

Database Systems CSE 414

Lectures 11 – 12:
Basics of Query Optimization and
Cost Estimation
(Ch. 15.{1,3,4,6,6} & 16.4-5)

CSE 414 - Spring 2017

1

Announcements

- HW3 is due Tuesday
- WQ4 is due Thursday
- Midterm on Friday
 - we'll talk more about it on Monday
- Husky Football spring game tomorrow

CSE 414 - Spring 2017

2

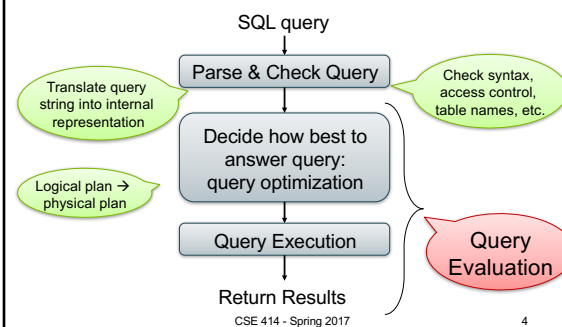
Motivation

- 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 - Spring 2017

3

Recap: Query Evaluation



CSE 414 - Spring 2017

4

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 - Spring 2017

5

Query Optimizer Overview

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

CSE 414 - Spring 2017

6

Rest of Today

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

CSE 414 - Spring 2017

7

Cost of Reading Data From Disk

CSE 414 - Spring 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 - Spring 2017

9

Selectivity Factors for Conditions

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

CSE 414 - Spring 2017

10

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

$$\begin{array}{l} T(R) = 100,000 \\ V(R, A) = 20 \end{array}$$

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 - Spring 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 - Spring 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 !

Cost of Executing Operators (Focus on Joins)

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

Join Algorithms

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

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 $B(R) \leq M$
 - more disk access also when $B(R) > M$

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

Hash Join Example

Patient \bowtie Insurance

Memory M = 21 pages

Showing pid only

Large enough

This is one page with two tuples

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

Hash Join Example

Step 1: Scan Patient and **build** hash table in memory

Memory M = 21 pages

Hash h: pid % 5

Input buffer

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

Hash Join Example

Step 2: Scan Insurance and **probe** into hash table

Memory M = 21 pages

Hash h: pid % 5

Input buffer

Output buffer

Write to disk

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

Hash Join Example

Step 2: Scan Insurance and **probe** into hash table

Memory M = 21 pages

Hash h: pid % 5

Input buffer

Output buffer

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

Hash Join Example

Step 2: Scan Insurance and **probe** into hash table

Memory M = 21 pages

Hash h: pid % 5

Input buffer

Output buffer

Keep going until read all of Insurance

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

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

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 since S is read many times

CSE 414 - Spring 2017

Block-at-a-time Refinement

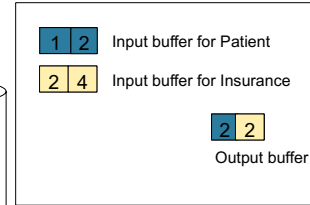
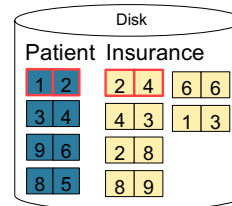
```

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)
    
```

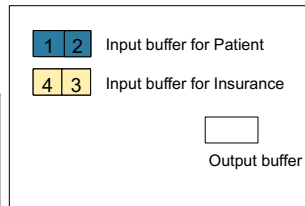
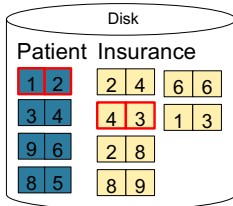
- Cost: $B(R) + B(R)B(S)$

What is the Cost?

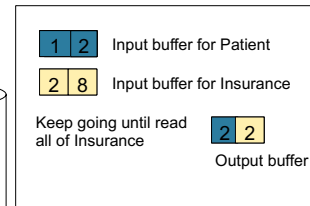
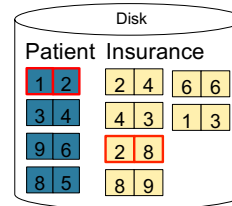
Block-at-a-time Refinement



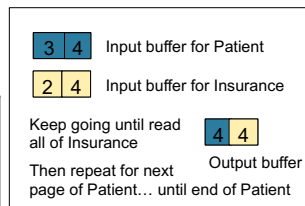
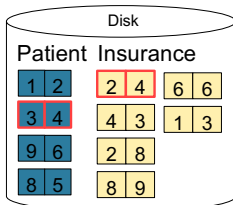
Block-at-a-time Refinement



Page-at-a-time Refinement



Block-at-a-time Refinement



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

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?

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 - Spring 2017 31

Sort-Merge Join Example

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

Memory M = 21 pages

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

Disk

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

32

Sort-Merge Join Example

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

Memory M = 21 pages

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

Disk

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

33

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance

Memory M = 21 pages

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

Disk

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

Output buffer

1	1
---	---

34

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance

Memory M = 21 pages

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

Disk

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

Output buffer

2	2
---	---

35

Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance

Memory M = 21 pages

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

Disk

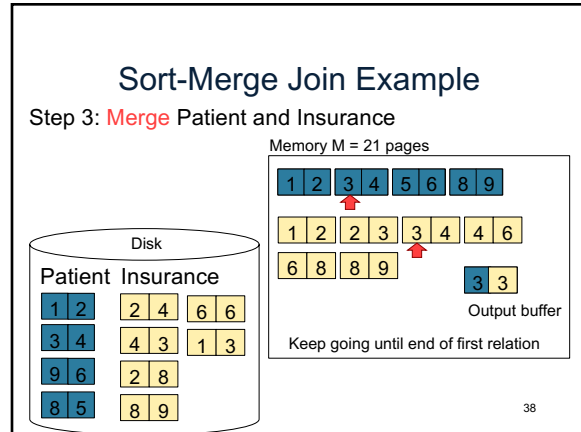
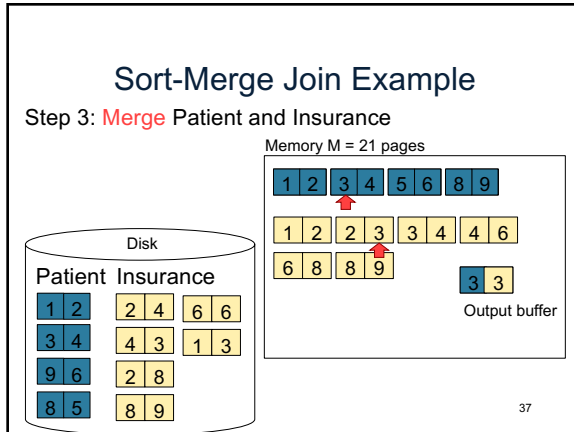
Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

Output buffer

2	2
---	---

Using PK, so only one can match

36



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 - Spring 2017 39

Cost of Query Plans

CSE 414 - Spring 2017 40

