

# ValScope: Value-Semantics-Aware Metamorphic Testing for Detecting Logical Bugs in DBMSs

## Abstract

Pending...

## 1 Introduction

### 2 SQL Query Approximation

A fundamental challenge in DBMS testing lies in the test oracle problem: given two semantically related SQL queries, it is often unclear whether their outputs should be identical or follow a predictable metamorphic relationship. Traditional testing frameworks typically rely on strict equivalence or simple set inclusion checks, which are insufficient to capture subtle semantic deviations introduced by optimizer transformations or operator mutations [1]. To address this limitation, we introduce a more comprehensive notion of SQL Query Approximation, which unifies two complementary perspectives of query behavior: the **set-semantic** and **value-semantic** dimensions. The set-semantic dimension characterizes differences in the returned tuple sets, while the value-semantic dimension captures monotonic variations in the computed or aggregated values over those tuples. Together, these two dimensions form a unified framework for expressing and reasoning about semantic consistency, enabling the construction of more expressive test oracles that can detect both structural and value-level inconsistencies in DBMS behavior.

#### 2.1 Set-Semantic Approximation

In this section, we first formalize the notion of approximation at the set level. This relation captures inclusion or containment among query result sets.

**Definition 2.1** (Set-Semantic Approximation Relation). Given a database  $D$ , let  $q_1$  and  $q_2$  be two SQL queries whose result sets are  $R(q_1, D)$  and  $R(q_2, D)$ , respectively. We say that  $q_1$  is the *set-based under-approximation* of  $q_2$  over  $D$ , denoted by  $q_1 \leq_D^s q_2$ , if and only if:

$$R(q_1, D) \subseteq R(q_2, D)$$

Conversely,  $q_1$  is the *set-based over-approximation* of  $q_2$  over  $D$ , denoted by  $q_1 \geq_D^s q_2$ , if and only if:

$$R(q_1, D) \supseteq R(q_2, D)$$

Here,  $R(q, D)$  represents the multi-set returned by evaluating query  $q$  on database  $D$ , and  $\subseteq$  and  $\supseteq$  denote inclusion and containment relations between two multi-sets.

Intuitively, the set-semantic approximation forms a partial order over queries:  $q_1 \leq_D^s q_2$  means that  $q_1$  produces a narrower or more restrictive result than  $q_2$ , while  $q_1 \geq_D^s q_2$  means that  $q_1$  yields a broader or less restrictive result. These

two relations are inverses of each other and together define the lattice of set-based approximations.

**Example 2.1.** Consider a database  $D = \{t_1\}$ , where  $t_1(c_1) = \{-1, 0, 1\}$ . Let the following queries be defined:

$$q_1 : \text{SELECT } c_1 \text{ FROM } t_1 \text{ WHERE } c_1 \leq 0$$

$$q_2 : \text{SELECT } c_1 \text{ FROM } t_1 \text{ WHERE TRUE}$$

$$q_3 : \text{SELECT } c_1 \text{ FROM } t_1 \text{ WHERE } c_1 < 0$$

We have:

$$R(q_3, D) = \{-1\}, \quad R(q_1, D) = \{-1, 0\}, \quad R(q_2, D) = \{-1, 0, 1\}$$

Hence,

$$R(q_3, D) \subseteq R(q_1, D) \subseteq R(q_2, D)$$

which gives the approximation chain:

$$q_3 \leq_D^s q_1 \leq_D^s q_2$$

Intuitively,  $q_3$  is a stricter version of  $q_1$ , and  $q_1$  a stricter version of  $q_2$ , each progressively expanding the selection condition and thus broadening the result set.

#### 2.2 Value-Semantic Approximation

To overcome the limitation of purely set-based inclusion, we extend the approximation relation from the result-set level to the value level. Unlike the set-based relation that focuses on tuple inclusion, the value-semantic relation captures the monotonic variation of target columns – the columns whose values are directly affected by functional or aggregation operations. This allows the framework to detect logical bugs where queries return identical tuples but diverge in their value semantics, such as incorrect computations in aggregation or updates.

**Definition 2.2** (Value-Semantic Approximation Relation). Given a database  $D$ , let  $q_1$  and  $q_2$  be two SQL queries whose result sets are  $R(q_1, D)$  and  $R(q_2, D)$ , respectively. Let  $C_t \subseteq \text{Cols}(R(q_1, D)) \cap \text{Cols}(R(q_2, D))$  denote the target columns whose values will be compared. Let  $G$  denote the grouping or ordering basis, determined as follows:

If the query contains a GROUP BY clause,  $G$  corresponds to the group-by keys. Otherwise,  $G$  represents a deterministic ordering over non-target columns (e.g., primary key or lexicographic ordering of attributes) to align tuples for comparison.

We say that  $q_1$  is the *value-based over-approximation* of  $q_2$  over  $D$ , denoted by  $q_1 \geq_D^v q_2$ , if and only if:

$$\forall g \in G^*, \forall c \in C_t, V_{q_1}(g, c) \geq V_{q_2}(g, c)$$

where  $G^*$  is the set of all comparable tuple groups under  $G$ , and  $V_q(g, c)$  denotes the value of column  $c$  in group  $g$  (or tuple position) produced by query  $q$ .

Conversely,  $q_1$  is the *value-based under-approximation* of  $q_2$ , denoted by  $q_1 \leq_D^v q_2$ , if and only if:

$$\forall g \in G^*, \forall c \in C_t, V_{q_1}(g, c) \leq V_{q_2}(g, c)$$

This definition unifies two cases: group-wise comparison for aggregation queries, and order-aligned comparison for non-aggregated results.

Intuitively, the *set-semantic approximation* ( $\leq_D^s$ ) describes inclusion of tuples, while the *value-semantic approximation* ( $\leq_D^v$ ) reflects monotonicity among the values of corresponding tuples. [lin: Check why can coexist?] In practice, the two forms often coexist:

$$q_1 \leq_D^s q_2 \wedge q_1 \leq_D^v q_2$$

which indicates that  $q_1$  returns a subset of  $q_2$ 's tuples and the corresponding values in the target columns are not larger.

**Example 2.2.** Consider a table  $t_1(c_1, c_2)$  as follows:

	$c_2$	$c_1$
	A	10
$t_1 =$	A	20
	B	5
	B	7

Let the following two queries be defined:

$q_1 : \text{SELECT } c_2, \text{MAX}(c_1) \text{ FROM } t_1 \text{ GROUP BY } c_2$

$q_2 : \text{SELECT } c_2, \text{MIN}(c_1) \text{ FROM } t_1 \text{ GROUP BY } c_2$

The results are:

$$R(q_1, D) = \begin{array}{c|c} c_2 & \text{MAX}(c_1) \\ \hline A & 20 \\ B & 7 \end{array} \quad R(q_2, D) = \begin{array}{c|c} c_2 & \text{MIN}(c_1) \\ \hline A & 10 \\ B & 5 \end{array}$$

Under the grouping basis  $G = \{c_2\}$  and target column  $C_t = \{c_1\}$ , we have for each  $g \in G^* = \{A, B\}$ :

$$V_{q_1}(g, c_1) \geq V_{q_2}(g, c_1)$$

Hence,  $q_1 \geq_D^v q_2$ . Intuitively, both queries return identical group sets (thus  $q_1 \equiv_D^s q_2$ ), but differ monotonically in their value semantics: the aggregated value of  $q_1$  in each group is no smaller than that of  $q_2$ .

### 2.3 Approximation Propagation

The approximation relations introduced in the previous section capture the semantic correspondence between two complete SQL queries by comparing their result sets or value outputs. However, in practical DBMS testing, a mutation usually affects only a local part of the query—for instance, a predicate, an operator, or an aggregation function—rather than the entire query. To understand how such a local change influences the final query result, we extend the discussion from the semantic level of full-query comparison to the structural level of SQL. Specifically, we define the concept of

*approximation propagation*, which describes how a local approximation relation established at one node of the query's abstract syntax tree (AST) can be transmitted through its parent operators and clauses, thereby determining how a single mutation impacts the overall approximation behavior of the query.

**Definition 2.3** (Approximation Propagation). Let  $D$  be a database, and let  $n_1, n_2$  denote two semantically comparable nodes (e.g., subqueries, predicates, or expressions) in the SQL AST. We use the unified notation  $n_1 \leq_D^\alpha n_2$  to represent an *approximation relation* of type  $\alpha \in \{s, v\}$ , where  $s$  and  $v$  correspond to the set-semantic and value-semantic levels, respectively. The relations defined in §2.1 and §2.2 describe query-level approximations between complete queries. In contrast, approximation propagation extends these relations to the structural level, capturing how local approximations between AST nodes can influence or induce approximations at higher layers of the query.

Formally, each operator  $op$  is characterized by two semantic properties: a *mapping* ( $\alpha_{in} \rightarrow \alpha_{out}$ ), which specifies how the operator transforms between set-level and value-level semantics, and a *direction*  $\sigma(op) \in \{+1, -1\}$ , where  $+1$  indicates that the operator preserves the approximation direction (monotone increasing) and  $-1$  indicates that it reverses the direction (monotone decreasing or negating). Based on these properties, the propagation of  $n_1 \leq_D^\alpha n_2$  can be classified into four canonical forms:

- (**Set → Set**): If a subquery or predicate  $p_1 \leq_D^s p_2$  is embedded under a higher-level set operator  $op_s$  (e.g., EXISTS, NOT EXISTS, logical NOT), then the resulting relation satisfies:

$$R(n_1, D) \leq_D^{s \cdot \sigma(op_s)} R(n_2, D)$$

where operators such as EXISTS are monotone increasing ( $\sigma = +1$ ), while NOT EXISTS or NOT are monotone decreasing ( $\sigma = -1$ ), reversing inclusion ( $\subseteq \leftrightarrow \supseteq$ ).

- (**Set → Value**): If an aggregation or mapping function  $f$  is applied to two relations that satisfy  $R(n_1, D) \leq_D^s R(n_2, D)$ , then the corresponding value-level results satisfy:

$$f(R(n_1, D)) \leq_D^{v \cdot \sigma(f)} f(R(n_2, D))$$

where  $\sigma(f) = +1$  for monotone-increasing functions (e.g., MAX, SUM, COUNT), and  $\sigma(f) = -1$  for monotone-decreasing ones (e.g., MIN).

- (**Value → Value**): If an expression or scalar operator  $op_v$  is transformed to another form with monotonic direction  $\sigma(op_v)$ , then the resulting value-level outputs satisfy:

$$V_{n_1}(g, c) \leq_D^{v \cdot \sigma(op_v)} V_{n_2}(g, c)$$

This covers arithmetic transformations (+,  $\times 2$  with  $c > 0$ ) and functional ones (MAX → MIN).

- (**Value → Set**): If a value expression  $V_n$  feeds into a predicate or filtering operator  $op_s$ , and  $V_{n_1} \leq_D^v V_{n_2}$ , then the

induced output relations satisfy:

$$R(n_1, D) \leq_D^{s \cdot \sigma(\text{op}_s)} R(n_2, D)$$

where  $\sigma(\text{op}_s) = +1$  for monotone-increasing predicates (e.g.,  $x > c$ , where larger values of  $x$  make the condition more likely to hold and thus expand the result set), and  $\sigma(\text{op}_s) = -1$  for monotone-decreasing ones (e.g.,  $x < c$  or NOT EXISTS, where larger values of  $x$  make the condition less likely to hold, causing the result set to shrink).

**Remark.** In this definition,  $n_1$  and  $n_2$  are not restricted to complete queries. They can represent corresponding subqueries, expressions, or predicates within a single query or across two query variants. The relation  $\leq_D^\alpha$  thus captures how a local semantic approximation propagates through SQL operators according to their monotonic behavior, bridging the value- and set-level semantics within the same unified framework.

Intuitively, the propagation mechanism provides the semantic bridge between *tuple-level inclusion* and *value-level monotonicity*. Set-based approximations can trigger value changes through monotone operators, while value-based changes can, in turn, alter the query result set when the affected values participate in predicates. This bidirectional propagation enables comprehensive reasoning over multi-layer SQL dependencies.

**Example 2.3** (Set → Value Propagation). Consider two queries over a table  $t_1(c_1, c_2)$ :

$q_1 : \text{SELECT MAX}(c_1) \text{ FROM } t_1 \text{ WHERE } c_2 < 100$

$q_2 : \text{SELECT MAX}(c_1) \text{ FROM } t_1 \text{ WHERE } c_2 < 200$

In the query structure, let  $n_1$  and  $n_2$  denote the WHERE clause nodes of  $q_1$  and  $q_2$ , respectively. The condition  $c_2 < 100$  in  $n_1$  is stricter than  $c_2 < 200$  in  $n_2$ , so the rows selected by  $n_1$  form a subset of those selected by  $n_2$ . The parent node of these filters is the aggregation operator MAX, which is monotone increasing: when more rows are included, the maximum value of  $c_1$  can only increase or remain the same. As a result, the difference at the set level (fewer or more tuples) propagates upward to a difference at the value level (smaller or larger aggregated value).

Intuitively,  $n_1$  and  $n_2$  illustrate how a local change in the filter condition at the set level can influence the aggregated result value, demonstrating the propagation from Set to Value.

**Example 2.4** (Value → Set Propagation). Consider two semantically related queries over a table  $t_1(c_1)$ :

$q_1 : \text{SELECT * FROM (SELECT MAX}(c_1) \text{ AS } x \text{ FROM } t_1) \text{ AS subq WHERE } x > 100$

$q_2 : \text{SELECT * FROM (SELECT MIN}(c_1) \text{ AS } x \text{ FROM } t_1) \text{ AS subq WHERE } x > 100$

The two queries differ only in the inner aggregation. Let  $n_1$  and  $n_2$  denote the aggregation nodes MAX( $c_1$ ) and MIN( $c_1$ ),

respectively. Changing MAX to MIN decreases the derived value  $x$ . Since the outer predicate  $x > 100$  is monotone increasing in  $x$ , smaller  $x$  values make the condition harder to satisfy, resulting in fewer output tuples. Consequently, the result of  $q_2$  becomes a subset of  $q_1$ , showing a typical Value → Set propagation.

These propagation behaviors connect the two approximation dimensions, allowing a single mutation at any AST node (e.g., MAX→MIN) to yield predictable, analyzable effects on both result structure and result values. The unified propagation model forms the semantic foundation of our testing framework.

### 3 Approach

#### 3.1 Overview

#### 3.2 Construction of the Original Query

#### 3.3 Approximate Mutators

#### 3.4 Approximation Propagation Analysis

In this section, we further develop an executable algorithm to determine how local semantic changes propagate through the SQL AST. The main purpose of this algorithm is to formalize the top-down reasoning process introduced in Definition 2.3 into a systematic, bottom-up propagation procedure that connects local node mutations with their global semantic consequences at the query level.

Algorithm 1 presents the overall propagation process. Initially (Line 1–3), the algorithm receives the mutated node information  $node\_info$  and the complete AST of the query. It first identifies the mutated node  $n_{mut}$  and constructs its ancestor chain from the mutation site to the query root, represented as list  $= [n_{mut}, \dots, n_{root}]$ . This structure enables the algorithm to traverse each parent operator sequentially and reason about how the mutation propagates upward. The algorithm then initializes the local semantic level of the mutation (Line 4–6): if the mutated node involves predicates or subqueries, it starts at the set level ( $\alpha(n_{mut}) = s$ ); otherwise, for expressions or aggregations, it starts at the value level ( $\alpha(n_{mut}) = v$ ). A direction accumulator  $sign$  is also initialized to +1 to indicate that propagation initially preserves directionality. In the propagation stage (Line 8–17), the algorithm iteratively traverses each parent node of  $n_{mut}$  in a bottom-up manner. For each parent operator  $op_i$  (Line 9), it consults Table 1 to retrieve the corresponding semantic mapping ( $\alpha_{in} \rightarrow \alpha_{out}$ ) and its monotonic direction  $\sigma(op_i) \in \{+1, -1\}$  (Line 10–12). The semantic level  $\alpha$  is updated according to the operator’s input-output mapping (e.g., Set→Value for aggregation or Value→Set for predicate filters), while the cumulative direction  $sign$  is updated multiplicatively ( $sign \leftarrow sign \cdot \sigma(op_i)$ ), preserving or reversing the relation based on operator polarity (Line 13–15). This recursive process effectively tracks the path of semantic transformation from the mutation site to the query output. Finally, at the root

---

**Algorithm 1** Approximation Propagation across SQL AST

---

**Input:** Mutated node info  $node\_info$ , SQL AST  $AST$   
**Output:** Final query-level approximation between original and mutated queries

1: // **Step 1: Initialization**  
2: Identify mutated node  $n_{mut}$ .  
3: Build ancestor chain list =  $[n_{mut}, \dots, n_{root}]$ .  
4: // **Step 2: Local relation at mutation site**  
5: Decide initial level  $\alpha(n_{mut}) \in \{s, v\}$  by node type.  
6: Set local relation  $\mathcal{R}(n_{mut}) \leftarrow (\leq^{\alpha(n_{mut})})$ .  
7: Set direction accumulator  $sign \leftarrow +1$ .  
8: // **Step 3: Bottom-up propagation**  
9: **for** each parent node  $n_i$  in  $list$  (from child to root) **do**  
10:   Determine operator type  $op_i$  at  $n_i$ .  
11:   Lookup rule in Table 1 for  $op_i$ .  
12:   Obtain mapping  $(\alpha_{in} \rightarrow \alpha_{out})$ .  
13:   Obtain monotonic direction  $\sigma(op_i) \in \{+1, -1\}$ .  
14:   Update level:  $\alpha(n_{i+1}) \leftarrow \alpha_{out}$ .  
15:   Update direction:  $sign \leftarrow sign \cdot \sigma(op_i)$ .  
16:   Propagate symbolically:  
      $\mathcal{R}(n_{i+1}) \leftarrow \mathcal{R}(n_i)$  with  $sign$  applied.  
17: **end for**  
18: // **Step 4: Materialize root-level relation**  
19: **if**  $\alpha(n_{root}) = s$   
20:   **return**  $R(q_1, D) \leq_D^s R(q_2, D)$  with sign +,  
     or  $R(q_1, D) \geq_D^s R(q_2, D)$  with sign -.  
21: **else**  
22:   **return**  $V_{q_1}(g, c) \leq_D^v V_{q_2}(g, c)$  with sign +,  
     or  $V_{q_1}(g, c) \geq_D^v V_{q_2}(g, c)$  with sign -.  
23: **end if**

---

node (Line 18–25), the algorithm derives the final query-level relation—returning either a set-level ( $R(q_1, D) \leq_D^s R(q_2, D)$ ) or value-level ( $V_{q_1}(g, c) \leq_D^v V_{q_2}(g, c)$ ) approximation, determined by the accumulated propagation direction.

### 3.5 Results Checking

## References

- [1] 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.

**Table 1. Operator Type and Semantic Propagation Rules**

Operator Type	Semantic Mapping	Monotonic $\sigma(op)$	Direction	Example and Interpretation
Aggregation	Set $\rightarrow$ Value	+1: MAX, SUM, COUNT -1: MIN		Input set expands $\Rightarrow$ aggregated value increases (MAX/SUM) or decreases (MIN).
Predicate / Filter	Value $\rightarrow$ Set	+1: $x > c, x \geq c$ -1: $x < c, x \leq c$		Increasing value makes predicate easier ( $x > c$ ) or harder ( $x < c$ ) to satisfy, thus expanding or shrinking result set.
Logical Operator	Set $\rightarrow$ Set	+1: EXISTS -1: NOT, NOT EXISTS		Negation flips inclusion direction (EXISTS $\leftrightarrow$ NOT EXISTS).
Arithmetic / Expression	Value $\rightarrow$ Value	Depends on sign		$x + 1$ (+1) preserves direction; $-x$ (-1) reverses monotonicity.
Join / Projection / Subquery	Set $\rightarrow$ Set	+1		Structural operators preserve tuple inclusion without reversing semantic direction.