

Scalable Approximate Query Processing With The DBO Engine

C. Jermaine et al.
Presented by: Yi Zhang
Jan 29, 2008

1

Problem

- Need for online estimation of query answer
 - Updates as query processing progresses
 - Statistically meaningful bounds
 - Accurate final answer of course
 - Application:
 - Interactive data exploration over large datasets
- Query style
 - SELECT FROM WHERE GROUP BY

2

Nested Loop Join

- Long delay between successive updates
- Limited contribution to shortening confidence interval

3

Ripple Join

4

Estimation & Statistics

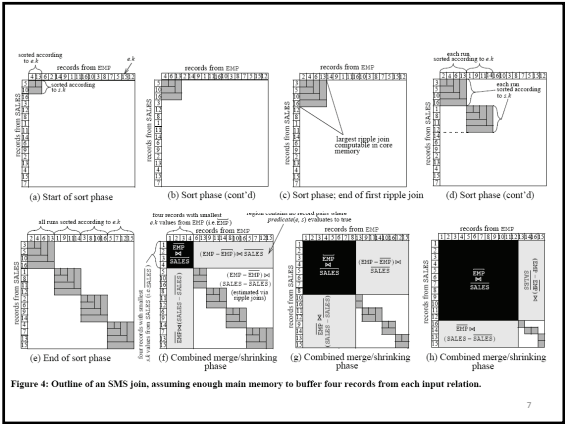
- An example
 - SELECT SUM(expression) FROM R,S WHERE p
 - Estimated by $\frac{|R| \cdot |S|}{|R| \cdot |S|} \sum_{r \in R, s \in S} \text{expression}_p(r,s)$
- Estimator unbiased, consistent
- Capable of giving tight confidence intervals
- Variance can also be calculated
- But...
 - When samples go out of memory...
 - Significant disk I/Os

5

SMS Join

- Generalization of ripple join
 - Identical to ripple join when all fits in memory
- Continues to provide confidence bounds when operating on disk
- Phases
 - Sort: read input relations in parallel, hash ripple join, and sort on join attributes
 - Interleave reads and writes across runs
 - Merge: similar to SMJ, parallel with Shrink
 - Shrink: update estimator after removing merged data space

6



Scalability!

- Two phase model of sort-merge does not scale well

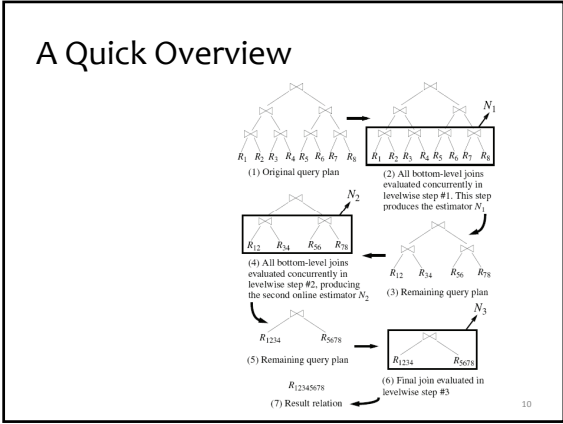
```

SELECT SUM(T.C)
FROM R,S,T
WHERE R.A=S.A AND S.B=T.B
  
```

- Fixes
 - Two separate joins: Can we do better?
 - Pipelining
 - Bad randomness
 - Hard to estimate e.g. intermediate relation size

How DBO Does It

- We need both scalability and progressive accurate estimation
- Design differences from traditional databases
 - Information sharing across relational operations
 - Processing abstracted into levelwise steps
 - Each levelwise step maintains and updates an online estimator
 - Results used as input for next levelwise step
 - Levelwise estimators combined with all lower level estimators to produce a single estimator for final answer



Levelwise Step

- Scan & Merge
- Scan phase
 - Concurrent for input relations in a levelwise step
 - Records of a relation divided into runs
 - In-memory join for unbiased guess
 - Randomized sort order
 - Sort on H(id, R.key)
 - Round-robin processing of runs

Example

- ...WHERE R1.B=R2.C AND R2.E=R3.F AND R3.G=R4.H

Result table discovered on-iteratively

R1	R2	R3	R4
A B	C D E	F G	H I
10 20	10 10 10	10 10	10 10 10
20 20	20 20 20	20 20	20 20 20
30 30	30 30 30	30 30	30 30 30
40 40	40 40 40	40 40	40 40 40
50 50	50 50 50	50 50	50 50 50
60 60	60 60 60	60 60	60 60 60
70 70	70 70 70	70 70	70 70 70
80 80	80 80 80	80 80	80 80 80
90 90	90 90 90	90 90	90 90 90

Sorted on H(A,B)

R1	R2	R3	R4
A B	C D E	F G	H I
10 20	10 10 10	10 10	10 10 10
20 20	20 20 20	20 20	20 20 20
30 30	30 30 30	30 30	30 30 30
40 40	40 40 40	40 40	40 40 40
50 50	50 50 50	50 50	50 50 50
60 60	60 60 60	60 60	60 60 60
70 70	70 70 70	70 70	70 70 70
80 80	80 80 80	80 80	80 80 80
90 90	90 90 90	90 90	90 90 90

Run 1

R1	R2	R3	R4
A B	C D E	F G	H I
10 20	10 10 10	10 10	10 10 10
20 20	20 20 20	20 20	20 20 20
30 30	30 30 30	30 30	30 30 30
40 40	40 40 40	40 40	40 40 40
50 50	50 50 50	50 50	50 50 50
60 60	60 60 60	60 60	60 60 60
70 70	70 70 70	70 70	70 70 70
80 80	80 80 80	80 80	80 80 80
90 90	90 90 90	90 90	90 90 90

Run 2

R1	R2	R3	R4
A B	C D E	F G	H I
10 20	10 10 10	10 10	10 10 10
20 20	20 20 20	20 20	20 20 20
30 30	30 30 30	30 30	30 30 30
40 40	40 40 40	40 40	40 40 40
50 50	50 50 50	50 50	50 50 50
60 60	60 60 60	60 60	60 60 60
70 70	70 70 70	70 70	70 70 70
80 80	80 80 80	80 80	80 80 80
90 90	90 90 90	90 90	90 90 90

Merge Phase

Shaded tuples are in memory

R ₁	R ₂	R ₃	R ₄	Output so far:
A B	C D E	F G	H I	R ₁₂ R ₃₄
8 7 4	4 3 1 0	7 2 0	0 5 1	A B C D E F G H I
9 7 4	6 4 9 3	8 0 1	7 0 4	0 1 2 3 4
1 0 5	7 1 5 4	2 5 3	8 0 5	4 8 8 2 7
2 3 8	0 5 4 5	1 1 9	1 1 9	5 0 0 5
2 4 0	8 9 2 1	8 0 1	2 4 0	8 0 0 8
4 8 1	1 7 2 9	4 3 8	0 8 1	5 0 0 5
0 9 8	0 1 8 8	6 9 7	3 9 3	8 0 0 8
5 5 9	5 8 3 9	3 4 8	4 9 8	5 0 0 5

$H_1(B)^2$ $H_1(C)^2$ $H_2(G)^2$ $H_2(I)^2$ (a)

R ₁	R ₂	R ₃	R ₄	Output so far:
A B	C D E	F G	H I	R ₁₂ R ₃₄
8 7 4	4 3 1 0	7 2 0	0 5 1	A B C D E F G H I
9 7 4	6 4 9 3	5 0 1	7 0 4	0 1 2 3 5
1 0 5	7 1 5 4	7 5 9	8 5 5	4 8 8 2 7
2 3 8	0 5 4 5	1 1 9	1 1 9	5 0 0 5
2 4 0	8 9 2 1	8 0 1	1 4 9	8 0 0 8
4 8 1	1 7 2 9	8 0 1	2 5 5	5 0 0 5
0 9 8	0 1 8 8	6 9 7	0 2 1	8 0 0 8
5 5 9	5 8 3 9	6 9 7	4 9 8	5 0 0 5

$H_1(B)^2$ $H_1(C)^2$ $H_2(G)^2$ $H_2(I)^2$ (b)

R ₁	R ₂	R ₃	R ₄	Output so far:
A B	C D E	F G	H I	R ₁₂ R ₃₄
8 7 4	4 3 1 0	7 2 0	0 5 1	A B C D E F G H I
9 7 4	6 4 9 3	5 0 1	7 0 4	0 1 2 3 5
1 0 5	7 1 5 4	2 5 3	8 0 5	4 8 8 2 7
2 3 8	0 5 4 5	1 1 9	1 1 9	5 0 0 5
2 4 0	8 9 2 1	8 0 1	2 4 0	8 0 0 8
4 8 1	1 7 2 9	4 3 8	0 8 1	5 0 0 5
0 9 8	0 1 8 8	6 9 7	3 9 3	8 0 0 8
5 5 9	5 8 3 9	3 4 8	4 9 8	5 0 0 5

$H_1(B)^2$ $H_1(C)^2$ $H_2(G)^2$ $H_2(I)^2$ (c)

Estimation

- First levelwise step
 - Sum of aggregate of tuples discovered so far * size ratio

$$\alpha = \sum_{a=1}^p \left[\sum_{i=1}^n f(i) \right] \cdot \frac{\beta}{\sum_{a=1}^p [T(a, 1, n)] + \sum_{a=2}^p [T(a-1, b+1, n)]}$$

$$\beta = \frac{[R_1 \times R_2 \times \dots \times R_d]}{\sum_{a=1}^p [T(a, 1, n)] + \sum_{a=2}^p [T(a-1, b+1, n)]}$$

- What about subsequent levels?
 - Intermediate results grouped by join key – not completely random
 - Cardinality of intermediate results unknown due to pipelining

Solution to Semi-Randomness

- “Clumps”
 - View group of tuples w/ same join key from a merge phase as a single, indivisible output tuple
 - Estimate using clumps

$$f(i_1 \cdot i_2 \cdot \dots \cdot i_n) = \sum_{i_1 \in I_1} \sum_{i_2 \in I_2} \dots \sum_{i_n \in I_n} f(i_1' \cdot i_2' \cdot \dots \cdot i_n')$$

Solution to Unknown Relation Size

- We don't need to estimate the unknown size!
- Partition H's output space into p equi-sized ranges
 - Size ratio $\beta = \frac{p^n}{1+(p-1)n}$
 - Round-robin scan searches 1+(r-1)n combinations of runs

Putting All Estimators Together

- d+1 unbiased estimators for d level
- Take weighted average

$$N = \sum_{i=1}^{d+1} w_i N_i$$

$$\sigma^2(N) = \sum_{i=1}^{d+1} w_i^2 \sigma^2(N_i)$$

Different hash seeds for different levels, so estimators independent

- To choose weights:
 - Lagrangian multiplier: $f(w_1, \dots, w_d, \lambda) = \sigma^2(N) + \lambda(\sum_{i=1}^d w_i - 1)$
 - Variance minimized by

$$w_i = 1 / \left\{ \sigma^2(N_i) \sum_{j=1}^{d+1} \frac{1}{\sigma^2(N_j)} \right\}$$

Other Cases

- Other aggregates
 - GROUP BY: separate estimators for each group
 - So far only SUM aggregate considered
 - COUNT: also a simple version of SUM
 - AVG: function of SUM and COUNT
 - STD_DEV: function of AVG and COUNT
- Non-bushy query plan trees
 - SELECT SUM (R..A)
 - FROM R, S, T, U
 - WHERE R.A = S.A AND
 - R.B = T.B AND
 - R.C = U.C

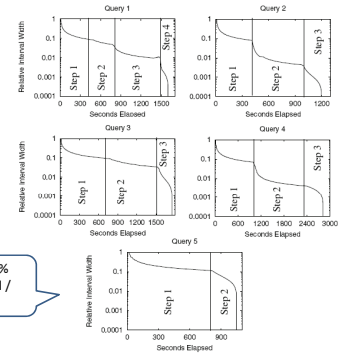
Special scan and re-randomize operator

Experiments

- Five queries over TPC-H schema
- 10GB generated data

19

Decrease of confidence interval width



y axis: width of 95% confidence interval / current estimate

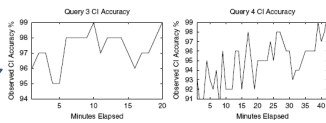
20

Running Time and Accuracy

Query Execution Time

Query	DBO	Postgres
Q ₁	26m42s	43m47s
Q ₂	20m08s	34m27s
Q ₃	29m12s	37m40s
Q ₄	47m05s	88m28s
Q ₅	17m28s	46m31s

Fraction of confidence intervals that do contain actual answer



21

22

Problems and Discussion

- Be organized or random?
 - Randomness provides unbiased online estimation
 - Indexing provides fast access
 - Can we make use of index? Can we access random data fast?
- Current system limits itself in many aspects
 - Sort-merge style join
 - Equality predicates
 - Are we losing the opportunity of optimization by this specific parallel way of query evaluation?
 - ...