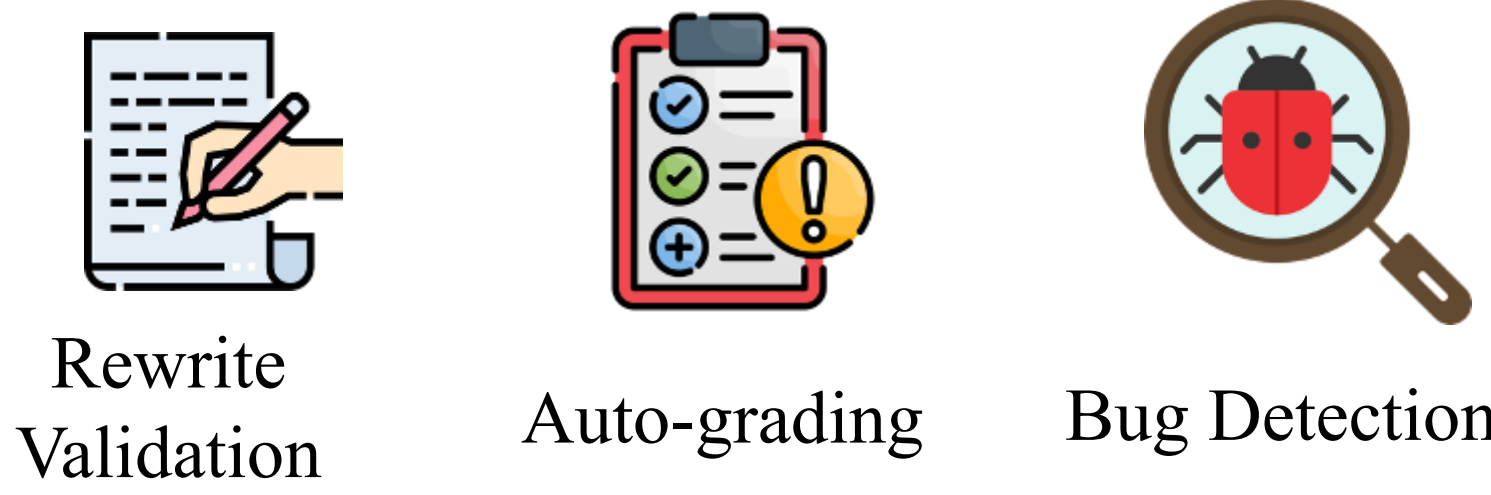


Motivation

Why is equivalence checking of SQL queries important?
 This task has a wide variety of application scenarios.

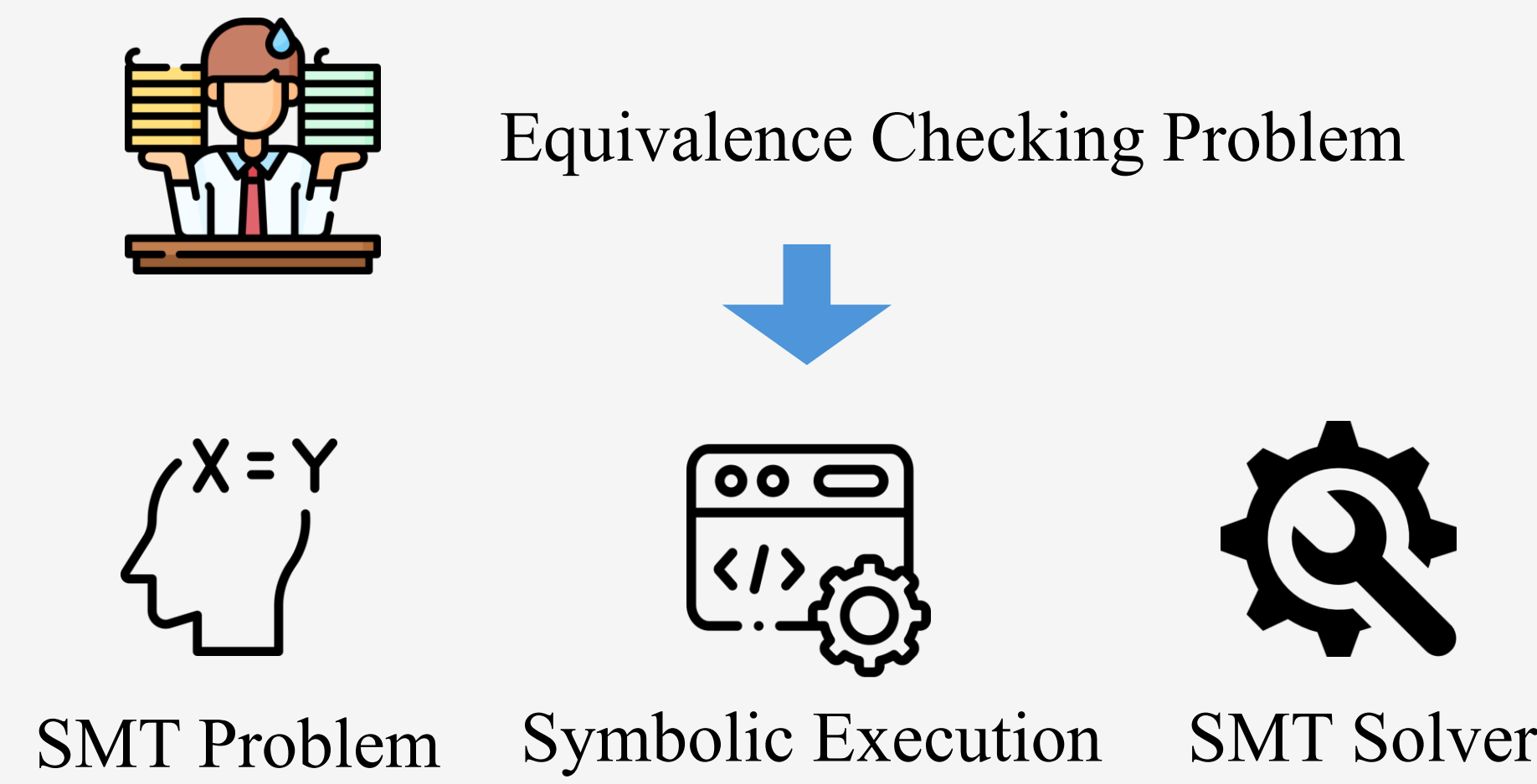


Problem Statement

Bounded equivalence verification: Given two SQL queries under a database schema, VeriEQL aims to verify whether these two queries always produce identical results on all possible input databases up to a bounded size that conform to the schema.

Technique

Key idea: VeriEQL reduces the equivalence checking problem into an SMT problem using symbolic reasoning and utilizes off-the-shelf constraint solvers to determine the satisfiability of SMT formulas.



Highlights

Expressive query language. VeriEQL supports SPJ, GROUP BY, aggregate functions, three-valued semantics, set/bag operations, conditional statements, etc.

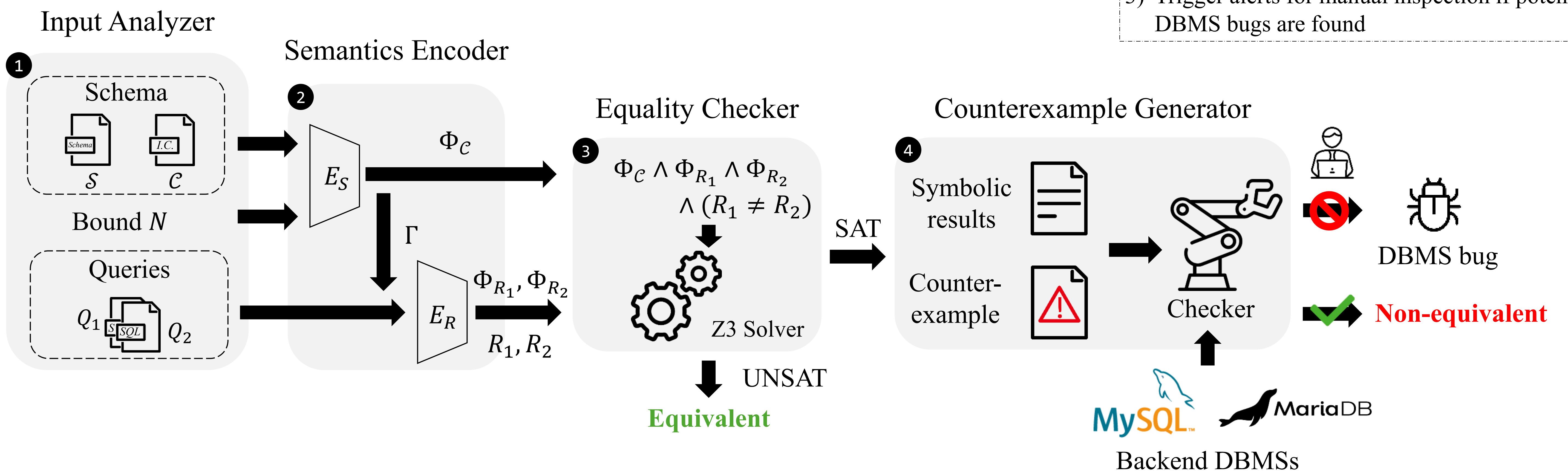
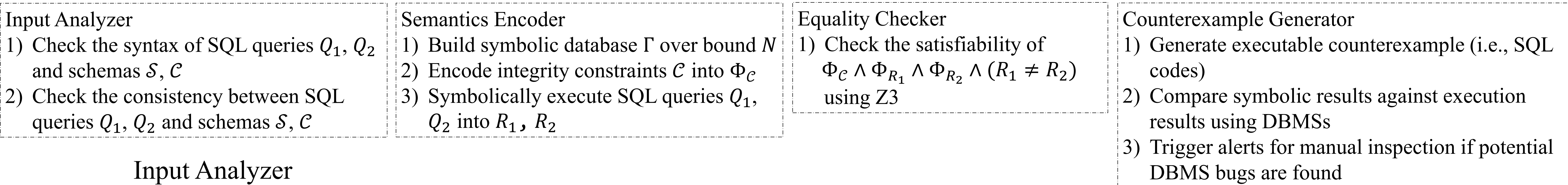
Genuine counterexample. VeriEQL refutes non-equivalent SQL queries with concrete database instances and SMT formulas.

Dialects. VeriEQL supports different SQL dialects, i.e., MySQL, MariaDB, Oracle, and PostgreSQL.

Good scalability. VeriEQL can check 70% of the 15,200 benchmarks with bound 5 in 5 minutes.

Small-Scope Hypothesis. 96% of non-equivalent benchmarks are refuted with less than 3 tuples.

System Overview



Schema \mathcal{S} :
 EMP: {id: int, age: int, ...}

Integrity constraints \mathcal{C} :
 PK: EMP.id

SQL queries:
 Q_1 : SELECT id FROM EMP WHERE age > 30
 Q_2 : SELECT id FROM EMP WHERE age >= 30

Symbolic table Γ ($N = 2$):

	EMP		
	id	age	...
t_1	$t_1.id$	$t_1.age$...
t_2	$t_2.id$	$t_2.age$...

Constraint formula Φ_C :
 $-2^{31} \leq t_1.id \leq 2^{31} - 1 \wedge -2^{31} \leq t_1.age \leq 2^{31} - 1$
 $\wedge -2^{31} \leq t_2.id \leq 2^{31} - 1 \wedge -2^{31} \leq t_2.age \leq 2^{31} - 1$
 $\wedge t_1.id \neq t_2.id \wedge t_1.id \neq \text{Null} \wedge t_2.id \neq \text{Null}$

Φ_{R_1} : $(t_1.age > 30 \rightarrow \neg \text{Del}(t'_1) \wedge t'_1.id = t_1.id) \wedge (t_1.age \leq 30 \rightarrow \text{Del}(t'_1))$
 $\wedge (t_2.age > 30 \rightarrow \neg \text{Del}(t'_2) \wedge t'_2.id = t_2.id) \wedge (t_2.age \leq 30 \rightarrow \text{Del}(t'_2))$
 Φ_{R_2} : $(t_3.age \geq 30 \rightarrow \neg \text{Del}(t'_3) \wedge t'_3.id = t_1.id) \wedge (t_3.age < 30 \rightarrow \text{Del}(t'_3))$
 $\wedge (t_4.age \geq 30 \rightarrow \neg \text{Del}(t'_4) \wedge t'_4.id = t_2.id) \wedge (t_4.age < 30 \rightarrow \text{Del}(t'_4))$
 where Del is an uninterpreted function.

Symbolic results

	R_1		R_2	
	id	id	id	id
t'_1	$t'_1.id$	t'_3	$t'_3.id$	
t'_2	$t'_2.id$	t'_4	$t'_4.id$	

Demonstration Scenarios

VeriEQL can validate optimized SQL queries.

Query	SQL
Q_1	SELECT DEPTNO, COUNT(*) FILTER (WHERE JOB = 'CLERK') FROM (SELECT * FROM EMP WHERE DEPTNO = 10 UNION ALL SELECT * FROM EMP WHERE DEPTNO > 20) AS t3 GROUP BY DEPTNO
Q_2	SELECT DEPTNO, COALESCE(SUM(EXPR\$1), 0) FROM (SELECT DEPTNO, COUNT(*) FILTER (WHERE JOB = 'CLERK') AS EXPR\$1 FROM EMP WHERE DEPTNO = 10 GROUP BY DEPTNO UNION ALL SELECT DEPTNO, COUNT(*) FILTER (WHERE JOB = 'CLERK') AS EXPR\$1 FROM EMP WHERE DEPTNO > 20 GROUP BY DEPTNO) AS t12 GROUP BY DEPTNO

Q_1 : the optimized query, Q_2 : the original query

The testPushCountFilterThroughUnion test case of Apache Calcite.

Time to check query equivalence on different input sizes for validating query optimizations.

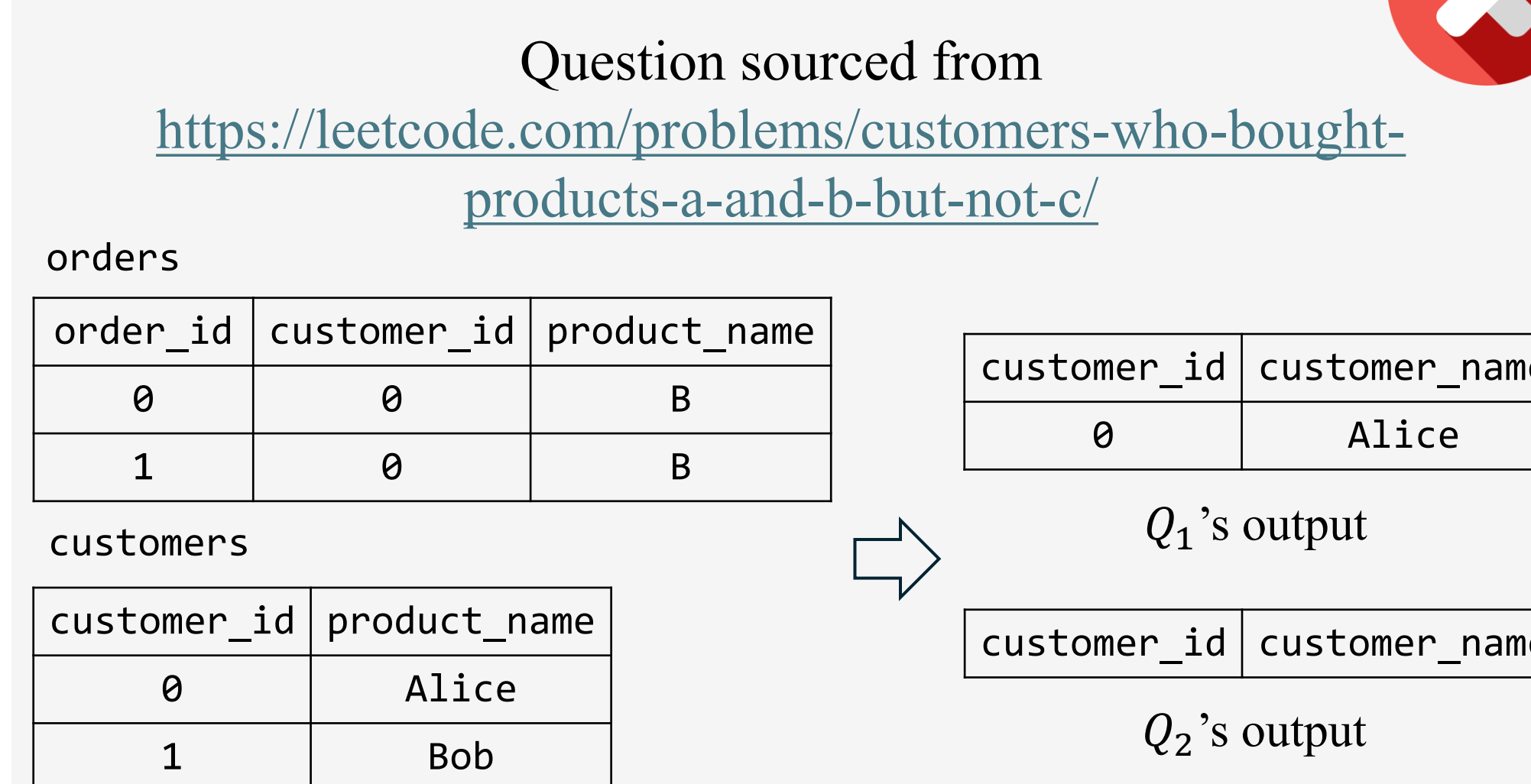
Size	1	2	3	4	5	6	7	8	9
Time (s)	0.2	0.4	0.6	1.0	2.4	6.6	19.7	98.5	118.2

Good scalability

VeriEQL can auto-grade queries from LeetCode.

Query	SQL
Q_1	WITH temp AS (SELECT DISTINCT A.customer_id, B.customer_id, B.customer_name SUM(CASE WHEN A.product_name IN ('A', 'B') THEN 1 ELSE 0 END) AS AB, SUM(CASE WHEN A.product_name = 'C' THEN 1 ELSE 0 END) AS C, FROM orders A JOIN customers B ON A.customer_id = B.customer_id GROUP BY A.customer_id) SELECT customer_id, customer_name FROM temp WHERE AB >= 2 AND C = 0
Q_2	SELECT customer_id, customer_name FROM customers WHERE customer_id IN (SELECT DISTINCT customer_id FROM orders WHERE product_name = 'A') AND customer_id IN (SELECT DISTINCT customer_id FROM orders WHERE product_name = 'B') AND customer_id NOT IN (SELECT DISTINCT customer_id FROM orders WHERE product_name = 'C') ORDER BY customer_id

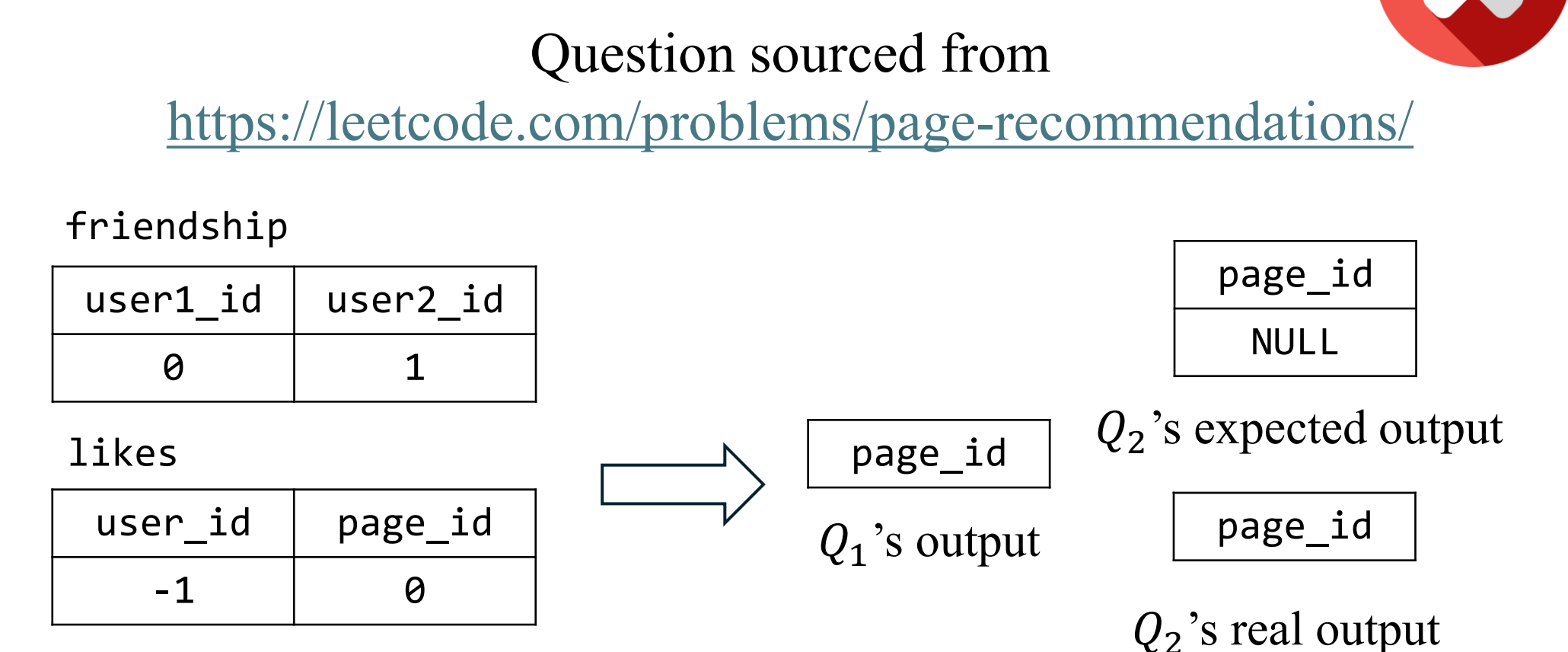
Q_1 : a user-provided answer, Q_2 : the ground-truth



VeriEQL can find bugs of DBSMs.

Query	SQL
Q_1	SELECT DISTINCT page_id AS recommended_page FROM (SELECT CASE WHEN user1_id = 1 THEN user2_id WHEN user2_id = 1 THEN user1_id ELSE NULL END AS user_id FROM friendship) AS tb1 JOIN likes AS tb2 ON tb1.user_id = tb2.user_id WHERE page_id NOT IN (SELECT page_id FROM likes WHERE user_id = 1)
Q_2	SELECT DISTINCT page_id AS recommended_page FROM (SELECT b.user_id, b.page_id FROM friendship a LEFT JOIN likes b ON (a.user2_id = b.user_id OR a.user1_id = b.user_id) AND (a.user1_id = 1 OR a.user2_id = 1) WHERE b.page_id NOT IN (SELECT DISTINCT page_id FROM likes WHERE user_id = 1)) T

Q_1 : a user-provided answer, Q_2 : the ground-truth



An implementation bug in MySQL v8.0.32.
<https://bugs.mysql.com/bug.php?id=110244>