

Database Systems CSE 414

Lectures 16 – 17:
Basics of Query Optimization and
Cost Estimation
(Ch. 15.{1,3,4,6,6} & 16.4-5)

CSE 414 - Fall 2017 1

Announcements

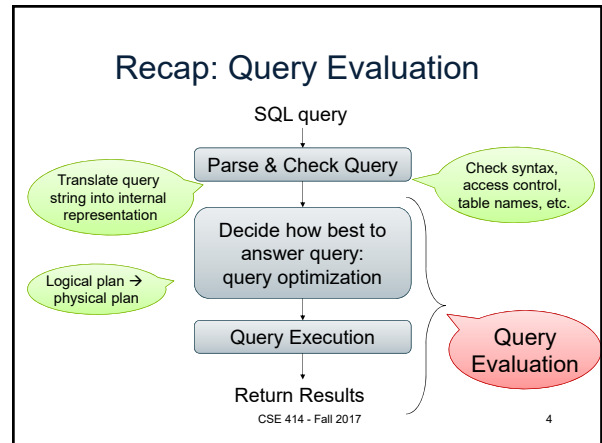
- WQ4 is due Friday 11pm
- HW3 is due next Tuesday 11pm
- Midterm is next Monday

CSE 414 - Fall 2017 2

Motivation

- To understand performance, need to understand a bit about how a DBMS works
 - my database application is too slow... why?
 - one of the queries is very slow... why?
- Under your direct control: index choice
 - understand how that affects query performance

CSE 414 - Fall 2017 3



Query Optimizer Overview

- **Input:** Parsed & checked SQL
- **Output:** A good physical query plan
- **Basic query optimization algorithm:**
 - Enumerate alternative plans (logical and physical)
 - Compute estimated cost of each plan
 - Compute number of I/Os
 - Optionally take into account other resources
 - Choose plan with lowest cost
 - This is called cost-based optimization

CSE 414 - Fall 2017 5

Query Optimizer Overview

- There are exponentially many query plans
 - exponential in the size of the query
 - simple SFW with 3 joins does not have too many
- Optimizer will consider many, many of them
- Worth substantial cost to avoid **bad plans**

CSE 414 - Fall 2017 6

Rest of Today

- Cost of reading from disk
- Cost of single RA operators
- Cost of query plans

CSE 414 - Fall 2017 7

Cost of Reading Data From Disk

CSE 414 - Fall 2017 8

Cost Parameters

- **Cost = Disk I/O + CPU + Network I/O**
 - We will focus on Disk I/O
- **Parameters:**
 - **B(R)** = # of blocks (i.e., pages) for relation R
 - **T(R)** = # of tuples in relation R
 - **V(R, A)** = # of distinct values of attribute A
 - When **A** is a key, **V(R, A) = T(R)**
 - When **A** is not a key, **V(R, A)** can be anything < **T(R)**
- Where do these values come from?
 - DBMS collects **statistics** about data on disk

CSE 414 - Fall 2017 9

Selectivity Factors for Conditions

- **A = c** /* $\sigma_{A=c}(R)$ */
 - Selectivity = $1/V(R, A)$
- **A < c** /* $\sigma_{A<c}(R)$ */
 - Selectivity = $(c - Low(R, A)) / (High(R, A) - Low(R, A))$
- **c1 < A < c2** /* $\sigma_{c1<A<c2}(R)$ */
 - Selectivity = $(c2 - c1) / (High(R, A) - Low(R, A))$

CSE 414 - Fall 2017 10

Example: Selectivity of $\sigma_{A=c}(R)$

T(R) = 100,000
V(R, A) = 20

How many records are returned by $\sigma_{A=c}(R)$ = ?

Answer: $X * T(R)$, where X = selectivity...
... $X = 1/V(R, A) = 1/20$

Number of records returned = $100,000/20 = 5,000$

CSE 414 - Fall 2017 11

Cost of Index-based Selection

- Sequential scan for relation R costs **B(R)**
- Index-based selection
 - Estimate selectivity factor **X** (see previous slide)
 - Clustered index: $X * B(R)$
 - Unclustered index $X * T(R)$

Note: we are ignoring I/O cost for index pages

CSE 414 - Fall 2017 12

Example: Cost of $\sigma_{A=c}(R)$

- Example:

$B(R) = 2000$
$T(R) = 100,000$
$V(R, A) = 20$

cost of $\sigma_{A=c}(R) = ?$
- Table scan: $B(R) = 2,000$ I/Os
- Index based selection:
 - If index is clustered: $B(R)/V(R, A) = 100$ I/Os
 - If index is unclustered: $T(R)/V(R, A) = 5,000$ I/Os

Lesson: Don't build unclustered indexes when $V(R, A)$ is small !

CSE 414 - Fall 2017 13

Cost of Executing Operators (Focus on Joins)

CSE 414 - Fall 2017 14

Outline

- Join operator algorithms**
 - One-pass algorithms (Sec. 15.2 and 15.3)
 - Index-based algorithms (Sec 15.6)
- Note about readings:
 - In class, we discuss only algorithms for joins
 - Other operators are easier: read the book

CSE 414 - Fall 2017 15

Join Algorithms

- Hash join
- Nested loop join
- Sort-merge join

CSE 414 - Fall 2017 16

Hash Join

Hash join: $R \bowtie S$

- Scan R, build buckets in main memory
- Then scan S and join
- Cost: $B(R) + B(S)$
- One-pass algorithm when $B(R) \leq M$ (memory size)
 - more disk access also when $B(R) > M$

CSE 414 - Fall 2017 17

Hash Join Example

Patient(pid, name, address)
Insurance(pid, provider, policy_nb)
Patient \bowtie Insurance

Patient	Insurance
1 'Bob' 'Seattle'	2 'Blue' 123
2 'Ela' 'Everett'	4 'Prem' 432
3 'Jill' 'Kent'	4 'Prem' 343
4 'Joe' 'Seattle'	3 'GrpH' 554

Two tuples per page

CSE 414 - Fall 2017 18

Hash Join Example

Patient \bowtie Insurance

Memory M = 21 pages

Large enough

Showing pid only

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

This is one page with two tuples

CSE 414 - Fall 2017 19

Hash Join Example

Step 1: Scan Patient and build hash table in memory

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

CSE 414 - Fall 2017 20

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer: 2 4

Output buffer: 2 2

Write to disk

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

CSE 414 - Fall 2017 21

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer: 2 4

Output buffer: 4 4

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

CSE 414 - Fall 2017 22

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer: 4 3

Output buffer: 4 4

Keep going until read all of Insurance

Cost: B(R) + B(S)

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

CSE 414 - Fall 2017 23

Nested Loop Joins

- Tuple-based nested loop $R \bowtie S$
- R is the outer relation, S is the inner relation

```

for each tuple t1 in R do
  for each tuple t2 in S do
    if t1 and t2 join then output (t1, t2)
    
```

What is the Cost?

- Cost: B(R) + T(R) B(S)
- Multiple-pass because S is read many times

CSE 414 - Fall 2017 24

Block-at-a-time Refinement

for each block of tuples r in R do
 for each block of tuples s in S do
 for all pairs of tuples t₁ in r, t₂ in s
 if t₁ and t₂ join then output (t₁, t₂)

- Cost: B(R) + B(R) B(S) What is the Cost?

CSE 414 - Fall 2017 25

Block-at-a-time Refinement

CSE 414 - Fall 2017 26

Block-at-a-time Refinement

CSE 414 - Fall 2017 27

Page-at-a-time Refinement

CSE 414 - Fall 2017 28

Block-at-a-time Refinement

Cost: B(R) + B(R)B(S)

CSE 414 - Fall 2017 29

Block-Nested-Loop Refinement

for each group of M-1 pages r in R do
 for each page of tuples s in S do
 for all pairs of tuples t₁ in r, t₂ in s
 if t₁ and t₂ join then output (t₁, t₂)

- Cost: B(R) + B(R)B(S)/(M-1) What is the Cost?

CSE 414 - Fall 2017 30

Sort-Merge Join

Sort-merge join: $R \bowtie S$

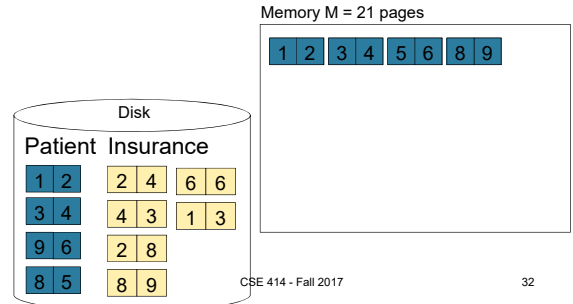
- Scan R and sort in main memory
- Scan S and sort in main memory
- Merge R and S
- Cost: $B(R) + B(S)$
- One pass algorithm when $B(S) + B(R) \leq M$
- Typically, this is NOT a one pass algorithm

CSE 414 - Fall 2017

31

Sort-Merge Join Example

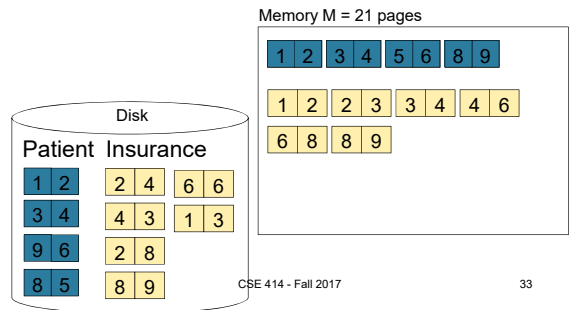
Step 1: Scan Patient and **sort** in memory



32

Sort-Merge Join Example

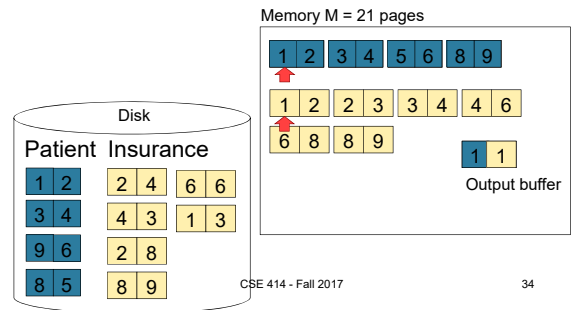
Step 2: Scan Insurance and **sort** in memory



33

Sort-Merge Join Example

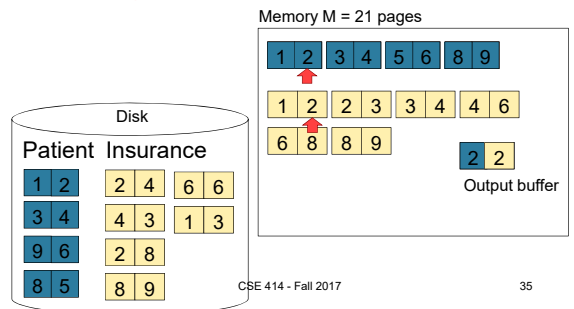
Step 3: **Merge** Patient and Insurance



34

Sort-Merge Join Example

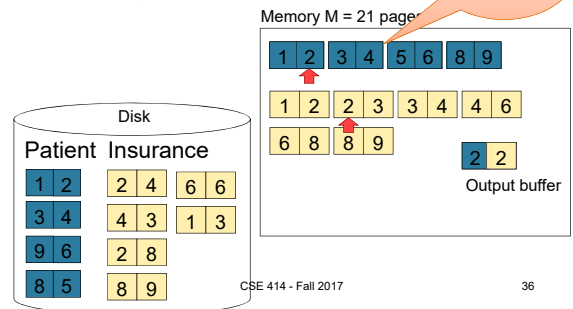
Step 3: **Merge** Patient and Insurance



35

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance



36

Sort-Merge Join Example

Step 3: Merge Patient and Insurance

Memory M = 21 pages

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

Memory M = 21 pages

1 2 3 4 5 6 8 9
1 2 2 3 3 4 4 6
6 8 8 9

Output buffer

3 3

CSE 414 - Fall 2017 37

Sort-Merge Join Example

Step 3: Merge Patient and Insurance

Memory M = 21 pages

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

Memory M = 21 pages

1 2 3 4 5 6 8 9
1 2 2 3 3 4 4 6
6 8 8 9

Output buffer

3 3

Keep going until end of first relation

CSE 414 - Fall 2017 38

Index Nested Loop Join

$R \bowtie S$

- Assume S has an index on the join attribute
- Iterate over R, for each tuple, fetch corresponding tuple(s) from S

Cost:

- If index on S is clustered: $B(R) + T(R)B(S)/V(S, A)$
- If index on S is unclustered: $B(R) + T(R)T(S)/V(S, A)$

CSE 414 - Fall 2017 39

Cost of Query Plans

CSE 414 - Fall 2017 40

$T(\text{Supplier}) = 1000$ $B(\text{Supplier}) = 100$ $V(\text{Supplier, scity}) = 20$ $M = 11$
 $T(\text{Supply}) = 10,000$ $B(\text{Supply}) = 100$ $V(\text{Supplier, sstate}) = 10$
 $V(\text{Supply, pno}) = 2,500$

Physical Query Plan 1

(On the fly) π_{sname} Selection and project on-the-fly -> No additional cost.

(On the fly) $\sigma_{\text{scity}='Seattle' \wedge \text{sstate}='WA' \wedge \text{pno}=2}$

(Nested loop) $\text{sno} = \text{sno}$

Supplier (File scan) Supply (File scan)

Total cost of plan is thus cost of join:
 $= B(\text{Supplier}) + B(\text{Supplier}) * B(\text{Supply})$
 $= 100 + 100 * 100$
 $= 10,100 \text{ I/Os}$

CSE 414 - Fall 2017 41

$T(\text{Supplier}) = 1000$ $B(\text{Supplier}) = 100$ $V(\text{Supplier, scity}) = 20$ $M = 11$
 $T(\text{Supply}) = 10,000$ $B(\text{Supply}) = 100$ $V(\text{Supplier, sstate}) = 10$
 $V(\text{Supply, pno}) = 2,500$

Physical Query Plan 2

(On the fly) π_{sname} (d)

(Sort-merge join) $\text{sno} = \text{sno}$ (c)

(Scan write to T1) (a) $\sigma_{\text{scity}='Seattle' \wedge \text{sstate}='WA'}$ (Scan write to T2) (b) $\sigma_{\text{pno}=2}$

Supplier (File scan) Supply (File scan)

Total cost
 $= 100 + 100 * 1/20 * 1/10$ (a)
 $+ 100 + 100 * 1/2500$ (b)
 $+ 2$ (c)
 $+ 0$ (d)
 Total cost $\approx 204 \text{ I/Os}$

CSE 414 - Fall 2017 42

