

CMU SCS

15-721 DB Sys. Design & Impl.

R* Optimizer

Christos Faloutsos
www.cs.cmu.edu/~christos

CMU SCS

Roadmap

- 1) Roots: System R and Ingres
- 2) Implementation: buffering, indexing, q-opt
- 3) Transactions: locking, recovery
- 4) Distributed DBMSs
 - intro
 - R* architecture
 - R* optimizer
- 5) Parallel DBMSs: Gamma, Alphasort
- 6) OO/OR DBMS
- 7) Data Analysis - data mining

15-721 C. Faloutsos 2

CMU SCS

Citation

Mackert & Lohman, "R* Optimizer Validation and Performance Evaluation for Distributed Queries" VLDB 1986.

15-721 C. Faloutsos 3

CMU SCS

Problem – definition

ideally: `connect to distr-LA; exec sql select * from EMPL;`

15-721 C. Faloutsos 4

CMU SCS

Overview

- Distr. computation and Optimization
- Instrumentation
- Distr. join results
- Alternative join strategies
- Conclusions

15-721 C. Faloutsos 5

CMU SCS

Introduction - targets of study

- Validation of q-optimizer
 - how often is the chosen plan sub-optimal?
 - which are the most influential parms?
 - sensitivity analysis
- Other questions
 - other improvements for distr. joins?
 - other promising techniques?

15-721 C. Faloutsos 6

CMU SCS

Distr. comp. & opt.

- Definitions / assumptions
 - each table: stored at one site (no partitioning)
 - dfn: query site; master site; apprentice sites

15-721 C. Faloutsos 7

CMU SCS

Distr. comp. & opt.

- in centralized q-opt:
 - join order; method; access path
- now: ++ join site
- how to send inner rel. to join site?
 - “ship whole” (W)
 - “fetch matches” (F) (~ semijoin)

15-721 C. Faloutsos 8

CMU SCS

Distr. comp. & opt.

- “fetch matches” (F)

15-721 C. Faloutsos 9

CMU SCS

Overview

- Distr. computation and Optimization
- Instrumentation
- Distr. join results
- Alternative join strategies
- Conclusions

15-721 C. Faloutsos 10

CMU SCS

Instrumentation

- ?

15-721 C. Faloutsos 11

CMU SCS

Instrumentation

- distributed EXPLAIN
 - (PLAN_TABLE @ master, + local ones)
- COLLECT COUNTERS
 - RSS stats, I/Os, buffer look-ups, comm. stats)
- FORCE OPTIMIZER

15-721 C. Faloutsos 12

CMU SCS

General measurements

Cost function:

- $\text{cpu} + \text{I/O} + \#\text{msgs} + \#\text{bytes}$
- w_{CPU} 0.0004 msec/instr (2.5MIPS)
- $w_{\text{I/O}}$ 23msec (actually: 17msec)
- w_{msg} 11msec (actually: 16msec)
- w_{byte} 0.002 msec/byte (4Mbits/sec)

15-721 C. Faloutsos 13

CMU SCS

Overview

- Distr. computation and Optimization
- Instrumentation
- ➔ • Distr. join results
- Alternative join strategies
- Conclusions

15-721 C. Faloutsos 14

CMU SCS

Distr. join results

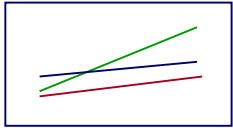
- Q1: how to ship data - ship 'whole', or 'fetch', or?
- A1:

15-721 C. Faloutsos 15

CMU SCS

Distr. join results

- Q1: how to ship data - ship 'whole', or 'fetch', or?
- A1: ship outer to inner and back



time

card. of outer table

Fetch whole
outer->inner & back

15-721 C. Faloutsos 16

CMU SCS

Distr. join results

- Q2: Distr. vs local joins
- A2:

15-721 C. Faloutsos 17

CMU SCS

Distr. join results

- Q2: Distr. vs local joins
- A2: Distr. have more overhead, but better parallelism and buffer contention
- (distinguish: response time vs resource time)

15-721 C. Faloutsos 18

Distr. join results

- Q3: relative importance of cost components?
(2.5MIPS, 20ms/IO, 11ms/msg, 4Mbps)
- A3:

Distr. join results

- Q3: relative importance of cost components?
- A3: local cost is important
 - high speed net.: comm. cost < 10%
 - medium speed net: local cost still not negligible
(50msec/msg, 40Kbps)

Distr. join results

- Q4: Opt. evaluation - how often sub-optimal?
- A4: usually OK - problems with join selectivity est.
 - solution:?

Distr. join results

- Q4: Opt. evaluation - how often sub-optimal?
- A4: usually OK - problems with join selectivity est.
 - solution: use estimate from previous execution!

Overview

- Distr. computation and Optimization
- Instrumentation
- Distr. join results
- ➡ • Alternative join strategies
- Conclusions

Alternative join strategies

- 1) temporary indices
- 2) semi-joins
- 3) “Bloom”-joins

CMU SCS

1) Temporary indices

- ship table T; build local index

15-721 C. Faloutsos 25

CMU SCS

2) Semijoins

- Idea: reduce the tables before shipping

S1

s#	...
s1	
s2	
s5	
s11	

(s1,s2,s5,s11)

s#	p#
s1	p1
s2	p1
s3	p5
s2	p9

S3

SUPPLIER Join SHIPMENT = ?

15-721 C. Faloutsos 26

CMU SCS

Semijoins

- Idea: reduce the tables before shipping
- Eg., to reduce 'SHIPMENT'
 - send distinct values of SUPPLIER.S#

15-721 C. Faloutsos 27

CMU SCS

Semijoins

- Formally:
- SHIPMENT' = SHIPMENT \bowtie SUPPLIER

15-721 C. Faloutsos 28

CMU SCS

A brilliant idea: two-way semijoins

- (not in book, not in final exam)
- reduce both relations with one more exchange: [Kang, '86]
- ship back the list of keys that didn't match

15-721 C. Faloutsos 29

CMU SCS

Two-way Semijoins

S1

s#	...
s1	
s2	
s5	
s11	

(s1,s2,s5,s11)

s#	p#
s1	p1
s2	p1
s3	p5
s2	p9

(s5,s11)

S3

SUPPLIER Join SHIPMENT = ?

S2

15-721 C. Faloutsos 30

CMU SCS

Two-way semijoins

- ship back the list of keys that didn't match
- CAN NOT LOSE! (why?)
- further improvement:
 - or the list of ones that matched – whatever is shorter!

15-721 C. Faloutsos 31

CMU SCS

3) 'Bloom-joins'

- how to ship the projection, say, of SUPPLIER.s#, even cheaper?
- A: Bloom-filter [Lohman+] =
 - quick&dirty membership testing

15-721 C. Faloutsos 32

CMU SCS

Bloom-join

- Idea: reduce table - using only, say, 10 bits? !!

15-721 C. Faloutsos 33

CMU SCS

Bloom-join

- idea: use a bit-string and hashing
 - may have “false alarms” (OK!)

15-721 C. Faloutsos 34

CMU SCS

Bloom-join

- idea: use a bit-string and hashing
 - may have “false alarms” (OK!)

15-721 C. Faloutsos 35

CMU SCS

Bloom-join

- idea: use a bit-string and hashing
 - may have “false alarms” (OK!)

15-721 C. Faloutsos 36

CMU SCS

Bloom filters

- could set $m > 1$ bits per value
- used for text retrieval (“Zato-coding”, in ‘49(!); signature files)
- differential files [Lohman+Severance]
- UNIX’s spell checker [McIlroy - IEEE COM’82]
- membership testing, in general

15-721 C. Faloutsos 37

CMU SCS

Bloom join

- Q1: How many false alarms, if we have
 - F bits, and
 - Ds (# of distinct values in table ‘S’)
 - Ct (# of tuples in table ‘T’)
 - SCt (# of tuples in ‘T’ that match)
- Q1’: How many ‘1’s, in the Bloom filter?
- A1’:

15-721 C. Faloutsos 38

CMU SCS

Bloom join

- Q1: How many false alarms, if we have
 - F bits, and
 - Ds (# of distinct values in table ‘S’)
 - Ct (# of tuples in table ‘T’)
 - SCt (# of tuples in ‘T’ that match)
- Q1’: How many ‘1’s, in the Bloom filter?
- A1’: $bits_S = F(1 - (1 - 1/F)^{Ds}) \sim F(1 - \exp(-Ds/F))$

15-721 C. Faloutsos 39

CMU SCS

Bloom join

- Q1: How many false alarms, if we have
 - F bits, and
 - Ds (# of distinct values in table ‘S’)
 - Ct (# of tuples in table ‘T’)
 - SCt (# of tuples in ‘T’ that match)
- A1: $BCt = SCt + bits_S/F * (Ct - SCt)$
– (book: slightly different formula)

15-721 C. Faloutsos 40

CMU SCS

Comparison

- A join B
- A: 1,000 tuples (query site)
- B: 100 - 6,000 tuples
- F=4Kb for bloom filter
- high/medium speed network
- R*, R*+temp-index, semijoin, bloom-join

15-721 C. Faloutsos 41

CMU SCS

Comparison

time

size of ‘A’

card. of ‘B’

semijoins

R* (dist)

temp ind

R* (local)

bloomj.

15-721 C. Faloutsos 42

Comparison

- Q: why bloom joins are better than semijoins?

Comparison

- Q: why bloom joins are better than semijoins?
- A: lower **local** processing! (simple scan of 'B')
- (similar results for 'B' being the query site)

Comparison - slower network

- bloom joins outperform all distributed algo's
 - (fewer bytes shipped)

Conclusions

- ship whole inner table wins
- R* optimizer: accurate
- distribution of queries: often helps, due to
 - parallelism
 - more buffers
- local cost: not negligible