

Adaptive ordering of pipelined stream Filters

S. Babu, R. Motwani, K. Munagala,
I. Nishizawa, and J. Widom
SIGMOD 2004

Presented by: Shimin Guo



Outline

- Introduction
- The Filter Ordering Problem
- The *A-Greedy* Algorithm
- The *Sweep* Algorithm
- The *Independent* Algorithm
- The *LocalSwaps* Algorithm
- Multiway Joins
- Experimental Evaluation
- Summary & Discussion

2



Introduction

- Streams processed by a set of *commutative* filters
- Overall processing costs depends on how the filters are ordered
- The best orderings are dependent on current stream and filter characteristics, which may change over time
- The selectivity of a filter depends on the filters before it
- Three-way tradeoff: convergence (C), run-time overhead (O), and speed of adaptivity (S)

3



The Filter Ordering Problem



- n commutative filters: F_1, F_2, \dots, F_n
- $f(\cdot)$: mapping from positions in the filters ordering to the indices of the filters at those positions
- O : an ordering $F_{f(1)}, F_{f(2)}, \dots, F_{f(n)}$
- $d(i|j)$: conditional probability that $F_{f(i)}$ will drop a tuple e , given that e was not dropped by any of $F_{f(1)}, F_{f(2)}, \dots, F_{f(j)}$
- t_i : expected time for F_i to process one tuple
- D_i : percentage of tuples that passed the first $i-1$ filters
- The goal: minimize $\sum_{i=1}^n t_{f(i)} D_i$

4

The A-Greedy Algorithm



- A greedy algorithm based on stable statistics:

☞ Choose the filter F_i with the highest $d(i|0)t_i$ as the first filter.

☞ Among the remaining filters, choose the filter F_j with the highest $d(j|1)t_j$ as the second filter.

☞ And so on.

- Greedy Invariant (GI):

$$\frac{d(i|i-1)}{t_{f(i)}} \geq \frac{d(j|i-1)}{t_{f(j)}}, 1 \leq i \leq j \leq n$$

- Goal of A-Greedy: maintain an ordering that satisfies the GI in an online manner

5

The A-Greedy Algorithm



- Two logical components of A-Greedy:

- profiler: continuously collects and maintains statistics about filter selectivities and processing costs
- reoptimizer: detects and corrects violations of the GI in the current filter ordering

- Challenge faced by the profiler: there are $n2^{n-1}$ conditional selectivities for n filters

- It's impractical for the profiler to maintain online estimates of all these selectivities

- Solution: *profile* of recently dropped tuples

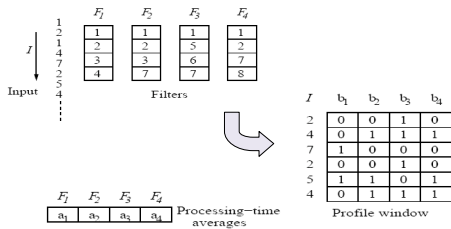
6

The A-Greedy Profiler

- Profile: a sliding window of *profile tuples*
- A profile tuple contains n boolean attributes b_1, \dots, b_n corresponding to the n filters
- Dropped tuples are sampled with some probability, called *drop-profiling probability*
- If a tuple e is chosen for profiling, it will be tested by all remaining filters
- A new profile tuple inserted into the profile window, where $b_i=1$ if F_i drops e and $b_i=0$ otherwise

7

The A-Greedy Profiler



8

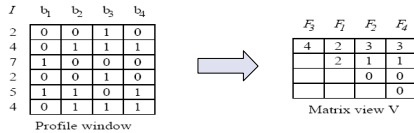
The A-Greedy Reoptimizer

- The reoptimizer maintains an ordering O such that O satisfies the GI for statistics estimated from the profile window
- How does the reoptimizer make use of the profile to derive estimates of conditional selectivities?
- It incrementally maintains a view over the profile window

9

The A-Greedy Reoptimizer

- View over the profile window: $n \times n$ upper triangular matrix V
- $V[i, j]$: number of tuples in the profile window that were dropped by $F_{f(j)}$ but not dropped by $F_{f(1)}, F_{f(2)}, \dots, F_{f(i-1)}$



10

The A-Greedy Reoptimizer

- $V[i, j]$ is proportional to $d(j|i-1)$
- Greedy Invariant:

$$\frac{d(i|i-1)}{t_{f(i)}} \geq \frac{d(j|i-1)}{t_{f(j)}}, \quad 1 \leq i < j \leq n$$

$$\frac{V[i, i]}{a_{f(i)}} \geq \frac{V[i, j]}{a_{f(j)}}, \quad 1 \leq i < j \leq n$$

(to avoid thrashing)

$$\frac{V[i, i]}{a_{f(i)}} \geq \alpha \frac{V[i, j]}{a_{f(j)}}, \quad 1 \leq i < j \leq n$$

11

The A-Greedy Reoptimizer

```

1. ** Input:  $O$  violates the Greedy Invariant at position  $i^*$ 
2.  $maxr = i$ 
3. for ( $r = i; r \leq maxr; r = r + 1$ ) {
4.   ** Greedy Invariant holds at positions  $1, \dots, r - 1$  and at
    $maxr + 1, \dots, n$ . Compute  $V$  entries for row  $r$ 
5.   for ( $c = r; c \leq n; c = c + 1$ )
6.      $V[r, c] = 0$ 
7.   for each profile tuple  $(b_1, b_2, \dots, b_n)$  in the profile window {
8.     ** Ignore tuples dropped by  $F_{f(1)}, F_{f(2)}, \dots, F_{f(i-1)}$ 
9.     if ( $b_{f(i)} == 0$  and  $b_{f(i+1)} == 0$  and  $\dots$  and  $b_{f(i-1)} == 0$ )
10.      for ( $c = r; c \leq n; c = c + 1$ )
11.        if ( $b_{f(c)} == 1$ )  $V[r, c] = V[r, c] + 1$ ;
12.      }
13.   }
14. ** Find the column  $maxc$  with maximum  $\frac{V[r, maxc]}{a_{f(maxc)}}$ 
15.  $maxc = r$ ;
16. for ( $m = r + 1; m \leq n; m = m + 1$ )
17.   if ( $\frac{V[r, m]}{a_{f(m)}} > \frac{V[r, maxc]}{a_{f(maxc)}}$ ) {
18.      $maxc = m$ ;
19.   }
20. }
21. if ( $r \neq maxc$ ) {
22.   ** Current filter  $F_{f(maxc)}$  becomes the new  $F_{f(i)}$ .
23.   ** We swap the filters at positions  $maxr$  and  $r$ .
24.   for ( $k = 0; k \leq r; k = k + 1$ )
25.     Swap  $V[k, r]$  and  $V[k, maxc]$ ;

```

Figure 5: Correcting a violation of the Greedy Invariant

12

Convergence Properties



THEOREM 4.1. *When stream and filter characteristics are stable, the cost of a filter ordering satisfying the GI is at most four times the cost of the optimal filter ordering.* □

- Constant factor depends on number of filters, e.g., 2.35, 2.61, and 2.8 for 20, 100, and 200 filters, respectively
- Usually finds the optimal ordering in practice

13

Run-time Overhead



- Profile-tuple creation: needs additional $n-i$ evaluations for a tuple dropped by $F_{R(i)}$. Creation frequency determined by drop-profiling probability
- Profile-window maintenance: insertion and deletion of profile tuples. Also needs to maintain running averages of filter processing times
- Matrix-view update: every update would cause access to up to $n^2/4$ entries
- Violation detection: access to up to n entries
- Violation correction: up to $n-i$ full scans of the profile window to correct a GI violation at position i

14

Speed of Adaptivity



- Any GI violation will be detected and corrected immediately
- Thus, A-Greedy is a very rapidly adapting algorithm

15

The A-Greedy Algorithm

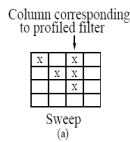
- A-Greedy has good convergence properties and extremely fast adaptivity, but it imposes significant run-time overhead
- Can we sacrifice some of A-Greedy's convergence properties or adaptivity speed to reduce its run-time overhead?



16

The Sweep Algorithm

- Proceeds in stages
- During one stage, only checks for GI violations involving the filter at one specific position j
- Does not need to maintain the entire matrix view
- Only $b_{f(1)}, \dots, b_{f(j)}$ are required in the profile window
- For each profiled tuple, needs to additionally evaluate $F_{f(j)}$ only
- By rotating j over $2, \dots, n$, eventually detects and corrects all GI violations



17

C, O, and S of Sweep

- Detects and corrects all GI violations
→ same convergence properties as A-Greedy
- Reduced view and need for additional evaluations
→ less overhead
- Only one filter is profiled in each stage
→ slower adaptivity



18

The *Independent* Algorithm

- Assumes filters are independent
- Only needs to maintain estimates of unconditional selectivities



19

C, O, and S of *Independent*

- Convergence: dependent on whether assumption holds
 - if so, optimal
 - otherwise, can be $O(n)$ times worse than GI orderings
- Lower view maintenance overhead
- Fast adaptivity

20

The *LocalSwaps* Algorithm

- Monitors “local” violations only, i.e., violations involving adjacent filters
- Intuitively, *LocalSwaps* detects situations where a swap between adjacent filters in the current ordering would improve performance
- Only needs to maintain two diagonals of the view
- For each profiled tuple dropped by $F_{f(i)}$, only needs to additionally evaluate $F_{f(i+1)}$



21

C, O, and S of LocalSwaps



- Convergence: path-dependent
 - Best case: converges to GI orderings
 - Worse case: can be $O(n)$ times worse than GI orderings
 - May get stuck in local maxima
- Lower profiling and view-maintenance overhead
- Restricted to local moves → takes longer to converge

22

Comparison of the four algorithms



	A-Greedy	Sweep	Independent	LocalSwaps
C	Good	Good	Optimal if independence assumption holds, $O(n)$ worse in general	Path-dependent Best case: same as A-Greedy Worse case: $O(n)$ worse
O	High	Low	Low	Low
S	Fast	Slow	Fast	Slow

23

Multiway Joins

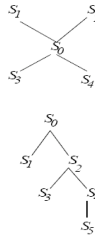


- MJoins maintain an ordering of $\{S_0, S_1, \dots, S_{n-1}\} - \{S_i\}$ for each stream S_i
- New tuples arriving from S_i is joined with other stream windows in that order
- Two-phase join algorithm
 - Drop-probing phase: the new tuple is used to probe all other windows in the specified order. If any window drops it, no further processing will be needed for it
 - output-generation phase: if no window drops the tuple, proceeds as conventional MJoins
- Drop probing resembles pipelined filters
- A-Greedy and its variants can be used to determine the orderings

24

Multiway Joins

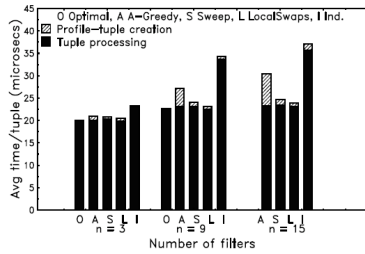
- Star Joins:
 - Tuple arrives from S_0 : straightforward
 - Tuple arrives from S_i : join with S_0 first, then apply the two-phase join algorithm for each tuple in $S_i \bowtie S_0$
- Acyclic Joins
 - Join graph defines a partial order
 - Join orderings constrained by the partial order



25

Experimental Evaluation

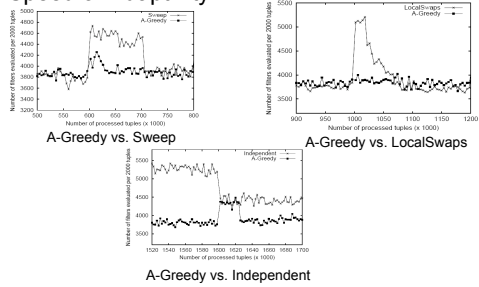
Convergence and overhead



26

Experimental Evaluation

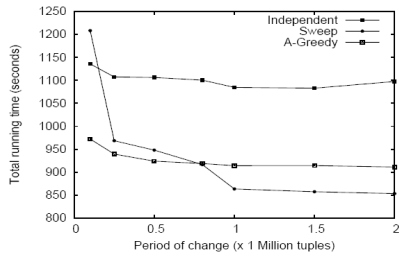
Speed of Adaptivity



27

Experimental Evaluation

Varying the rate of change



28

Summary

- A-Greedy handles correlated filters
- A-Greedy has good convergence properties, fast adaptivity, but incurs significant run-time overhead
- Three variants of A-Greedy are proposed, each lying at a different point along the tradeoff spectrum among convergence, run-time overhead, and speed of adaptivity

29

Discussion

- Given that each of the algorithms has different utility in different settings, can we add another level of adaptivity that adaptively choose the algorithm that best fits the current setting?
- How in reality can correlation among filters affect query optimizers that assume independent filters?
- Is online reordering feasible for join operators that maintain internal states?
- How to choose the size for the profile window?

30
