# Scaling Up Symbolic Reasoning for Relational Queries

or, speeding up debugging & verification of database queries

**Chenglong Wang**, Alvin Cheung, Ras Bodik University of Washington

> PAUL G. ALLEN SCHOOL of computer science & engineering

The language between human and relational databases (tables)





Select val From T Where color = red;

The language between human and relational databases (tables)







_		
Г		

Select val From T Where color = red;

 $T_1$  Join  $T_2$  On  $T_1$ .color= $T_2$ .color;

The language between human and relational databases (tables)







_		
Г		

Select val From T Where color = red;

 $T_1$  Join  $T_2$  On  $T_1$ .color= $T_2$ .color;

**Select** color, Sum(val) From T Group by color;

# "The Count Bug"

#### On Optimizing an SQL-like Nested Query

WON KIM IBM Research

SQL is a high-level nonprocedural data language which has received wide recognition in relational databases. One of the most interesting features of SQL is the nesting of query blocks to an arbitrary depth. An SQL-like query nested to an arbitrary depth is shown to be composed of five basic types of nesting. Four of them have not been well understood and more work needs to be done to improve their execution efficiency. Algorithms are developed that transform queries involving these basic types of nesting into semantically equivalent queries that are amenable to efficient processing by existing query-processing subsystems. These algorithms are then combined into a coherent strategy for processing a general nested query of arbitrary complexity.

Categories and Subject Descriptors: H.2.4 [Information Systems]: Systems-query processing

### 1982 "On Optimizing an SQL-like Nested Query" (Kim Won)



# "The Count Bug"

#### On Optimizing an SQL-like Nested Query

WON KIM IBM Research

SQL is a high-level nonprocedural data language which has received wide recognition in relational databases. One of the most interesting features of SQL is the nesting of query blocks to an arbitrary depth. An SQL-like query nested to an arbitrary depth is shown to be composed of five basic types of nesting. Four of them have not been well understood and more work needs to be done to improve their execution efficiency. Algorithms are developed that transform queries involving these basic types of nesting into semantically equivalent queries that are amenable to efficient processing by existing query-processing subsystems. These algorithms are then combined into a coherent strategy for processing a general nested query of arbitrary complexity.

Categories and Subject Descriptors: H.2.4 [Information Systems]: Systems-query processing

### 1982 "On Optimizing an SQL-like Nested Query" (Kim Won)



1987 "Optimization of Nested SQL Queries Revisited" (Ganski & Wong)

#### Rewrite rules for nested queries



PARTS: SUPPLY: PNUM QOH PNUM 3 6 3 10 1 3 8 0 10	QUAN SHIPDATE 4 7-3-79 2 10-1-78 1 6-8-78	qı	PARTS PNUM 10 8
8 0 10 10 8	1 6-8-78 2 8-10-81 5 5-7-83 [KIE 84:2]	<b>q</b> <sub>2</sub>	PARTS PNUM 10

# **Reasoning Tasks**



Verification

"Are two queries equivalent on ALL inputs"



Property Checking (for optimization) "Can the query return empty output on SOME input?"



Mutation testing / Grading "Find a distinguishing input between queries."





#### q(I) = empty





- complex analytical functions
- plays important roles in industry

#### can't afford 5 years to find a bug!

### Automatic Reasoning

Check whether  $q_1$  is equivalent to  $q_2$  (on ALL inputs)





## Automatic Reasoning

Check whether  $q_1$  is equivalent to  $q_2$  (on ALL inputs)



(unbounded) equivalence is undecidable





Check whether  $q_1$  is equivalent to  $q_2$  (on ALL inputs within a search space)





"Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent on ALL tables with at most k tuples"

(1) Target queries

 $q_1$ : Select id, val $q_2$ : Select id, valFrom TFrom TWhere id > IWhere id  $\neq I$ 

(2) Search space

tables with at most k rows

(3) Property  $q_1(T) \neq q_2(T)$ 



Target queries (1)

q<sub>1</sub>: Select id, val q<sub>2</sub>: Select id, val From T From T Where id > I Where  $id \neq I$ 

(2) Search space

tables with at most k rows

Property  $\left( 3\right)$  $q_1(T) \neq q_2(T)$ 

#### "Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent on ALL tables with at most k tuples"



"Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent on ALL inputs within size k"



Grouping & aggregation

" Select f(val) From T Group By id "

"Check whether  $q_1$ ,  $q_2$  are equivalent on ALL inputs within size k"



Grouping & aggregation

"Select f(val) From T Group By id "

Exponential ways to partition the table





"Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent on ALL inputs within size k"



Grouping & aggregation

" Select f(val) From T Group By id "

Exponential ways to partition the table



 $\begin{array}{l} Computationally\\ expensive\\ T_{out1} \subset T_{out2} \& T_{out2} \subset T_{out1} \end{array}$ 



id

id

2

id

2

k

**y**2

**y**k

• • •

"Check whether  $q_1$ ,  $q_2$  are equivalent on ALL inputs within size k"



Grouping & aggregation

"Select f(val) From T Group By id "

Exponential ways to partition the table

Computationally expensive  $T_{out1} \subset T_{out2} \& T_{out2} \subset T_{out1}$ 





"Check whether  $q_1$ ,  $q_2$  are equivalent on ALL inputs within size k"



Grouping & aggregation

" Select f(val) From T Group By id "

#### "Small Model"

A smaller search space to achieve same reasoning guarantee

Exponential ways to partition the table

Computationally expensive  $T_{out1} \subset T_{out2} \& T_{out2} \subset T_{out1}$ 



## Space Refinement



 $T_{out1} \neq T_{out2}$ (property)



S (search space)

#### "Small Model"

(1) If exists  $T \in S$  satisfying the property, we can find one in the S' too. (2) If none of tables in S' satisfying the property, then no T exists in S too.

"Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent on ALL inputs within size k"



## Space Refinement



(1) If exists  $T \in S$  satisfying the property, we can find one in the S' too. (2) If none of tables in S' satisfying the property, then no T exists in S too.

"Check whether  $q_1$ ,  $q_2$  are equivalent on ALL inputs within size k"

# Insight from Property

Many properties requires only one tuple in the output to invalidate.

"Check whether q<sub>1</sub>, q<sub>2</sub> are equivalent"  $q_1(T) \neq q_2(T)$ 

> Exists a row r with different *multiplicities in T<sub>out</sub> and T<sub>out</sub><sup>2</sup>*



 $r \in T_{out1}, r \notin T_{out2} \rightarrow q_1(T) \neq q_2(T)$ 

# Insight from the Property

*T* from search space *S* 



 $r \in T_{out}, r \notin T_{out}^2$ 

Many important properties requires only one tuple in the output to be invalidated.



 $r \in T_{out}, r \notin T_{out}$ 

T' can also distinguish q1 from q2!

# Assume r=(a, b) is the output tuple showing the difference between two queries $r \in T_{out1}, r \notin T_{out2}$

- q<sub>I</sub>: Select id, max(val) From T Group By id
- q<sub>2</sub>: Select id, min(val) From T Group By id



id	val
	• • •

id	val
а	Ь
	•••

# Assume r=(a, b) is the output tuple showing the difference between two queries $r \in T_{out1}, r \notin T_{out2}$

- q<sub>1</sub>: Select id, max(val) From T Group By id
- q<sub>2</sub>: Select id, min(val) From T Group By id



id	val
	•••

id	val
а	Ь
	• • •

# Assume r=(a, b) is the output tuple showing the difference between two queries $r \in T_{out1}, r \notin T_{out2}$

- q<sub>1</sub>: Select id, max(val) From T Group By id
- q<sub>2</sub>: Select id, min(val) From T Group By id



#### Assume r=(a, b) is the output tuple showing the difference between two queries $r \in T_{out}, r \notin T_{out}^2$

- q<sub>1</sub>: Select id, max(val) From T Group By id
- q<sub>2</sub>: Select id, min(val) From T Group By id





### Space Refinement



q<sub>1</sub>: Select id, max(val) From T Group By id

q<sub>2</sub>: Select id, min(val) From T Group By id





**q**2

tables with at most k rows

 $S' = \{T \in S \mid T \text{ contain only one group}\}$ 

id	val
0	2
Ι	3

id	val
0	1
Ι	3







q<sub>1</sub>: Select id, max(val) From T Group By id

q<sub>2</sub>: Select id, min(val) From T Group By id

## Symbolic Provenance Analysis

$$\begin{split} \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi_1) \land \dots \\ \hline \hline (q \sim \phi) \rightsquigarrow (T \sim \phi) \quad \text{(Table)} \quad & \frac{q = \text{Select}(q_1, f) \quad \text{schema}(q_1) = \bar{c}}{(q \sim \phi) \sim (q_1 \sim \phi \land [c_i \mapsto t.i]f)} \text{ (Select)} \\ \hline \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi) \quad \text{(Distinct)} \quad & \frac{q = \text{Proj}(\bar{v}, q_1) \quad \text{schema}(q_1) = \bar{c}}{(q \sim \phi) \rightsquigarrow (q_1 \sim \phi) \land (q_1 \sim \phi)} \\ \hline \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi) \quad \text{(Distinct)} \quad & \frac{q = \text{Proj}(\bar{v}, q_1) \quad \text{schema}(q_1) = \bar{c}}{(q \sim \phi) \rightsquigarrow (q_1 \sim (t_i \mapsto ([c_j \mapsto t.j]v_i)]\phi)} \\ \hline \hline q = \text{Join}(q_1, q_2) \quad \phi = \phi_1 \land \phi_2 \land \phi_3 \quad |\text{schema}(q_1)| = n_1 \\ \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi_1) \land (q_2 \sim [t.i \mapsto t.(i - n_1)]\phi_2) \quad \text{(Join)} \\ \hline \hline q = \text{LeftJoin}(q_1, q_2, f) \quad \phi = \phi_1 \land \phi_2 \land \phi_3 \quad \text{colRef}(\phi_i) \cap \text{schema}(q_i) = \emptyset_{(i=1,2)} \\ \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi) \land (q_2 \sim \phi) \quad \text{(Union)} \quad & \frac{q = \text{Rename}(q_1, name, \bar{c})}{(q \sim \phi) \rightsquigarrow (q_1 \sim \phi_1) \land (q_2 \sim \text{true})} \\ \hline \hline q = \text{Aggr}(\bar{c}, \bar{a}, \bar{c}_i, q_1) \\ \hline \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi_1) \land (q_2 \cap q_1) \land (q_2 \sim \text{true}) \quad \text{(LeftJoin)} \\ \hline \hline \phi_i = \phi_1 \land \phi_2 \quad \text{colRef}(\phi_1) \subseteq \{\bar{c}\} \\ \hline (q \sim \phi) \rightsquigarrow (q_1 \sim \phi_1) \land (q_2 \cap \phi_1) \rightsquigarrow (h_i (q_i \sim \phi_i) \rightsquigarrow (h_i (q_i'_k \sim \phi_{ik}')) \\ \hline \hline (h_i (q_i \sim \phi_i) \rightsquigarrow (h_i (q_i'_k \sim \phi_{ik}')) \quad \text{(Step)} \end{aligned}$$

Inductively define the analysis rules for different operators

How to combine provenance from *multiple queries* 

 $S \rightarrow S'$ 

If exists  $T \in S$  satisfying the property, we can find one in the S' too.

Analysis complexity: linear to the query size.







### Experiment

#### **Bounded Verification**

"Verify two queries are equivalent on ALL inputs."  $q_1(T) \equiv q_2(T)$ 

Benchmarks: 46 rules from Apache Calcite

#### **Test generation**

"Can the query return empty output?"

q(T) = empty

"Find a distinguishing input between queries."

 $q_1(T) \neq q_2(T)$ 

Benchmarks: 15 student submissions & prior work

### Experiment

#### **Bounded Verification**

"Verify two queries are equivalent on ALL inputs."  $q_1(T) \equiv q_2(T)$ 

Benchmarks: 46 rules from Apache Calcite

#### **Test generation**

"Can the query return empty output?"

q(T) = empty

"Find a distinguishing input between queries."

 $q_1(T) \neq q_2(T)$ 

Benchmarks: 15 student submissions & prior work

#### Process

Measure solving speed with and without space refinement

- 1. Increase search space size until hitting 10 minutes limit without refinement
- 2. Re-run the same search space with space refinement

### Experiment

#### **Bounded Verification**

"Verify two queries are equivalent on ALL inputs."  $q_1(T) \equiv q_2(T)$ 

Benchmarks: 46 rules from Apache Calcite

#### **Test generation**

"Can the query return empty output?"

q(T) = empty

"Find a distinguishing input between queries."

 $q_1(T) \neq q_2(T)$ 

Benchmarks: 15 student submissions & prior work

**Process** 

Measure solving speed with and without space refinement

- Increase search space size until hitting 1. 10 minutes limit without refinement
- 2. Re-run the same search space with space refinement

Cosette SQL Solver **Qex SQL Solver** 

#### **Bounded Verification**

"Verify that two queries are equivalent on ALL inputs with no more than k tuples."

#### Result



**Cosette** with and without refinement

### **Experiment — Verification**



**Qex** with and without refinement

#### **Bounded Verification**

"Verify that two queries are equivalent on ALL inputs with no more than k tuples."



**Cosette** with and without refinement

## **Experiment — Verification**



**Qex** with and without refinement

#### **Bounded Verification**

"Verify that two queries are equivalent on ALL inputs with no more than k tuples."

#### Result



## **Experiment — Verification**



## **Experiment - Test Generation**

#### **Test generation**

"Can the query return empty output on SOME input?"

"Find a distinguishing input between queries."



cosette with and without refinement



gex with and without refinement

## **Experiment - Test Generation**

#### **Test generation**

"Can the query return empty output on SOME input?"

"Find a distinguishing input between queries."

Result



case\_id

Property supported are those can be invalidate by one tuple in the output.

 $q_1(T) \neq q_2(T)$ " $q_1(T)$  contains exactly 5 tuples"  $q_I(T) \neq empty$ exists r,  $\phi(q(T))$ 

Improving provenance analysis precision.  $\bullet$ 

### Limitations

"every tuple in q(T) has same multiplicity"

Generalize to arbitrary property



### Summary

Scaling Up Symbolic Reasoning for Relational Queries

#### (3) Result



(1) smaller search space & easier to traverse (2) equivalent for reasoning

Low analysis overhead & over **100x** speed up in

bounded verification (1) test generation (2)



# Hidden Slides!

# Symbolic Provenance Analysis



Select id, min(val) From T Where val > 0Group By id



Multiple provenance exists for tout for a query q

(strong)

Choice of abstraction trades between the analysis overhead and pruning power