

Database System Internals  
Query Execution and Algorithms

Paul G. Allen School of Computer Science and Engineering  
University of Washington, Seattle

January 22, 2020 CSE 444 - Winter 2020 1

1

### Announcements

- Homework 2 released
  - Due January 31<sup>st</sup>
- 544 paper 1 report due Friday

January 22, 2020 CSE 444 - Winter 2020 2

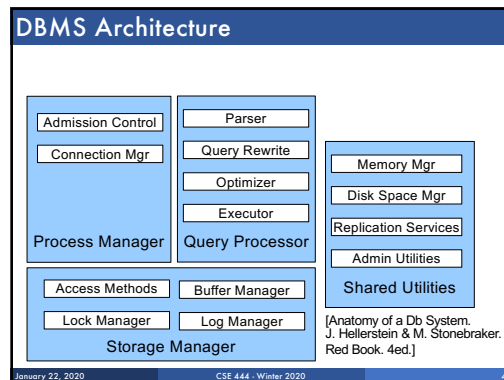
2

### What We Have Learned So Far

- Overview of the architecture of a DBMS
- Access methods
  - Heap files, sequential files, Indexes (hash or B+ trees)
- Role of buffer manager
- Practiced the concepts in hw1 and lab1

January 22, 2020 CSE 444 - Winter 2020 3

3



4

### Next Lectures

- How to answer queries **efficiently!**
  - Physical query plans and operator algorithms**
- How to automatically find good query plans
  - How to compute the cost of a complete plan
  - How to pick a good query plan for a query
  - i.e., Query optimization

January 22, 2020 CSE 444 - Winter 2020 5

5

### Query Execution Bottom Line

- SQL query transformed into **physical plan**
  - Access path selection** for each relation
  - Implementation choice** for each operator
  - Scheduling decisions** for operators
    - Single-threaded or parallel, pipelined or with materialization, etc.
- Execution of the physical plan is pull-based
- Operators **given a limited amount of memory**

January 22, 2020 CSE 444 - Winter 2020 6

6



### Operator Algorithms

Design criteria

- Cost: IO, CPU, Network
- Memory utilization
- Load balance (for parallel operators)

January 22, 2020 CSE 444 - Winter 2020 13

13

### Cost Parameters

- **Cost = total number of I/Os**
  - This is a simplification that ignores CPU, network
- **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  $< T(R)$

January 22, 2020 CSE 444 - Winter 2020 14

14

### Convention

- Cost = the cost of **reading** operands from disk
- Cost of **writing** the **final** result to disk is *not included*; need to count it separately when applicable

January 22, 2020 CSE 444 - Winter 2020 15

15

### Outline

- **Join operator algorithms**
  - One-pass algorithms (Sec. 15.2 and 15.3)
  - Index-based algorithms (Sec 15.6)
  - Two-pass algorithms (Sec 15.4 and 15.5)
- Note about readings:
  - In class, we discuss only algorithms for joins
  - Other operators are easier: book has extra details

January 22, 2020 CSE 444 - Winter 2020 16

16

### Join Algorithms

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

January 22, 2020 CSE 444 - Winter 2020 17

17

### 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  $B(R) \leq M$

January 22, 2020 CSE 444 - Winter 2020 18

18

### Hash Join Example

Patient(pid, name, address)  
Insurance(pid, provider, policy\_nb)  
Patient  $\bowtie$  Insurance

Two tuples per page

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

January 22, 2020 CSE 444 - Winter 2020 19

### Hash Join Example

Patient  $\bowtie$  Insurance

Memory M = 21 pages

Showing pid only

Some large-enough nb

This is one page with two tuples

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

January 22, 2020 CSE 444 - Winter 2020 20

### Hash Join Example

Step 1: Scan Patient and **build** hash table in memory  
Can be done in method open()

Memory M = 21 pages  
Hash h: pid % 5

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

Input buffer

January 22, 2020 CSE 444 - Winter 2020 21

### Hash Join Example

Step 2: Scan Insurance and **probe** into hash table  
Done during calls to next()

Memory M = 21 pages  
Hash h: pid % 5

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

Input buffer: 2 4

Output buffer: 2 2

Write to disk or pass to next operator

January 22, 2020 CSE 444 - Winter 2020 22

### Hash Join Example

Step 2: Scan Insurance and **probe** into hash table  
Done during calls to next()

Memory M = 21 pages  
Hash h: pid % 5

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

Input buffer: 2 4

Output buffer: 4 4

January 22, 2020 CSE 444 - Winter 2020 23

### Hash Join Example

Step 2: Scan Insurance and **probe** into hash table  
Done during calls to next()

Memory M = 21 pages  
Hash h: pid % 5

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

Input buffer: 4 3

Output buffer: 4 4

Keep going until read all of Insurance

Cost: B(R) + B(S)

January 22, 2020 CSE 444 - Winter 2020 24

## Nested Loop Joins

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

```
for each tuple  $t_1$  in  $R$  do
  for each tuple  $t_2$  in  $S$  do
    if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

What is the Cost?

25

## Nested Loop Joins

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

```
for each tuple  $t_1$  in  $R$  do
  for each tuple  $t_2$  in  $S$  do
    if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

What is the Cost?

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

26

## Page-at-a-time Refinement

```
for each page of tuples  $r$  in  $R$  do
  for each page of tuples  $s$  in  $S$  do
    for all pairs of tuples  $t_1$  in  $r$ ,  $t_2$  in  $s$ 
      if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

What is the Cost?

27

## Page-at-a-time Refinement

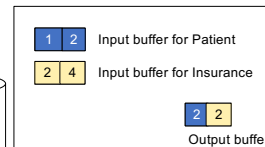
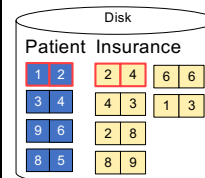
```
for each page of tuples  $r$  in  $R$  do
  for each page of tuples  $s$  in  $S$  do
    for all pairs of tuples  $t_1$  in  $r$ ,  $t_2$  in  $s$ 
      if  $t_1$  and  $t_2$  join then output  $(t_1, t_2)$ 
```

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

What is the Cost?

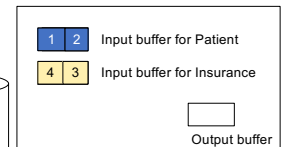
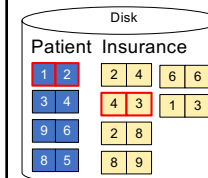
28

## Page-at-a-time Refinement



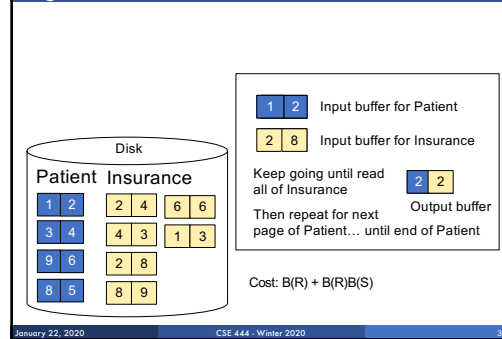
29

## Page-at-a-time Refinement



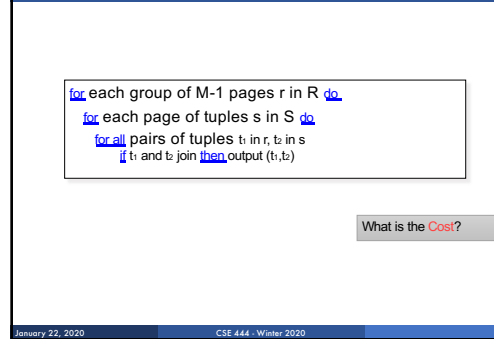
30

## Page-at-a-time Refinement



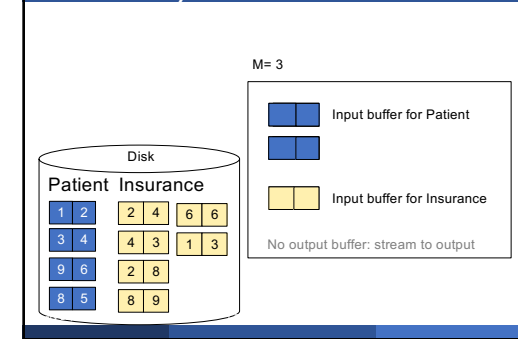
31

## Block-Memory Refinement



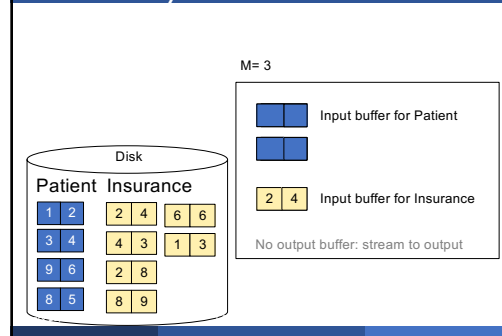
32

## Block Memory Refinement



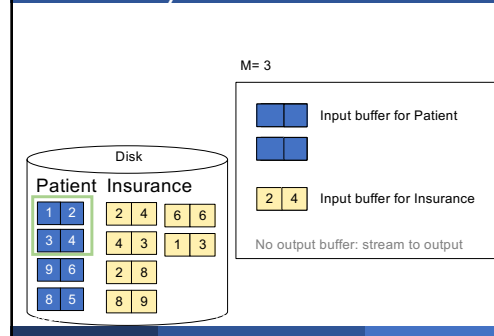
33

## Block Memory Refinement



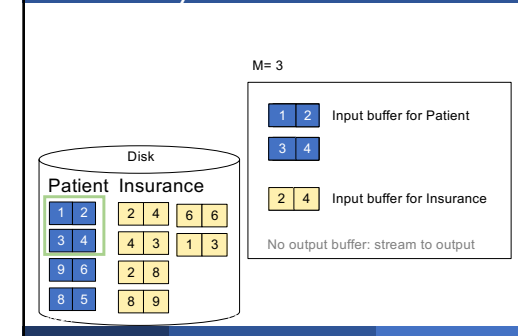
34

## Block Memory Refinement

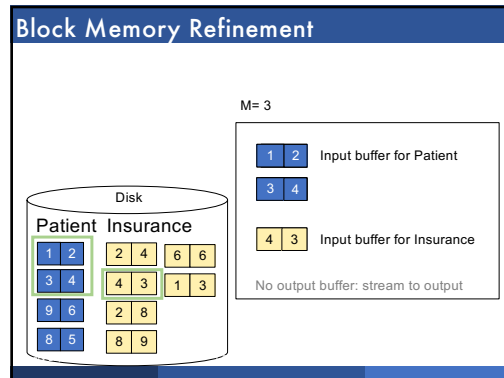


35

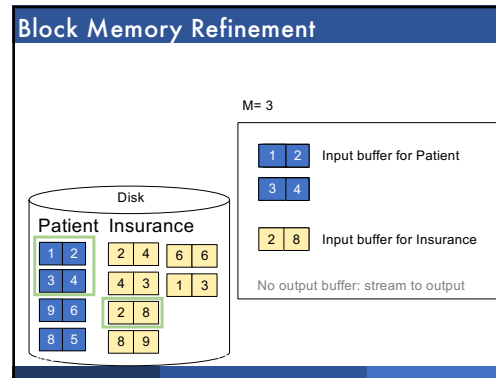
## Block Memory Refinement



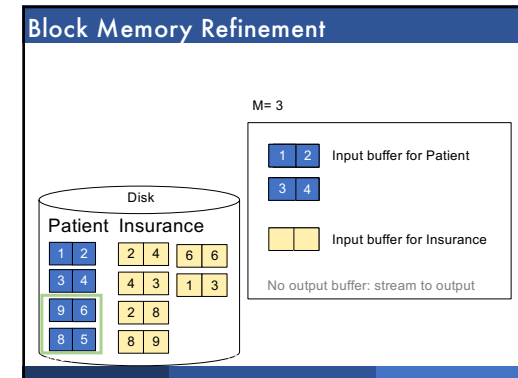
36



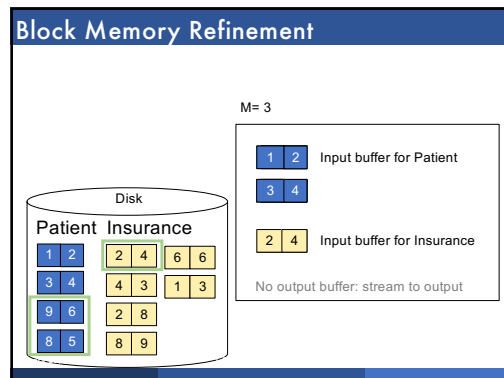
37



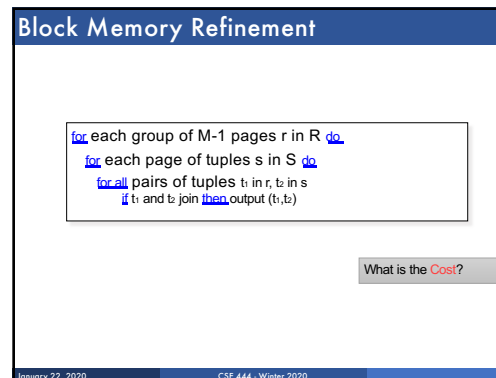
38



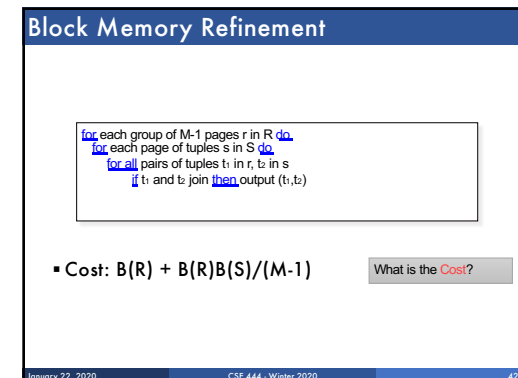
39



40



41



42

### Sort-Merge Join

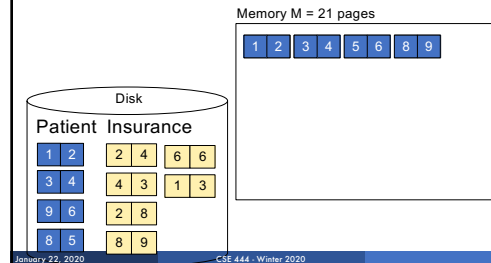
Sort-merge join:  $R \bowtie S$

- 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,
  - We'll see the multi-pass version next lecture

43

### Sort-Merge Join Example

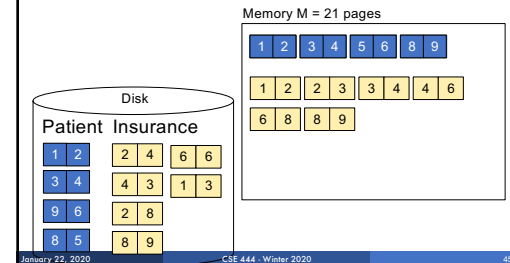
Step 1: Scan Patient and **sort** in memory



44

### Sort-Merge Join Example

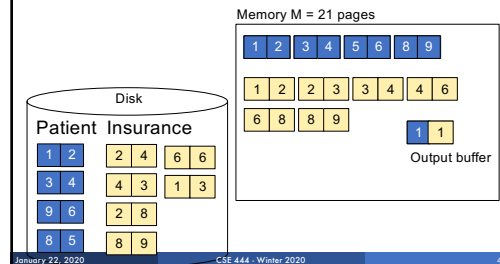
Step 2: Scan Insurance and **sort** in memory



45

### Sort-Merge Join Example

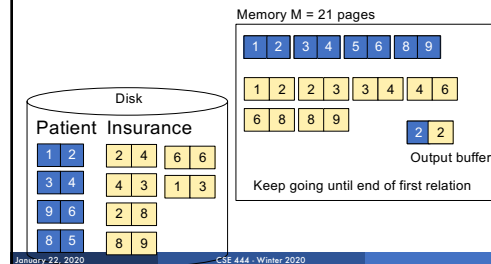
Step 3: **Merge** Patient and Insurance



46

### Sort-Merge Join Example

Step 3: **Merge** Patient and Insurance



47

### Outline

- **Join operator algorithms**
  - One-pass algorithms (Sec. 15.2 and 15.3)
  - Index-based algorithms (Sec 15.6)
  - Two-pass algorithms (Sec 15.4 and 15.5)

48



## Index Based Selection

Selection on equality:  $\sigma_{a=v}(R)$ 

- $B(R)$  = size of R in blocks
- $T(R)$  = number of tuples in R
- $V(R, a)$  = # of distinct values of attribute a

January 22, 2020 CSE 444 - Winter 2020 49

49

## Index Based Selection

Selection on equality:  $\sigma_{a=v}(R)$ 

- $B(R)$  = size of R in blocks
- $T(R)$  = number of tuples in R
- $V(R, a)$  = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a:
- Unclustered index on a:

January 22, 2020 CSE 444 - Winter 2020 50

50

## Index Based Selection

Selection on equality:  $\sigma_{a=v}(R)$ 

- $B(R)$  = size of R in blocks
- $T(R)$  = number of tuples in R
- $V(R, a)$  = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a:  $B(R)/V(R, a)$
- Unclustered index on a:  $T(R)/V(R, a)$

January 22, 2020 CSE 444 - Winter 2020 51

51

## Index Based Selection

Selection on equality:  $\sigma_{a=v}(R)$ 

- $B(R)$  = size of R in blocks
- $T(R)$  = number of tuples in R
- $V(R, a)$  = # of distinct values of attribute a

What is the cost in each case?

- Clustered index on a:  $B(R)/V(R, a)$
- Unclustered index on a:  $T(R)/V(R, a)$

Note: we ignore I/O cost for index pages

January 22, 2020 CSE 444 - Winter 2020 52

52

## Index Based Selection

Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ 

- Table scan:
- Index based selection:

January 22, 2020 CSE 444 - Winter 2020 53

53

## Index Based Selection

Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ 

- Table scan:  $B(R) = 2,000$  I/Os
- Index based selection:

January 22, 2020 CSE 444 - Winter 2020 54

54

## Index Based Selection

## ▪ Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ ▪ Table scan:  $B(R) = 2,000$  I/Os

▪ Index based selection:

- If index is clustered:
- If index is unclustered:

January 22, 2020 CSE 444 - Winter 2020 55

55

## Index Based Selection

## ▪ Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ ▪ 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:

January 22, 2020 CSE 444 - Winter 2020 56

56

## Index Based Selection

## ▪ Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ ▪ 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

January 22, 2020 CSE 444 - Winter 2020 57

57

## Index Based Selection

## ▪ Example:

$$\begin{array}{l} B(R) = 2000 \\ T(R) = 100,000 \\ V(R, a) = 20 \end{array}$$
cost of  $\sigma_{a=v}(R) = ?$ ▪ 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 !

January 22, 2020 CSE 444 - Winter 2020 58

58

## 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:  $B(R) + T(R)T(S)/V(S,a)$

January 22, 2020 CSE 444 - Winter 2020 59

59