts

on these systems have not yet saturated, we are still reporting more bugs. 180 bugs have already been fixed, indicating that the developers considered the bugs important. The feedback mechanism in the adaptive generator significantly increased the validity rate of statements; for SQLite, by 292.5%. Similarly, the prioritization approach has shown promise. In CrateDB, *SQLancer++* identified over 60K bug-inducing cases in an hour on average, which the bug prioritizer reduced to 35.8, with 11.4 being unique bugs.

*SQLancer++* is the first step towards an automated testing platform for DBMSs that can find bugs at a large scale, similar to platforms such as OSS-Fuzz [39] and syzkaller [14], which tackled the problem of scaling general-purpose greybox fuzzers like AFL, as well as kernel fuzzing. Overall, we hope that *SQLancer++* will have a significant impact in improving the reliability of DBMSs at large and, in particular, help smaller DBMS development teams that cannot invest significant resources into testing. To facilitate future research and adoption, we have released *SQLancer++* at https://doi.org/10.5281/zenodo.18289297.

In summary, we propose the following:

- At a conceptual level, we have identified the problem of scaling existing DBMS testing approaches and propose a new testing platform, *SQLancer++* to apply previously-proposed test oracles at scale.
- At a technical level, we propose an adaptive query generator, a model-based schema representation, as well as a bug-prioritization approach for logic bugs.
- At an empirical level, we show that these techniques are highly effective, and already found 196 unique, previously unknown bugs in popular DBMSs.

## 2  Background and Motivation

*SQLancer* is among the most popular automated testing tools for DBMSs, supports over 20 DBMSs, and is widely used by DBMS developers. While the *SQLancer* authors proposed various approaches built on it [2–4, 35–37], also other tools have been prototyped on it [30, 34, 44, 45]. Thus, we focus on *SQLancer* to motivate the limitations of existing approaches. Different from other fuzzing approaches, such as SQLsmith, *SQLancer* implements various state-of-the-art test oracles [2, 3, 35–37] for finding logic and performance issues.

**SQL generators.** At the core of automated DBMS testing tools like *SQLancer* are DBMS-specific, rule-based SQL generators. Listing 1 shows a simplified excerpt of a *SQLancer* generator for CREATE INDEX statements. Line 2 uses a string to represent a CREATE statement, to which, subsequently, other elements such as INDEX are added. To create an index, schema information is required. Specifically, line 7 determines an index name that is not yet present in the database, and line 9 determines an existing table, on which the index should be created. *SQLancer*—as well as other DBMS

**Listing 1.** Simplified excerpt of `PostgresIndexGenerator`.

```
1  public SQLStatement createIndex(PostgresState state) {
2    String stmt = "CREATE ";
3    if (state.nextBoolean()) {
4      stmt += "UNIQUE ";
5    }
6    stmt += "INDEX ";
7    stmt += state.getSchema().getFreeIndexName();
8    stmt += " ON ";
9    stmt += state.getSchema().getRandomTableName();
10   stmt += "(" + table.getRandColumnNames() + ")";
11   if (state.nextBoolean()) {
12     stmt += " WHERE ";
13     stmt += state.getRandomExpr(table.getColumns());
14   }
15   return new SQLStatement(stmt);
16 }
```

testing tools like SQLsmith or Griffin—query this information from the DBMS under test. As shown by the call to `random.nextBoolean()` in lines 3 and 11, it is randomly decided whether to create a unique index by including the `UNIQUE` keyword, or a partial index by including the `WHERE` keyword; typically, in *SQLancer*, the probabilities of choosing a feature are uniformly distributed.

***Case study.*** To motivate the challenges, we assumed that the CrateDB development team would want to adopt one DBMS testing tool to test their system. CrateDB is a PostgreSQL-compatible system; luckily, a PostgreSQL generator already exists within *SQLancer*, and it would be reasonable to assume that the CrateDB developers would want to adapt it, rather than creating one without reference. However, even given this best-possible scenario, our case study demonstrates that the adaptation effort is still high. In order to minimally adapt the generator, we added 309 and deleted 987 LOC, resulting in 1296 LOC being modified to avoid syntax errors and acquire meta-data for schema information. Performing these changes is non-trivial, requiring detailed knowledge of *SQLancer*. Several other DBMS testing tools exist, such as *Griffin* [9] for grammar-free testing, and *SQL-Right* [29] for detecting logic bugs. However, both tools were built upon AFL [50], making them inapplicable to DBMS written in Java, like CrateDB. While Java-based fuzzers like JQF [33] or Jazzer [1] could be used for Java-based DBMSs, these fuzzers are agnostic to SQL, and would thus struggle with generating meaningful sequences of SQL statements.

***C1. SQL dialects.*** First, the SQL generators are highly DBMS-specific as they must account for differences in statements, operators, functions, and data types. The above generator is specific to PostgreSQL, which supports partial indexes, which are, for example, unavailable in MySQL. The full versions of such generators typically contain many uncommon keywords, testing which is desirable, as unique functionalities are prevalent and might be affected by bugs. For example, index generators might generate additional DBMS-specific

index types (*e.g.*, HASH indexes in PostgreSQL), or use DBMS-specific operators in expressions. Thus, applying *SQLancer* or other DBMS testing tools to a new DBMS requires significant implementation effort in implementing new generators, as shown in Figure 1. Regarding our case study, despite CrateDB claiming to be dialect-compatible with PostgreSQL, we had to modify 846 LOC to avoid syntax errors. If these syntax errors were simply ignored, an experiment we performed suggests that the validity rate of queries would drop to less than 1%, significantly affecting *SQLancer*'s performance. Going beyond a minimal implementation, the CrateDB developers would subsequently likely also need to add CrateDB-specific functionality, incurring additional effort.
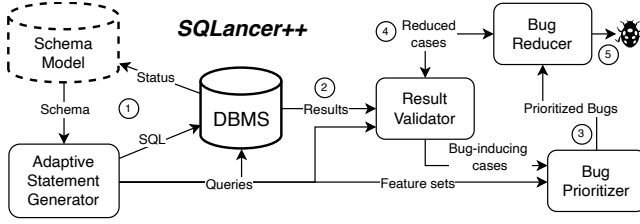
***C2. Schema state.*** Second, acquiring metadata for schema information differs across DBMSs. In SQLite, `sqlite_master` holds the schema information. However, PostgreSQL, for example, exposes similar information in `information_schema.tables`, and MySQL-like systems provide SQL statements like `SHOW TABLE`. In our case study, we found that even when DBMSs share the same metadata tables (*e.g.*, `information_schema` in CrateDB), the attributes of these tables can differ, requiring the logic that parses the schema to be changed. Besides, several features are related to metadata tables that are not fully supported by CrateDB (*e.g.*, collation information in `pg_collation`), and we manually omitted these. Note that, as with *SQLancer*, several other tools—such as Griffin [9] and DynSQL [20]—also acquire schema state through a similar interface.

***C3. Duplicate bug-inducing test cases.*** Third, *SQLancer* has no deduplication mechanism and is likely to repeatedly trigger the same underlying bugs. This is problematic since manual analysis is necessary to determine whether two bug-inducing test cases likely trigger the same underlying bug, to avoid overburdening developers by reporting duplicate bugs. In our case study, running the generator on a historic version of CrateDB for one hour would result in more than 400 bug-inducing test cases. Whether two bug-inducing test cases trigger the same logic bug is decided by the developers' verdict; it is not an inherent property of the test cases, meaning that the problem cannot be fully addressed. Existing techniques such as deduplication based on stack traces cannot be directly applied to logic bugs [6, 19]. Despite this, we believe a pragmatic solution is required to deal with this challenge in practice.

## 3  SQLancer++

*SQLancer++* is an automated DBMS testing platform aiming to find logic bugs in DBMSs (see Figure 2), addressing the scalability limitations of existing rule-based approaches. The core component of *SQLancer++* is an *adaptive statement generator*, which infers the supported SQL features of the

**Figure 2.** Architecture of the *SQLancer++*



**Figure 3.** Schema model after executing each DDL statement

target DBMS during execution and adaptively generates random SQL statements that can be processed by the DBMS under test. First, the adaptive generator generates SQL statements to create a database state while maintaining an internal model of the schema. Then, the generator generates random queries and executes them on the DBMS, and a feedback mechanism automatically prioritizes the supported features and deprioritize the unsupported features to ensure most queries are semantically valid. After retrieving the queries' results, an oracle validator checks whether the results are as expected. When detecting issues, *SQLancer++* compares the SQL features with the previously found bug-inducing test cases to prioritize the bug-inducing test cases for reporting. If a previous bug-inducing feature set is a subset of the new test case, it is marked as a potential duplicate.

**SQL features.** The notion of SQL *features* is an important, but abstract concept in both the generator and bug prioritizer. A feature refers to an element or property in the query language, which we expect to be either supported or unsupported by a given DBMS. More concretely, developers can mark any code block in a generator in *SQLancer++* as a potential feature. This instructs the generator to determine whether the feature is supported by the DBMS (*i.e.*, whether statements based on which the feature is generated execute successfully), and only if so, continue generating the feature. In addition, it instructs the prioritizer to use it for bug-prioritization. Features can be specified in different granularities. Often, a feature might be a specific keyword or operator. However, it can also refer to a class of functions, for example, string functions. Finally, it can be used to specify abstract properties, such as whether the DBMS provides implicit conversions between most types.

**Adaptive statement generator.** The core of *SQLancer++* is an adaptive statement generator detailed in Section 4, which tackles C1, the differences across DBMSs' SQL dialects. The generator produces SQL statements with SQL features that might not be supported by all DBMSs, or, whose constraints are difficult to meet, causing statements containing them to frequently fail. It learns through execution feedback from the DBMS under test. Users can set a minimum success probability threshold for a feature, which is estimated through *Bayesian inference* [11] by utilizing feedback from
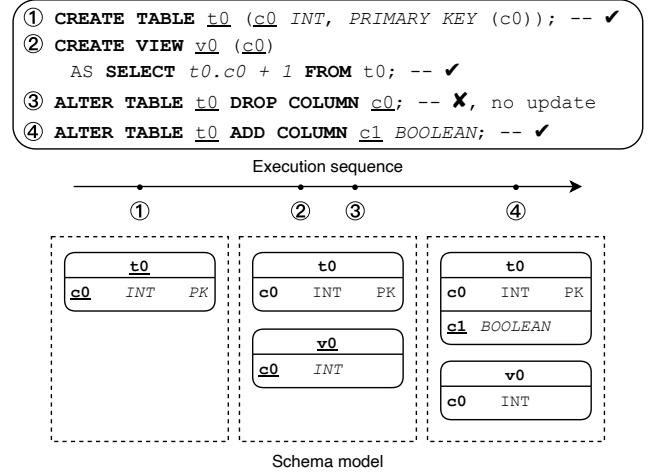
previous executions. This suppresses the generation of frequently failing SQL features to increase the validity rate of SQL statements.

**Schema model.** To address C2, the challenge of obtaining database schemas, the generator maintains an internal data model for the database, enabling it to obtain schema information without querying the DBMS. The internal data model matches the DBMSs' model by simulating the behavior of DDL statements generated and checking the execution status from the DBMS. Initially, the objects (*i.e.*, tables, views, or indexes) in the model are $O = \{\}$. When executing a statement $s$ that potentially creates an object $T$, such as a CREATE TABLE statement creating a table with a specific name, we obtain the execution status from the DBMS. If the object was successfully created, we add the object to our model of the schema such that $O = O \cup T$; otherwise, we consider the creation of the table $T$ unsuccessful.

As shown in Figure 3, when generating a CREATE TABLE statement (①), we create an internal schema model of the table, which records its name (t0) along with its columns (c0) and corresponding data types (INT). Once the statement is successfully executed on the DBMS, this abstract object is added to the schema model. During subsequent testing, the schema model can be queried to retrieve available tables and their columns. Potential bugs could cause DBMSs to maintain an incorrect schema state, in which case our internal schema model could deviate from it. However, during our fuzzing campaign, we have not observed any related false positives.

**Result validator.** *SQLancer++* can be combined with any test oracle that is not specific to a DBMS. We adopted two state-of-the-art [10] test oracles, *Ternary Logic Partitioning* (TLP) [36] and *Non-optimizing Reference Engine Construction* (NoREC) [35], which can find logic bugs by comparing the results of two equivalent queries constructed through
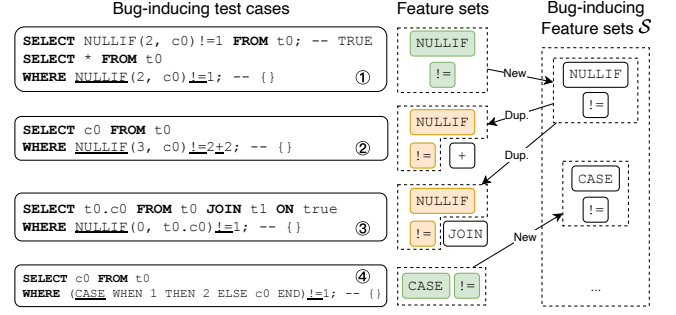
a syntactic transformation that applies to any DBMS. The validator fetches the results of two equivalent queries and validates if the results are the same. If not, it reports a bug.

***Bugs prioritizer.*** *SQLancer++* employs a bug prioritization approach to prioritize bug-inducing test cases for analysis and reporting, aiming to minimize the amount of duplicate bugs being reported, tackling C3. We prioritize the bugs by comparing SQL feature sets from previous bug-inducing test cases with the feature set from the newly found bug-inducing test case. We define the *feature set S* as a set of features $\{F_1, F_2, \cdots, F_n\}$, where each $F_i$ denotes one SQL feature that was enabled when generating a SQL statement. Specifically, assume that $\mathcal{S}$ represents the bug-inducing feature sets of the historical bugs, and $S_{new}$ denotes the feature set of the newly detected bug-inducing test case. Should there exist an $S$ within $\mathcal{S}$ such that $S \subseteq S_{new}$, then this test case is considered a potential duplicate, based on the intuition that the root cause of the bugs might be the faulty implementation of the SQL features in the bug-inducing test cases triggering the bugs. Figure 4 shows an example of the bug prioritization. When triggering a bug with feature set $S_{new}$ = {NULLIF, !=}, we retrieve $\mathcal{S}$ and find no set is a subset of $S_{new}$. Thus, we consider this a new bug and update $\mathcal{S}$. Subsequently, feature sets of test cases ② and ③ contain {NULL, !=}, and we mark these cases as potential duplicated, meaning that would analyze and report them only once existing bugs have been fixed by the developers.

Given two bug-inducing test cases, the generator might misclassify them both as triggering the same underlying bug when they are not, or misclassify them as triggering different bugs when they trigger the same. For example, the expression NULLIF(2, c0)<>1 may trigger a bug involving the feature set {NULLIF, <>}, which would be considered new although its root cause might be identical to a previously reported issue. However, we believe our method is effective for reducing the large volume of reported bugs—potentially hundreds or thousands in one hour—when applying *SQLancer++* to an untested system (see Section 5.5). False negatives can occur, which means two different bugs affect the statements with the same feature sets. However, these do not result in missed bugs; after developers fix the prioritized bugs, they can run the reproducer to reproduce the detected bugs again quickly to check if there are any false negatives.

## 4 Adaptive Statement Generator

Figure 5 illustrates the steps of adaptive statement generation, which is at the core of *SQLancer++*. First, we initialize the generator with each feature sharing the same probability as the other options in the respective context (see step ①). Next, we randomly generate statements and queries, whose results are checked by the provided test oracle. Expectedly, some of these statements might contain features that are not supported by the DBMS under test. To identify these, we



**Figure 4.** Bug prioritization. We omitted the unoptimized queries in test cases ② to ④ for simplicity.

generate each statement one by one, execute it on the DBMS, and obtain its execution status, which indicates whether the statement could be executed successfully, whether it failed (see step ②). Each feature in the feature set is marked as valid (shown in green) if the statement succeeds, and as failed otherwise (shown in red). We repeat this process for an empirically-determined number of iterations. In step ③, we calculate the estimated probability of successfully executing each SQL feature, using *Bayesian inference* [11] based on previous execution feedback. Features are marked as *unsupported*—unlikely to execute successfully—for the target system if their estimated probability falls below a specified threshold; otherwise, they will be regarded as *supported*. In step ④, we update the generator based on the state of each feature. Unsupported features are subsequently suppressed, and the probabilities for other alternatives are distributed uniformly. The probabilities from step ④ can be persisted in a file and loaded in step ① of future executions.

***SQL statement generation.*** The generator generates random SQL statements based on the specified rule and records the corresponding feature set. By default, the selection of each feature under each grammar rule is uniformly distributed (see step ①). For example, <=> is a feature in the generator for generating a null-safe comparison expression in the SQL statement. Correspondingly, we add this feature to the feature set such that $S = S \cup \{<=>\}$.

***Validity feedback.*** We execute the generated statements in step ② and record their execution statuses as feedback. If the DBMS encounters a syntax or semantic error—whether caused by unsupported features, constraints, or any other issue—it returns an error message, causing that statement to fail. For example, in Figure 5, the DBMS under test does not support indexes, which is why the first statement fails. The third SQL statement executes successfully, and it contains features CASE, !=, SIN and SIN1INT. Here, SIN1INT is a composite feature for data types, which we will describe later. Consequently, we mark the feature set {INDEX} as *failed* and {CASE, !=, SIN, SIN1INT} as *succeeded*. The same feature might fail in one context but succeed in another
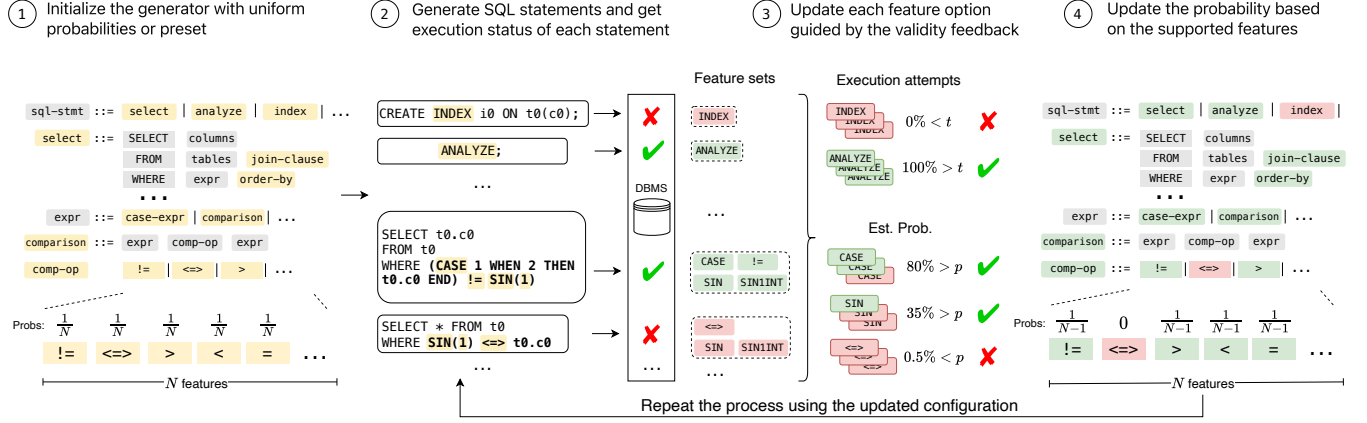
**Figure 5.** Overview of the adaptive statement generator

(*e.g.*, `ASIN(1)` can succeed in PostgreSQL while `ASIN(2)` will throw an error). If a feature that rarely results in a valid SQL statement is generated frequently, it lowers the overall success rate of statements, preventing the DBMS's core logic from being exercised, which is critical for finding logic bugs. Thus, the next steps aim to analyze the feature sets to infer whether features are *supported* (*i.e.*, executed successfully).

**Feedback mechanism.** In step ③, we estimate the probability that the respective features are supported by the DBMS via validity feedback. Since our bug detection operates in two phases—first establishing the database state via DDL (Data Definition Language) or DML (Data Manipulation Language) statements, and then issuing a large number of random queries—we evaluate features in these two categories separately. For DDL/DML statements, we apply straightforward tests: if a feature repeatedly fails beyond a user-specified number of attempts, it is deemed unsupported. For queries, we employ a simple *Bayesian inference* [11] model to estimate the posterior probability distribution indicating whether each query feature is *supported*. Users may specify a threshold $p$ (*e.g.*, 1%) for the least query successful probability, and if the posterior probability mass for a given feature predominantly falls below this threshold, *SQLancer++* marks this feature as unsupported.

**Statistical model.** For each iteration, we record each feature's total number of executions $N$ and successful executions $y$. Queries are generated independently, unlike test case generation methods with guidance (*e.g.*, code coverage). Consequently, each feature has a constant probability of being included in any given query. Assuming that every feature behaves deterministically and leaves the database state unaltered, the probability of success for each feature is independent and identically distributed among each iteration. This scenario can be modeled using a binomial distribution,

where the probability mass function is defined as

$$p(y|\theta) = \binom{N}{y}\theta^y(1-\theta)^{N-y} \tag{1}$$

with $\theta$ representing the probability that the feature is executed successfully. To assess whether a feature can be executed successfully, we calculate the posterior probability distribution $p(\theta|y)$. In the worst-case scenario, we assume a uniform prior distribution for $\theta$ over the interval $[0, 1]$, implying that $p(\theta) = 1$. Applying Bayes' theorem, we obtain

$$p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)} \tag{2}$$

where $p(y) = \int_0^1 p(\theta)p(y|\theta)d\theta$. Substituting Equation 1 into Equation 2, we find

$$p(\theta|y) \propto \theta^y(1-\theta)^{N-y} \tag{3}$$

indicating that $\theta|y$ follows a Beta distribution: $\theta|y \sim \text{Beta}(y+1, N-y+1)$. Consider a user-specified threshold $p = 0.01$, which indicates that the success probability should be at least 1%. In the case where $y = 0, N = 400$—meaning the feature has been executed 400 times without any successful executions—the posterior distribution is then $\text{Beta}(1, 401)$. The 95% credible interval of this posterior is approximately $[6 \times 10^{-5}, 0.009]$. Since more than 95% of the posterior probability mass is below 0.01, the feature is deemed unsupported. Lowering the threshold $p$ would necessitate additional executions to achieve a similar level of confidence.

**Generator specification.** In step ④, we update the probability of the generator choices based on the states of the SQL features. Specifically, our goal is to prevent unsupported features from being generated. For example, the comparison operator "`<=>`" is not supported for the DBMS under test in the example shown in Figure 5. We assign zero probability to "`<=>`". Considering that the total number of features in this production rule is $N$, and that the other alternatives are uniformly distributed, each receives a probability of $1/(N-1)$.

**Table 1.** Comparison of related DBMS testing approaches. Base tools are shown in bold. *Eval.* denotes the number of DBMSs reported in the research paper. ○ indicates close-to-zero manual effort for adapting the tool to a new DBMS, while the ● indicates that manual adaptation is necessary.

| Tool | Bug Type | | DBMS | | | Manual |
| | Crash | Logic | C/C++ | Non C | Eval. | Efforts |
|---|---|---|---|---|---|---|
| **AFL** [50] | ✓ | - | ✓ | - | - | ○ |
| - *Griffin* [9] | ✓ | - | ✓ | - | 4 | ○ |
| - *WingFuzz* [27] | ✓ | - | ✓ | - | 12 | ○ |
| - *SQLRight* [29] | ✓ | ✓ | ✓ | - | 3 | ● |
| **SQLSmith** [38] | ✓ | - | ✓ | ✓ | - | ● |
| - *EET* [22] | ✓ | ✓ | ✓ | ✓ | 5 | ● |
| **SQLancer** [37] | ✓ | ✓ | ✓ | ✓ | - | ● |
| - *TLP* [36] | ✓ | ✓ | ✓ | ✓ | 6 | ● |
| - *NoREC* [35] | ✓ | ✓ | ✓ | ✓ | 4 | ● |
| - *CODDTest* [52] | ✓ | ✓ | ✓ | ✓ | 5 | ● |
| *SQLancer++* | ✓ | ✓ | ✓ | ✓ | 18 | ○ |

Based on the updated rule in step ④, the generator has a higher probability of selecting other comparison operators, avoiding generating unsupported "<=>", which increases the validity rate of statements.

## 5 Evaluation

We sought to understand how effective and efficient *SQLancer++* is in finding bugs in DBMSs at a large scale from these perspectives: the overall bug-finding effectiveness (see Section 5.1), the impact of SQL features (see Section 5.2), bug-finding efficiency (see Section 5.3), feedback mechanism effectiveness (see Section 5.4), and bug-prioritization effectiveness (see Section 5.5).

**Baselines.** We implemented *SQLancer++* in Java in 8.4K LOC. We compared *SQLancer++* with *SQLancer* and *SQLRight*, which are state-of-the-art logic bug detection tools. In comparison, *SQLancer* has 83K lines of Java code, and the core of *SQLRight* has 30K lines of C++ code. Table 1 demonstrates the other DBMS testing tools. Our aim for *SQLancer++* was to detect logic bugs in various DBMSs with different dialects; thus, for a given DBMS that is supported by *SQLancer*, we would expect it to perform better than *SQLancer++* due to its manually-written generators specific to that system. Our main use case is to apply *SQLancer++* to SQL dialects that are not supported by existing DBMSs. We did not compare with other AFL-based fuzzers, including Griffin [9], WingFuzz [27], and BUZZBEE [49], since they cannot detect logic bugs or require instrumentation, which is limited to C/C++ systems, whereas many DBMSs are written in other languages (*e.g.*, CrateDB is written in Java).

**Experiments Setup.** We conducted the experiments using a server with a 64-core AMD EPYC 7763 CPU at 2.45GHz and 512GB memory running Ubuntu 22.04. In the bug-finding

campaign, we enabled multi-threading; however, for other experiments, we evaluated *SQLancer++* using a single thread to prevent feedback sharing across threads, as this could mask the specific impact of individual threads. We set the maximum expression depth to three and randomly created up to two tables and one view, respectively, which are the standard settings for *SQLancer*. We also compared *SQLancer++* under different feedback mechanisms. "*SQLancer++*" denotes the enabling of feedback, and "*SQLancer++$_{Rand}$*" indicates the absence of feedback guiding the generator.

**DBMS selection.** We considered 18 DBMSs in our evaluation (see Table 2). We first included DBMSs that had been extensively studied in prior work (DuckDB, H2, MariaDB, MonetDB, MySQL, Percona MySQL, SQLite, TiDB, and Virtuoso). We then added systems that were not covered by *SQLancer*, but were popular according to DB-Engines rank (at most 50) or GitHub stars (at least 4K), including Firebird, Oracle, CrateDB, Dolt, RisingWave, and Vitess, and further included Umbra and CedarDB as prominent recent academic systems. To test the effectiveness of *SQLancer++* (see Section 5.1), we tested them in a bug-finding campaign, using their latest available development versions to avoid reporting bugs that had already been fixed. For bug-prioritization effectiveness and bug-finding efficiency, we aimed to select DBMS that could be continuously tested for a longer period while also having a number of bugs that are fixed in the latest version, allowing us to determine whether a bug was unique based on its fix commit. CrateDB met these criteria, as all bugs we found are logic bugs, and we selected CrateDB 5.5.0, a historic version, for testing. To compare with other tools, we selected SQLite 3.45.2 as the baseline for measuring coverage and validity rate and selected PostgreSQL 14.11 for measuring validity rate. To the best of our knowledge, no logic bug-finding tools could be applied to CrateDB. SQLite and PostgreSQL are robust systems, which have been extensively tested and targeted by many bug-finding tools [35, 54]. We did not compare *SQLancer++* with *SQLancer* by testing MySQL, although both tools could find bugs in it, since most bugs found by previous automated testing work remain unfixed [10, 15, 37].

### 5.1 Effectiveness

We continuously tested the 18 DBMSs for about four months of intensive testing, followed by several months of intermittent testing. Typically, we ran *SQLancer++* for several seconds up to multiple hours, and then further processed the automatically-reduced and prioritized bug-inducing test cases. Before reporting them to the developers, we further reduced them manually and checked whether any similar bugs had already been reported to avoid duplicate issues. For the developers of DBMSs who were actively fixing bugs, we reported up to three bug-inducing test cases. For some DBMSs where bugs were not fixed, like for MySQL (as also

pointed out by prior work [35, 36]), we refrained from reporting bugs. After previous bugs were fixed, we started another testing run. Similar testing campaigns have also been conducted to evaluate the effectiveness of testing approaches for DBMSs [21, 29, 37] as well as in other contexts [25, 26, 53].

**Bug statistics.** Table 2 provides detailed statistics on the bugs we reported. In total, we created 196 bug reports, and 180 bugs have been fixed as a direct response to our bug reports, which demonstrates that the DBMS developers considered most of the bugs important. Note that 140 of the bugs are logic bugs, found by the TLP [36] and NoREC [35] oracles, with 133 and 7 found respectively. We found more bugs using TLP, because we first implemented and used it during testing. In comparison, the original TLP and NoREC papers reported 77 and 51 logic bugs. Our general *SQLancer++* showed promising ability in finding bugs, even though the generator was not specifically developed for the DBMS under test. Although we tested the latest development versions, many of the bugs we found also affected the stable release version at the time of testing, such as all SQLite bugs and eight out of ten DuckDB bugs. We reported only 4 duplicate bugs, which shows the effectiveness of our bug prioritization approach, that we were careful with avoiding duplicate issues, and closely worked with the DBMS developers.

**Developer reception.** Developer feedback is an important indicator of whether the bugs reported matter to developers. Our efforts received encouraging feedback. For example, developers from CrateDB were curious about our method and wished to integrate our tool into their development cycle: *"Yes, those bug reports are very helpful. We are interested in an automatic fuzz testing tool and we would integrate it into our development cycle. Please keep us updated!"*[2] Developers from Dolt repeatedly expressed their gratitude, responding with comments such as: *"This is very interesting. Keep the bugs coming. They are awesome."*[3] Many of our reports were confirmed using replies such as *"Thank you for the report"*, from MonetDB[4] and MySQL.[5] We also reported bugs to CedarDB and Vitess, and their developers were highly receptive, expressing interest in adopting our tool to detect bugs. The developers' reception shows the practical impact of our tool. To illustrate this, we discuss multiple selected bugs below.

**SQLite bug in REPLACE function.** Listing 2 shows a bug in the implementation of the REPLACE function, which returned the first argument not a string but an intermediate object due to faulty implementation in SQLite, causing subsequent wrong comparison with the text. We found this bug using the TLP [36] oracle. The original query without a WHERE filter returns exactly one row in the table; however,

---

[2] https://github.com/crate/crate/issues/15083
[3] https://github.com/dolthub/dolt/issues/7018
[4] https://github.com/MonetDB/MonetDB/issues/7429
[5] https://bugs.mysql.com/bug.php?id=113298

**Listing 2.** A bug in the REPLACE function remained undetected in SQLite for 10 years.

```
1  CREATE TABLE t0(c0 TEXT, PRIMARY KEY(c0));
2  INSERT INTO t0 (c0) VALUES (1);
3
4  SELECT * FROM t0 WHERE t0.c0=REPLACE(1, '', 0);
5  -- 1 ✔
6  SELECT * FROM t0
7      WHERE NOT t0.c0=REPLACE(1, '', 0);
8  -- 1 🐛
```

**Listing 3.** A bug in SQLite when handling multiple subqueries.

```
1  CREATE TABLE t0(c0 INT);
2  CREATE TABLE t1(c0 INT);
3  INSERT INTO t0 (c0) VALUES (1);
4  CREATE VIEW v0(c0) AS
5      SELECT 0 FROM t1 RIGHT JOIN t0 ON 1;
6
7  SELECT t0.c0
8      FROM v0 LEFT JOIN (
9          SELECT 'a' AS col0 FROM v0 WHERE false
10     ) AS sub0 ON v0.c0,
11         t0 RIGHT JOIN (
12             SELECT NULL AS col0 FROM v0
13     ) AS sub1 ON t0.c0; -- 1 ✔
14
15 SELECT t0.c0
16     FROM v0 LEFT JOIN (
17         SELECT 'a' AS col0 FROM v0 WHERE false
18     ) AS sub0 ON v0.c0,
19         t0 RIGHT JOIN (
20             SELECT NULL AS col0 FROM v0
21     ) AS sub1 ON t0.c0 WHERE t0.c0; -- {} 🐛
```

the derived partitioned query, whose results should be composed to form the same result, returns two rows. The reason is that the query with the predicate and its negation both returns one row. The fix was twofold: a direct fix to ensure the REPLACE function returns a TEXT value, and a deeper fix to make the comparison operator work normally even when one operand has both text and numeric types, which is a problem that could have been bisected to a commit in 2014. This bug has been hidden for about ten years. Surprisingly, *SQLancer++* found it despite the efforts of *SQLancer* and other logic bug detection tools.

**SQLite bug with subquery.** Listing 3 shows another bug we found in SQLite, also detected by the TLP oracle. The second query should return one row since the predicate t0.c0 should be evaluated to 1. The root cause of the bug is that the query flattener incorrectly changed an ON clause term to a WHERE clause term. *SQLancer* failed to detect this bug since it did not support subqueries, which is a feature
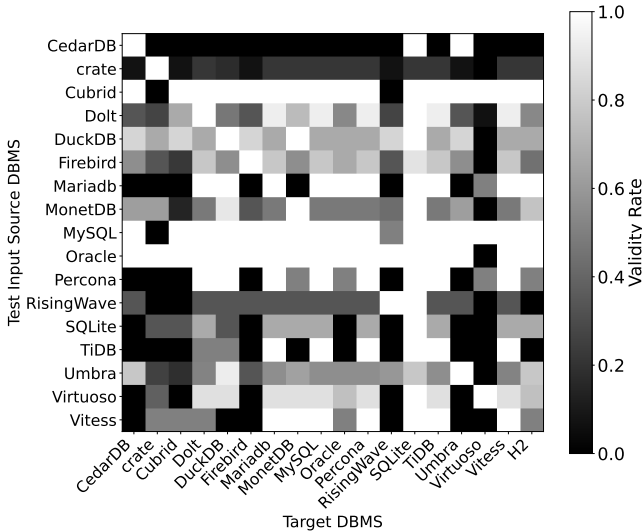
**Table 2.** *SQLancer++* allowed us to find and report 196 bugs in 18 systems, of which 180 were fixed by the developers, and of which 140 were logic bugs. DBMSs are sorted alphabetically.

| DBMS Names | All | Bug Status | | | Test Oracle | | DB-Engines Rankings | GitHub Stars | Released/ Published | Lines of Code | SQLancer Support | C/C++ System |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Fixed | Conf. | Dupl. | Logic | Other | | | | | | |
| CedarDB | 4 | 4 | 0 | 0 | 1 | 3 | - | - | 2024 | - | - | ✓ |
| CrateDB | 28 | 26 | 0 | 2 | 28 | 0 | 229 | 4.0k | 2017 | 597K | - | - |
| Cubrid | 1 | 1 | 0 | 0 | 1 | 0 | 169 | 0.3k | 2008 | 1105K | - | ✓ |
| Dolt | 28 | 27 | 1 | 0 | 16 | 12 | 197 | 16.9k | 2018 | 380K | - | - |
| DuckDB | 10 | 9 | 1 | 0 | 6 | 4 | 76 | 16.4k | 2019 | 1496K | ✓ | ✓ |
| Firebird | 11 | 9 | 1 | 1 | 9 | 2 | 31 | 1.2k | 2000 | 1643K | - | ✓ |
| H2 | 2 | 2 | 0 | 0 | 1 | 1 | 50 | 4.0k | 2005 | 297K | ✓ | - |
| MariaDB | 2 | 0 | 2 | 0 | 2 | 0 | 18 | 5.7k | 2009 | 1087K | ✓ | ✓ |
| MonetDB | 36 | 36 | 0 | 0 | 22 | 14 | 148 | 0.3k | 2004 | 411K | ✓ | ✓ |
| MySQL | 2 | 0 | 2 | 0 | 2 | 0 | 2 | 10.2k | 1995 | 5532K | ✓ | ✓ |
| Oracle | 1 | 0 | 1 | 0 | 1 | 0 | 1 | - | 1977 | - | - | ✓ |
| Percona MySQL | 2 | 0 | 2 | 0 | 2 | 0 | 121 | 1.1k | 2008 | 4182K | - | ✓ |
| RisingWave | 4 | 4 | 0 | 0 | 3 | 1 | 245 | 6.2k | 2022 | 624K | - | - |
| SQLite | 3 | 3 | 0 | 0 | 3 | 0 | 10 | 5.4k | 2000 | 372K | ✓ | ✓ |
| TiDB | 3 | 1 | 2 | 0 | 2 | 1 | 72 | 36.1k | 2016 | 1398K | ✓ | - |
| Umbra | 47 | 46 | 0 | 1 | 31 | 16 | - | - | 2018 | - | - | ✓ |
| Virtuoso | 10 | 10 | 0 | 0 | 8 | 2 | 83 | 0.8k | 1998 | 2659K | - | ✓ |
| Vitess | 2 | 2 | 0 | 0 | 2 | 0 | 200 | 18.5k | 2011 | 1533K | - | - |
| **Total** | 196 | 180 | 12 | 4 | 140 | 56 | | | | | | |



**Figure 6.** The validity rate of executing bug-inducing test cases across different DBMSs. At each intersection of *bug source* and *target DBMS*, the color represents the average success rate for executing bugs found in *bug source* on the *target DBMS*.

newly integrated into *SQLancer++*. Surprisingly, according to the developers' commit log,[6] we found this is a follow-up bug of one bug that EET [22] oracle detected a year ago.

---

[6] https://www.sqlite.org/src/info/bdd408a2557ff05c

## 5.2 SQL Feature Study

We conducted a case study on the 140 logic bugs that we found across different DBMSs to test our hypothesis that even for test cases with mostly common SQL features, still only a fraction of these features are supported by most DBMSs. Intuitively, if we generated only the most common features, we would expect most DBMSs to support them all; however, our experiment results show that most of the bug-inducing features are unsupported on more than half of the DBMSs. We considered only logic bugs since they returned incorrect results while executing without errors on the source DBMS. The H2 bug is excluded, since the developers fixed it by throwing an exception. We took bug-inducing test cases found on a specific DBMS and executed them on all target DBMSs. Bug-inducing statements that ran without error were marked as successful; otherwise, they were marked as unsuccessful, indicating a syntax or semantic error due to the feature not being supported.

*Result.* Figure 6 demonstrates the average validity rates when executing bug-inducing statements from each *bug source* DBMS on *target* DBMSs. We found that *none* of the bug-inducing tests can be executed successfully on *all* 18 DBMSs. Only 8% of the bug-inducing test cases can be executed successfully on more than 90% (17 of 18) DBMSs. The overall valid rate of execution of bug-inducing test cases in different DBMSs is 48%, indicating that features across DBMSs are mostly distinct and that generating these features is important for finding bugs. These results demonstrate that

most of the bugs cannot be executed successfully on other DBMSs, indicating that even though we generated mostly common features, they cannot be simply applied across all systems due to dialect deviations. Only 1 of 18 DBMSs can execute test cases successfully from more than half of the source DBMSs, showing that it is difficult to directly use existing testing generators and apply them to other DBMSs. The only one is SQLite which has a flexible type system, enabling most of the tests to be executed successfully.

## 5.3 Bug Finding Efficiency

We compared the efficiency in terms of bugs, unique query plan coverage and code coverage using multiple configurations of *SQLancer++* and *SQLancer*. First, we used the TLP oracle and evaluated *SQLancer++* with and without the validity feedback mechanism on CrateDB for 1 hour across five runs. We inspected the bugs found in one hour after prioritization. Second, we evaluated the line and branch coverage of *SQLancer++* under different mechanisms and *SQLancer* on three C/C++-based DBMSs, SQLite, PostgreSQL, and DuckDB. Code coverage is a common metric for fuzzing that assesses how much of a system might be tested and gives some indication of the features we covered. We observed that fuzzers for DBMSs can achieve higher code coverage than SQLancer and SQLancer++ even when executing only their seed corpus.

***Bugs detection.*** Table 5 shows the total number of bugs found within one hour. This result shows that *SQLancer++* could detect different unique logic bugs under different random seeds efficiently in a limited time. We observed that enabling the feedback mechanism resulted in more bugs being detected over the same period, indicating that the feedback mechanism can enhance bug-finding efficiency by increasing the validity rate of the generated queries.

***Code coverage.*** Table 3 shows the average percentage of code coverage across 10 runs in 24 hours. *SQLancer* achieves higher code coverage in all three DBMSs. This is because *SQLancer* generators are manually written for the specific DBMS and thus incorporate both features and heuristics specific to that DBMS, while *SQLancer++* supports only basic SQL statements and features; however, *SQLancer++* found 10 and 3 new bugs that were missed by *SQLancer* in DuckDB and SQLite respectively (see Table 2), demonstrating that the code coverage does not strongly relate to finding logic bugs. Typically, mutation-based coverage-guided fuzzers achieve higher coverage; for example, *SQLRight* outperformed *SQLancer* by 2% and 39% in line coverage for SQLite and PostgreSQL. We also found that *SQLancer* outperformed *SQLancer++* by 59% and 21% in line coverage for SQLite and PostgreSQL, likely due to developers putting more effort into making the generator more comprehensive; however, in DuckDB—where the generator might not be as mature—the gap was only 6%. This demonstrates that

**Table 3.** Coverage of *SQLancer++* and *SQLancer* on SQLite, PostgreSQL and DuckDB in 24 hours

| Approach | SQLite | | PostgreSQL | | DuckDB | |
|---|---|---|---|---|---|---|
| | Line | Branch | Line | Branch | Line | Branch |
| *SQLancer++* | 30.5% | 23.2% | 26.3% | 18.5% | 31.6% | 18.8% |
| *SQLancer++$_{Rand}$* | 30.0% | 22.8% | 26.1% | 18.5% | 31.4% | 18.5% |
| *SQLancer* | 46.6% | 36.4% | 32.3% | 23.3% | 33.4% | 20.9% |
| *SQLRight* | 47.5% | 37.7% | 44.9% | 34.2% | - | - |

**Table 4.** The validity rate of *SQLancer++* and *SQLancer* for all queries on SQLite, PostgreSQL, and DuckDB over 24 hours.

| Approach | SQLite | PostgreSQL | DuckDB |
|---|---|---|---|
| *SQLancer++* | 97.7% | 52.4% | 64.2% |
| *SQLancer++$_{Rand}$* | 24.9% | 21.6% | 24.6% |
| *SQLancer* | 98.0% | 25.1% | 35.5% |

our tool is especially beneficial for smaller developer teams, which, even if they adopt SQLancer, may be unable to incorporate many DBMS-specific features. For example, although MonetDB adopted SQLancer,[7] we still successfully uncovered 36 unique unknown bugs.

## 5.4 Feedback Mechanism

We further evaluated the impact of validity feedback by measuring the query validity rates.

***Methodology.*** We ran *SQLancer++* on PostgreSQL and SQLite and measured the validity rate of generated SQL statements using different feedback mechanisms for 24 hours across 10 runs. We calculated the validity rate by dividing the number of successfully executed test cases, each of which has two equivalent queries when using TLP oracle, by the total number of test cases. The validity rate converged in less than one minute due to the high throughput of these two DBMSs; therefore, we did not measure the efficiency of the feedback mechanism.

***Result.*** Table 4 shows the average validity rate across 10 runs after 24 hours when executing *SQLancer++* and *SQLancer* on SQLite, PostgreSQL, and DuckDB. In general, we observed a higher validity rate for SQLite due to its dynamic type system; conversely, PostgreSQL frequently raised errors for type mismatches and other errors and we thus observed a lower success rate. For SQLite, the validity rate was 97.7% when enabling feedback, increasing 292.5% compared with not enabling it. For PostgreSQL, the validity rate was 52.4%, increasing 121.9%. The rather lower validity rate of the *SQLancer* PostgreSQL generator might be due to its incorporation of additional DBMS-specific features, leading to a higher failure rate due to their complexity.

---

[7]https://github.com/MonetDB/sqlancer

**Table 5.** The average number of all, prioritized and unique bugs found in CrateDB across 5 runs in 1 hour

| Approach | Detected Bugs | Prioritized Bugs | Unique Bugs |
|---|---|---|---|
| *SQLancer++* | 67,878.2 | 35.8 | 11.4 |
| *SQLancer++$_{Rand}$* | 55,412.2 | 28.4 | 9.8 |

## 5.5 Bug Prioritization Effectiveness

To evaluate the effectiveness of our bug-prioritization approach, we determined how well it could identify bug-inducing test cases that triggered unique bugs.

***Methodology.*** We conducted a case study on an open-source DBMS, CrateDB, from which we could obtain the commit logs. This enabled us to distinguish bugs by reproducing them on different versions of the system to see if they have been fixed by distinct commits. We used the TLP oracle, and enabled the feedback mechanism. We evaluated *SQLancer++* on CrateDB for one hour of continuous running five times, recording any bugs that occurred. We ran the test for only one hour because CrateDB would run out of memory when we ran for multiple hours. Within the one-hour period, *SQLancer++* could still trigger numerous bugs. We analyzed the unique bugs identified by *SQLancer++* by bisecting them across different versions of the system.

***Result.*** Table 5 shows the results of bugs found by *SQLancer++* in one hour across five test runs. *SQLancer++* could detect more than 60K bug-inducing test cases. *SQLancer++* identified bugs that were potentially caused by the same reason and reported only 35.8 and 28.4 bugs, reducing more than 99% of the duplicated bugs. On average, among the bugs detected and prioritized with and without feedback mechanism, 11.4 and 9.8 bugs were unique respectively. More than half of the bugs were duplicated since bug-inducing test cases with different features may be due to the same reasons (*e.g.*, unequal operator "<>" and "!="). Despite this, the results show the effectiveness of prioritization, which could help developers prioritize the bugs.

## 6 Discussion

***Found bugs.*** As shown in Table 2, 140 of 196 bugs were logic bugs, which shows the effectiveness of our approach. Among the 56 other bugs, 2 of them were performance issues, 20 of them were due to unexpected errors (*e.g.*, connection failures), and 34 of them were crash bugs. We believe logic bugs are more severe and harder to detect than other kinds of bugs; crashes, and internal errors are immediately visible to users and developers, while logic bugs can go unnoticed. Furthermore, our current implementation includes mostly commonly used features, meaning that bugs could impact

users relying on these features. These reasons could explain the high fix rate for our reported bugs.

***Comparison with* SQLancer.** In addition to the bugs *SQLancer++* found in the previously untested DBMSs, it also found previously unknown bugs in the well-tested DBMSs, for example, MySQL and SQLite. Of the 196 bugs, 138 occurred in DBMSs not supported by *SQLancer*, while the remaining 58 occurred in DBMSs that *SQLancer* supported, but relied on features unsupported by its current generators. For example, *SQLancer* lacks the scalar function REPLACE, for which we found a related bug in SQLite (see Listing 2). Besides, *SQLancer* supports all bit operators except for the bitwise inversion (~) in TiDB, in which we discovered a bug.[8] In general, we expect *SQLancer* to find more bugs than *SQLancer++*, as *SQLancer* uses the same test oracles for finding logic bugs, but provides manually-written generators for each DBMS it supports. *SQLancer++* and *SQLancer* are thus complementary: *SQLancer++* scales to many DBMSs with minimal effort, whereas *SQLancer* can provide deeper testing for individual systems with manually written generators.

***Manual effort.*** Limited manual effort is still required to apply *SQLancer++* to new DBMSs. First, it is necessary to manually include the database driver and specify the connection string, indicating the port, user, and password. Second, for some DBMSs, after schema generation, *SQLancer++* must ensure that the data inserted can be read. SQL statements like COMMIT (implicitly specified by JDBC) or customized statements for distributed systems might need to be used, for example, the REFRESH statement in CrateDB. These are explicitly issued for DBMSs that require them. These adaptation efforts are minor, as we implemented support for 21 DBMSs with each 16 LOC per DBMS on average. For most DBMSs, only 4 LOC are required to specify the connection string. Finally, we also created Docker containers for each system under test, requiring additional effort. The automatic generation of fuzz drivers is an active research area [5, 18, 51], and we believe that some of their insights could also help the automation of this functionality in DBMS testing.

## 7 Related Work

***Testing platforms.*** Two other large-scale deployments of bug-finding platforms have inspired our work on *SQLancer++*, namely OSS-Fuzz [39] and syzkaller [14]. OSS-Fuzz supports common fuzzers such as libfuzz [31], AFL++ [8], Honggfuzz [12], and ClusterFuzz [13]. These fuzzers are grey-box fuzzers, which use code coverage feedback from the tested program to prioritize which inputs should be further mutated. Unlike *SQLancer++*, OSS-Fuzz can apply these fuzzers to any program for which a fuzz harness is written, but typically cannot find logic bugs, as the fuzzers rely on crashes or sanitizers [40, 41, 46], the latter

---

[8]https://github.com/pingcap/tidb/issues/53506

which detect undefined behavior in C/C++, as test oracles. Syzkaller is a coverage-guided kernel fuzzer, which finds bugs in operating systems by mutating sequences of system calls. As the above fuzzers, syzkaller aims to find crash bugs. In contrast to OSS-Fuzz and syzkaller, the key challenge in our context is that hundreds of DBMSs with different SQL dialects exist, which we aim to support. A previous study [27] investigated the challenges of applying fuzzing tools for DBMSs in a continuous integration process, while it focused on crash bugs, it highlighted related challenges such as those we tackled. We aim to find deeper kinds of bugs, such as logic bugs.

**DBMS test oracles.** *SQLancer++* can use any test oracle that is not based on DBMS-specific features. In this work, we have used TLP and NoREC. TLP [36] finds logic bugs by comparing an original query with an equivalent query derived from it by partitioning the query's results into three parts. NoREC [35] detects logic bugs by comparing the results of a query that is receptive to optimizations with an equivalent one that the DBMS is likely to fail to optimize. Various other test oracles could be adopted. Differential Query Execution (DQE) [44] uses the same predicate in different SQL statements, assuming that the predicate consistently evaluates to the same value for a given statement. Pinolo [15] generates pairs of queries that are in superset or subset relations and checks whether the expected relation is met. EET [22] transforms given queries and checks whether the transformed versions still produce the same results as the original query. However, several oracles might be difficult to support in *SQLancer++* due to dependencies on DBMS-specific features. For example, CERT [3] requires parsing query plans, which are typically DBMS-specific, to find unexpected discrepancies in the DBMSs' cardinality estimator. DQP [4] detects logic bugs by using DBMS-specific query hints, and Mozi [28] through DBMS-specific configurations, neither of which could be easily supported.

**Test-case generation for DBMSs.** Various approaches have been proposed to automatically generate test inputs for DBMSs. The most closely related work is Griffin [9], which was proposed as an alternative to grammar-based fuzzers by maintaining lightweight metadata related to the DBMS state to improve mutation correctness in fuzzing. While this enables it to fuzz a range of DBMSs, as a mutation-based fuzzer based on AFL++, it requires a diverse seed input corpus and aims to find only crash bugs, which limits the scope and scalability. BuzzBee [49] aims to fuzz various kinds of DBMSs, which include NoSQL systems. However, it shares a similar limitation with Griffin that it cannot be applied to find logic bugs. DynSQL [20] incorporates error feedback by the DBMSs to incrementally expand a given SQL query while maintaining its validity to find crash bugs in DBMSs. In contrast to DynSQL, our feedback applies to potentially unknown DBMSs and to entire SQL statements aiming to find

logic bugs. SQLRight [29], inspired by grey-box fuzzers, uses code-coverage feedback in combination with test oracles such as NoREC and TLP to find bugs. Query Plan Guidance (QPG) [2] uses query plans as a feedback mechanism to determine whether a given database state has saturated for finding potential bugs. Squirrel [54] is a mutation-based method to generate new queries for finding memory errors. None of the above approaches explored generators that could apply to the multitude of existing DBMSs aiming to find logic bugs.

## 8 Conclusion

In this paper, we have presented a simple and effective adaptive query generator, which is at the core of a new large-scale automated testing platform for DBMSs, called *SQLancer++*, which also includes new techniques to prioritize bugs as well as model the schema of the DBMS under test. As part of our initial efforts, *SQLancer++* has enabled developers of 18 DBMSs to fix 180 previously unknown bugs in their systems. The adaptive query generator is only a first step towards our vision of fully automated DBMS testing. Various challenges remain to be tackled, such as covering uncommon features in the generator, designing more sophisticated policies to utilize DBMS feedback, or prioritizing bug reports. We hope that *SQLancer++* will eventually become a standard tool in DBMS developers' toolboxes.

## References

[1] 2021. Jazzer: Coverage-guided, in-process fuzzing for the JVM. https://github.com/CodeIntelligenceTesting/jazzer.

[2] Jinsheng Ba and Manuel Rigger. 2023. Testing Database Engines via Query Plan Guidance. In *Proceedings of the 45th International Conference on Software Engineering* (Melbourne, Victoria, Australia) *(ICSE '23)*. IEEE Press, 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 *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (, Lisbon, Portugal,) *(ICSE '24)*. Association for Computing Machinery, New York, NY, USA, Article 133, 13 pages. doi:10.1145/3597503.3639076

[4] Jinsheng Ba and Manuel Rigger. 2024. Keep It Simple: Testing Databases via Differential Query Plans. *Proc. ACM Manag. Data* 2, 3, Article 188 (May 2024), 26 pages. doi:10.1145/3654991

[5] Domagoj Babić, Stefan Bucur, Yaohui Chen, Franjo Ivančić, Tim King, Markus Kusano, Caroline Lemieux, László Szekeres, and Wei Wang. 2019. FUDGE: fuzz driver generation at scale. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (Tallinn, Estonia) *(ESEC/FSE 2019)*. Association for Computing Machinery, New York, NY, USA, 975–985. doi:10.1145/3338906.3340456

[6] Yang Chen, Alex Groce, Chaoqiang Zhang, Weng-Keen Wong, Xiaoli Fern, Eric Eide, and John Regehr. 2013. Taming compiler fuzzers. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Seattle, Washington, USA) *(PLDI '13)*. Association for Computing Machinery, New York, NY, USA, 197–208. doi:10.1145/2491956.2462173

[7] The Business Research Company. 2024. *Database Software Global Market Report 2024*. Research and Markets. https://www.researchandmarkets.com/report/database-software

[8] Andrea Fioraldi, Dominik Maier, Heiko Eißfeldt, and Marc Heuse. 2020. {AFL++}: Combining incremental steps of fuzzing research. In *14th USENIX Workshop on Offensive Technologies (WOOT 20)*.

[9] Jingzhou Fu, Jie Liang, Zhiyong Wu, Mingzhe Wang, and Yu Jiang. 2023. Griffin: Grammar-Free DBMS Fuzzing. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering* (Rochester, MI, USA) *(ASE '22)*. Association for Computing Machinery, New York, NY, USA, Article 49, 12 pages. doi:10.1145/3551349.3560431

[10] Xiyue Gao, Zhuang Liu, Jiangtao Cui, Hui Li, Hui Zhang, Kewei Wei, and Kankan Zhao. 2023. A Comprehensive Survey on Database Management System Fuzzing: Techniques, Taxonomy and Experimental Comparison. arXiv:2311.06728 [cs.DB]

[11] Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. 1995. *Bayesian data analysis*. Chapman and Hall/CRC.

[12] Google. 2016. Honggfuzz. https://honggfuzz.dev/. Accessed: 2024-03-30.

[13] Google. 2018. ClusterFuzz - Scalable Fuzzing Infrastructure. https://google.github.io/clusterfuzz/. Accessed: 2024-03-30.

[14] Google. 2023. syzkaller: unsupervised, coverage-guided kernel fuzzer. https://github.com/google/syzkaller. Accessed: 2024-02-20.

[15] Zongyin Hao, Quanfeng Huang, Chengpeng Wang, Jianfeng Wang, Yushan Zhang, Rongxin Wu, and Charles Zhang. 2023. Pinolo: Detecting Logical Bugs in Database Management Systems with Approximate Query Synthesis. In *2023 USENIX Annual Technical Conference (USENIX ATC 23)*. 345–358.

[16] Benjamin Hilprecht, Carsten Binnig, and Uwe Röhm. 2020. Learning a Partitioning Advisor for Cloud Databases. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data* (Portland, OR, USA) *(SIGMOD '20)*. Association for Computing Machinery, New York, NY, USA, 143–157. doi:10.1145/3318464.3389704

[17] Gui Huang, Xuntao Cheng, Jianying Wang, Yujie Wang, Dengcheng He, Tieying Zhang, Feifei Li, Sheng Wang, Wei Cao, and Qiang Li. 2019. X-Engine: An Optimized Storage Engine for Large-scale E-commerce Transaction Processing. In *Proceedings of the 2019 International Conference on Management of Data* (Amsterdam, Netherlands) *(SIGMOD '19)*. Association for Computing Machinery, New York, NY, USA, 651–665. doi:10.1145/3299869.3314041

[18] Kyriakos K. Ispoglou, Daniel Austin, Vishwath Mohan, and Mathias Payer. 2020. FuzzGen: automatic fuzzer generation. In *Proceedings of the 29th USENIX Conference on Security Symposium (SEC'20)*. USENIX Association, USA, Article 128, 17 pages.

[19] Zhiyuan Jiang, Xiyue Jiang, Ahmad Hazimeh, Chaojing Tang, Chao Zhang, and Mathias Payer. 2021. Igor: Crash Deduplication Through Root-Cause Clustering. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security* (Virtual Event, Republic of Korea) *(CCS '21)*. Association for Computing Machinery, New

York, NY, USA, 3318–3336. doi:10.1145/3460120.3485364

[20] Zu-Ming Jiang, Jia-Ju Bai, and Zhendong Su. 2023. DynSQL: Stateful Fuzzing for Database Management Systems with Complex and Valid SQL Query Generation. In *32nd USENIX Security Symposium (USENIX Security 23)*. USENIX Association, Anaheim, CA, 4949–4965. https://www.usenix.org/conference/usenixsecurity23/presentation/jiang-zu-ming

[21] Zu-Ming Jiang, Si Liu, Manuel Rigger, and Zhendong Su. 2023. Detecting Transactional Bugs in Database Engines via Graph-Based Oracle Construction. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. USENIX Association, Boston, MA, 397–417. https://www.usenix.org/conference/osdi23/presentation/jiang

[22] 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 24)*. USENIX Association, Santa Clara, CA, 821–835. https://www.usenix.org/conference/osdi24/presentation/jiang

[23] Dario Korolija, Dimitrios Koutsoukos, Kimberly Keeton, Konstantin Taranov, Dejan S. Milojicic, and Gustavo Alonso. 2022. Farview: Disaggregated Memory with Operator Off-loading for Database Engines. In *12th Conference on Innovative Data Systems Research, CIDR 2022, Chaminade, CA, USA, January 9-12, 2022*. www.cidrdb.org. https://www.cidrdb.org/cidr2022/papers/p11-korolija.pdf

[24] Tim Kraska, Mohammad Alizadeh, Alex Beutel, Ed H. Chi, Ani Kristo, Guillaume Leclerc, Samuel Madden, Hongzi Mao, and Vikram Nathan. 2019. SageDB: A Learned Database System. In *9th Biennial Conference on Innovative Data Systems Research, CIDR 2019, Asilomar, CA, USA, January 13-16, 2019, Online Proceedings*. www.cidrdb.org. http://cidrdb.org/cidr2019/papers/p117-kraska-cidr19.pdf

[25] Vu Le, Mehrdad Afshari, and Zhendong Su. 2014. Compiler validation via equivalence modulo inputs. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Edinburgh, United Kingdom) *(PLDI '14)*. Association for Computing Machinery, New York, NY, USA, 216–226. doi:10.1145/2594291.2594334

[26] Bastien Lecoeur, Hasan Mohsin, and Alastair F. Donaldson. 2023. Program Reconditioning: Avoiding Undefined Behaviour When Finding and Reducing Compiler Bugs. *Proc. ACM Program. Lang.* 7, PLDI, Article 180 (jun 2023), 25 pages. doi:10.1145/3591294

[27] Jie Liang, Zhiyong Wu, Jingzhou Fu, Yiyuan Bai, Qiang Zhang, and Yu Jiang. 2024. WingFuzz: Implementing Continuous Fuzzing for DBMSs. In *2024 USENIX Annual Technical Conference (USENIX ATC 24)*. USENIX Association, Santa Clara, CA, 479–492. https://www.usenix.org/conference/atc24/presentation/liang

[28] Jie Liang, Zhiyong Wu, Jingzhou Fu, Mingzhe Wang, Chengnian Sun, and Yu Jiang. 2024. Mozi: Discovering DBMS Bugs via Configuration-Based Equivalent Transformation. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (Lisbon, Portugal) *(ICSE '24)*. Association for Computing Machinery, New York, NY, USA, Article 135, 12 pages. doi:10.1145/3597503.3639112

[29] Yu Liang, Song Liu, and Hong Hu. 2022. Detecting Logical Bugs of DBMS with Coverage-based Guidance. In *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, Boston, MA, 4309–4326. https://www.usenix.org/conference/usenixsecurity22/presentation/liang

[30] Shuang Liu, Chenglin Tian, Jun Sun, Ruifeng Wang, Wei Lu, Yongxin Zhao, Yinxing Xue, Junjie Wang, and Xiaoyong Du. 2024. Conformance Testing of Relational DBMS Against SQL Specifications. *arXiv preprint arXiv:2406.09469* (2024).

[31] LLVM Project. 2017. libFuzzer – a library for coverage-guided fuzz testing. https://llvm.org/docs/LibFuzzer.html. Accessed: 2024-03-30.

[32] Ryan Marcus, Parimarjan Negi, Hongzi Mao, Chi Zhang, Mohammad Alizadeh, Tim Kraska, Olga Papaemmanouil, and Nesime Tatbul. 2019. Neo: a learned query optimizer. *Proc. VLDB Endow.* 12, 11 (jul 2019),

1705–1718. doi:10.14778/3342263.3342644

[33] Rohan Padhye, Caroline Lemieux, Koushik Sen, Mike Papadakis, and Yves Le Traon. 2019. Semantic Fuzzing with Zest. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis* (Beijing, China) *(ISSTA 2019)*. Association for Computing Machinery, New York, NY, USA, 329–340. doi:10.1145/3293882.3330576

[34] PingCAP-Qe. [n. d.]. GitHub - PingCAP-QE/go-sqlancer: go-sqlancer. https://github.com/PingCAP-QE/go-sqlancer

[35] 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* (Virtual Event, USA) *(ESEC/FSE 2020)*. Association for Computing Machinery, New York, NY, USA, 1140–1152. doi:10.1145/3368089.3409710

[36] Manuel Rigger and Zhendong Su. 2020. Finding bugs in database systems via query partitioning. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 211 (nov 2020), 30 pages. doi:10.1145/3428279

[37] Manuel Rigger and Zhendong Su. 2020. Testing database engines via pivoted query synthesis. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 667–682.

[38] Andreas Seltenreich. 2022. Sqlsmith. https://github.com/anse1/sqlsmith

[39] Kostya Serebryany. 2017. OSS-Fuzz - Google's continuous fuzzing service for open source software. USENIX Association, Vancouver, BC.

[40] Konstantin Serebryany, Derek Bruening, Alexander Potapenko, and Dmitriy Vyukov. 2012. AddressSanitizer: A Fast Address Sanity Checker. In *2012 USENIX Annual Technical Conference (USENIX ATC 12)*. USENIX Association, Boston, MA, 309–318. https://www.usenix.org/conference/atc12/technical-sessions/presentation/serebryany

[41] Konstantin Serebryany and Timur Iskhodzhanov. 2009. ThreadSanitizer: data race detection in practice. In *Proceedings of the Workshop on Binary Instrumentation and Applications* (New York, New York, USA) *(WBIA '09)*. Association for Computing Machinery, New York, NY, USA, 62–71. doi:10.1145/1791194.1791203

[42] David Sidler, Zsolt Istvan, Muhsen Owaida, Kaan Kara, and Gustavo Alonso. 2017. doppioDB: A Hardware Accelerated Database. In *Proceedings of the 2017 ACM International Conference on Management of Data* (Chicago, Illinois, USA) *(SIGMOD '17)*. Association for Computing Machinery, New York, NY, USA, 1659–1662. doi:10.1145/3035918.3058746

[43] Donald R. Slutz. 1998. Massive Stochastic Testing of SQL. In *Proceedings of the 24rd International Conference on Very Large Data Bases (VLDB '98)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 618–622.

[44] 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 *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. 2072–2084. doi:10.1109/ICSE48619.2023.00175

[45] Jiansen Song, Wensheng Dou, Yu Gao, Ziyu Cui, Yingying Zheng, Dong Wang, Wei Wang, Jun Wei, and Tao Huang. 2024. Detecting Metadata-Related Logic Bugs in Database Systems via Raw Database Construction. *Proc. VLDB Endow.* 17, 8 (may 2024), 1884–1897. doi:10.14778/3659437.3659445

[46] Evgeniy Stepanov and Konstantin Serebryany. 2015. MemorySanitizer: Fast detector of uninitialized memory use in C++. In *2015 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 46–55. doi:10.1109/CGO.2015.7054186

[47] 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, Article 55 (may 2023), 26 pages. doi:10.1145/3588909

[48] Vitess. 2024. Vitess Fuzzing Summer 2023 Internship. https://vitess.io/blog/2024-04-08-vitess-fuzzing-summer-internship/. Accessed: 2024-10-09.

[49] Yupeng Yang, Yongheng Chen, Rui Zhong, Jizhou Chen, and Wenke Lee. 2024. Towards Generic Database Management System Fuzzing. In *33rd USENIX Security Symposium (USENIX Security 24)*. USENIX Association, Philadelphia, PA, 901–918. https://www.usenix.org/conference/usenixsecurity24/presentation/yang-yupeng

[50] Michal Zalewski. 2020. American Fuzzy Lop (AFL). https://github.com/google/AFL. Accessed: 2024-04-16.

[51] Cen Zhang, Mingqiang Bai, Yaowen Zheng, Yeting Li, Xiaofei Xie, Yuekang Li, Wei Ma, Limin Sun, and Yang Liu. 2023. Understanding Large Language Model Based Fuzz Driver Generation. arXiv:2307.12469 [cs.CR]

[52] Chi Zhang and Manuel Rigger. 2025. Constant Optimization Driven Database System Testing. *Proc. ACM Manag. Data* 3, 1, Article 24 (Feb. 2025), 24 pages. doi:10.1145/3709674

[53] Qirun Zhang, Chengnian Sun, and Zhendong Su. 2017. Skeletal Program Enumeration for Rigorous Compiler Testing. *SIGPLAN Not.* 52, 6 (jun 2017), 347–361. doi:10.1145/3140587.3062379

[54] 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 *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security* (Virtual Event, USA) *(CCS '20)*. Association for Computing Machinery, New York, NY, USA, 955–970. doi:10.1145/3372297.3417260

# A  Implementation

We illustrate the key implementation details of the adaptive statement generator.

## A.1  Supported Features

We summarize the SQL features supported by our generator. We considered basic and mostly standardized SQL features that we believe could successfully execute on different DBMSs. We summarize and explain the features in the following paragraphs; the exact set of features is detailed in the artifact. We have identified four common concrete granularities of *SQL features*, namely *statements*, *clauses*, expressions (*functions* and *operators*), *data types* as well as *abstract properties* (see Table 6).

We overall implemented only six common statements in the generator, including CREATE TABLE, CREATE INDEX, CREATE VIEW, INSERT, ANALYZE, and SELECT. Note that even such basic statements are not supported by all DBMSs (*e.g.*, CrateDB does not support CREATE INDEX).

Statements often have some optional *keywords* or *clauses* associated with them, such as UNIQUE and INNER JOIN in the above example. The keywords correspond to concrete strings in the statements, and the clauses usually contain various expressions. We assume the WHERE clause should be supported by every DBMS. However, some types of join clauses are not, for example, RIGHT JOIN was only supported in SQLite since 2022. We support six types of join.

Expressions consist of various *operators*, *functions*, as well as constants and column references, such as the comparison operators (=) and math function (SIN) shown in the table. Expressions are used in various contexts, such as in WHERE clauses or in ON for JOIN clauses. In total, we support 58 functions and 47 operators. Figure 7 shows the number
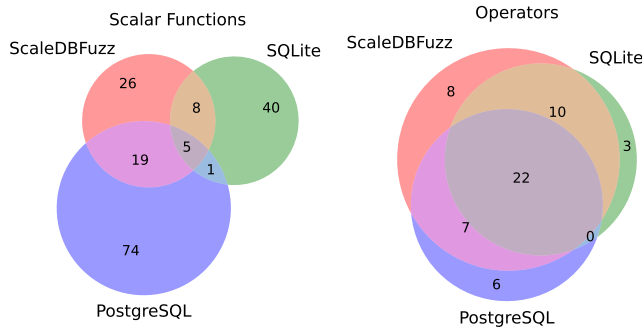
**Table 6.** SQL features

| Feature Type | | Number | Examples |
|---|---|---|---|
| Statement | | 6 | CREATE INDEX |
| Clause & Keyword | | 10 | RIGHT JOIN, SUBQUERY |
| Expression | Function | 58 | NULLIF, SIN |
| | Operator | 47 | +, =, AND, CASE−WHEN |
| Data type | | 3 | INTEGER |

**Listing 4.** Illustrative example of the adaptive statement generator. Design differences with non-adaptive generators are highlighted in red (non-adaptive) and green (adaptive).

```
 1  public SQLStatement createIndex(GeneralState
        state) {
 2    //...
 3    String stmt = "CREATE ";
 4  -   if (state.nextBoolean()) {
 5  +   if (state.shouldGenerate(Node.INDEX_UNIQUE)
 6  +     && state.nextBoolean()) {
 7      stmt += "UNIQUE ";
 8  +     state.generateFeature(Node.INDEX_UNIQUE);
 9    }
10    //...
11    return new SQLStatement(stmt, errors);
12  }
```



**Figure 7.** Venn diagram of SQL features shared by *SQLancer++*, and the SQLite and PostgreSQL *SQLancer* generators.

of distinct and shared features implemented in *SQLancer++* and the SQLite and PostgreSQL SQLancer generators. The leaf nodes of the expressions are constants and column references, which can have different *data type*s. *SQLancer++* supports generating three data types: integer, string, and boolean.

With respect to *abstract properties*, they mostly relate to the type system supported by the DBMS. Most importantly, we consider whether the DBMS' SQL dialect is statically typed or dynamically typed, and we represent these as features. For example, PostgreSQL is statically typed as it provides few implicit casts and rejects ill-typed statements, while SQLite is dynamically typed, as it can coerce most values

to the required data type at run time. This influences, for example, the generation of expressions. Consider a WHERE clause in a SELECT statement. For an untyped DBMS, the expression generator is free to generate an expression of any type. For a strictly typed system, it can only choose operators that produce a boolean expression; if this expression is, for example, a comparison operator, it must ensure that the compared operands have a compatible type. Rather than hard-coding the operand types of operators and argument types in functions, we also provide finer-grained features that learn the expected types. See Figure 5, where the identifier SIN1INT represents that the first argument of the function SIN is type integer. Further, if the generator generates SIN('a'), it will be captured as a feature SIN1STRING. Note that further properties are plausible; for example, we noticed that some features cannot co-occur, or some features require other features to be present. We currently do not capture such complex relationships.

### A.2 Feedback mechanism

Listing 4 demonstrates the simplified index generator of the adaptive generator, highlighting implementation differences compared to the non-adaptive one in *SQLancer*. We implemented the other generators similarly. When the generator selects new random alternatives, including statement keywords (line 7) and expression nodes, the feedback mechanisms indicate whether the option is supported, and record the feature selection for subsequent feedback updates. In addition to the interface shown in the example above, alternative interfaces allow querying the feedback generator for one of multiple supported features. In line 5, if the initial execution threshold has been reached, and the feature was found to not be supported by the DBMS under test, the call to shouldGenerate returns false. In line 8, we instruct the generator that we have selected the feature for generation. After execution of the SQL statement returned by line 11, the generator will use this information to update the validity rate. From the perspective of programmers extending *SQLancer++*'s generator, the internals are thus mostly transparent.

### A.3 Execution strategies

As with many other automated testing approaches, we have empirically determined suitable thresholds and strategies for our generator. First, the generator begins by generating expressions with low depth and gradually increases it. Specifically, the generator starts generating expressions at a depth of 1, and increases the depth after each $I$ executions by 1, up to a depth of 3. A lower depth corresponds to a reduced number of features, which enhances the learning efficiency of the generator since multiple features in one statement make it difficult to isolate the unsupported features. Besides, this enables it to identify simple bug-inducing features early and potentially deprioritize more complex ones later. Second,

the iterations *I* of executing test cases should be sufficient. The generator updates the probabilistic rules for each 100K test case, and we believe the number enables each feature to be sufficiently often executed. A larger *I* should be set if more features are incorporated or the probability threshold is more strict (*e.g.*, less than 1%).

## B  Artifact Appendix

### B.1  Abstract

This artifact contains the SQLancer++ implementation, Dockerized DBMS setups, and reproduction scripts. It supports reproducing the main results reported in the paper, including: bug-triggering test cases for each reported bug (Table 2), the SQL feature study (Figure 6), coverage and feedback effectiveness (Tables 3–4), and the bug prioritization evaluation (Table 5). Most experiments are executed inside Docker containers for isolation, although the experiment for Section 5.2 runs directly on the host machine.

### B.2  Artifact check-list (meta-information)

- **Program:** SQLancerPlusPlus/ (Java; Maven build)
- **Compilation:** Maven (mvn clean package -DskipTests)
- **Run-time environment:** Linux; Docker
- **Hardware:** x86_64; at least 32 cores, 64 GB RAM
- **Metrics:** Unique bugs found, code coverage, and validity rate.
- **Output:** logs/ and generated plots/tables
- **Experiments:** Reproduction for Sections 5.2–5.5
- **How much disk space required (approximately)?:** 50 GB
- **How much time is needed to prepare workflow (approximately)?:** About 1–2 hours
- **How much time is needed to complete experiments (approximately)?:** Around 1 day with parallel execution, or 2–3 days sequentially
- **Publicly available?:** Yes
- **Code licenses (if publicly available)?:** MIT
- **Archived (provide DOI)?:** https://doi.org/10.5281/zenodo.18289297

### B.3  Description

#### B.3.1  How to access.
All materials can be downloaded from https://doi.org/10.5281/zenodo.18289297.

#### B.3.2  Hardware dependencies.
An x86 CPU with at least 64 GB of RAM is required. Having more cores speeds up the experiments through increased parallelism. We used a 64-core AMD EPYC 7763 CPU at 2.45GHz and 512GB memory for our experiments.

#### B.3.3  Software Dependencies.
Our artifact requires Java 11 (or later), Python 3.10 (or later), Docker, and Maven.

#### B.3.4  Data sets.
The bug list is included in the artifact under the bugs directory.

### B.4  Installation

The installation steps and scripts assume a recent version of Ubuntu (*e.g.*, 22.04) and that the software dependencies (Java, Maven, Docker, and Python) above are already installed.

1. **Build SQLancer++:**
   ```
   cd SQLancerPlusPlus
   mvn clean package -DskipTests
   ```
2. **Build Docker images:**
   ```
   ./prepare_experiments_images.sh
   ```
3. **Prepare Python environment:**
   ```
   pip install matplotlib numpy pandas
   ```

### B.5  Experiment workflow

The experiments are launched by scripts on the host OS. These scripts start multiple Docker containers in isolation, save all collected results under the logs directory, and generate the tables and figures used in the paper. A typical workflow is:

1. Read README.md for an overview.
2. Build the Docker images.
3. Inspect the bug reports (for Section 5.1).
4. Analyze the bug-inducing test cases (for Section 5.2).
5. Run the experiments (for Sections 5.3–5.5) and inspect the results.

### B.6  Evaluation and expected results

The expected outputs include:

- **SQL feature study:** bugs/feature_heatmap.pdf
- **Coverage and validity rate:** the execution logs under logs/, and the generated files coverage.tex (for Table 3) and validity.tex (for Table 4)
- **Bug prioritization:** CrateDB logs and result files under logs/, used for Table 5

### B.7  Experiment customization

The provided scripts assume parallel execution with fixed time budgets as in the original paper, but they accept parameters that allow adjusting time budgets and switching to sequential execution.

### B.8  Methodology

Submission, reviewing, and badging methodology:

- https://www.acm.org/publications/policies/artifact-review-and-badging-current
- https://cTuning.org/ae