

## CSE 444: Database Internals

### Lecture 7 Query Execution and Operator Algorithms (part 1)

CSE 444 - Spring 2015

1

## Announcements

- Lab 2 / part 1 due Friday, 11pm
- CSE544M: review 2 due today, 11pm

CSE 444 - Spring 2015

2

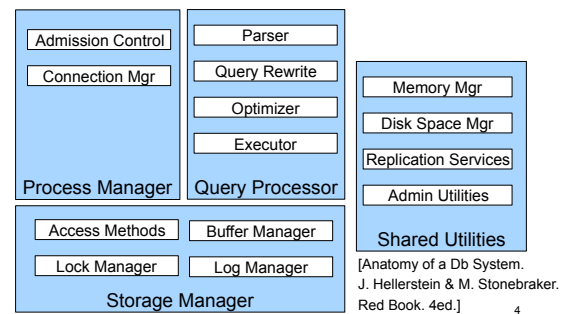
## What We Have Learned So Far

- 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

CSE 444 - Spring 2015

3

## DBMS Architecture



4

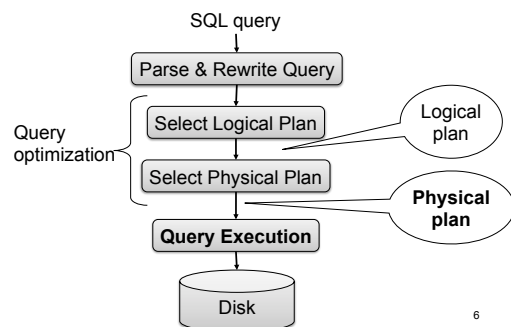
## Next Lectures

- 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

CSE 444 - Spring 2015

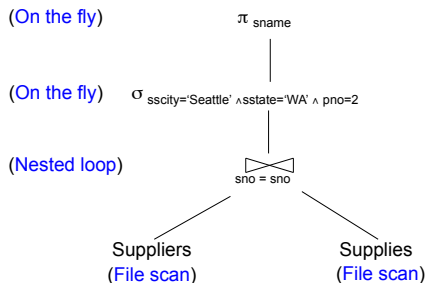
5

## Query Evaluation Steps Review



6

## Physical Query Plan



7

## Physical Query Plan

- **Access path selection** for each relation:
  - File scan, or
  - Index lookup with a predicate
- **Implementation choice** for each operator
  - We will learn different algorithms
- **Scheduling decisions** for operators
  - Pipelined execution, or
  - Intermediate tuple materialization

CSE 444 - Spring 2015

8

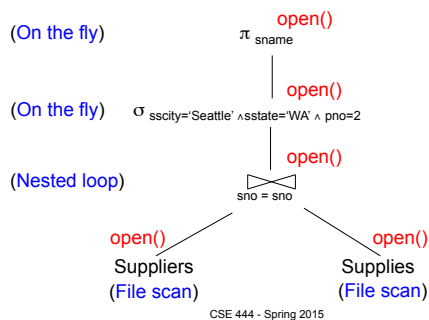
## Iterator Interface

- **open()**
  - Initializes operator state
  - Sets parameters such as selection condition
- **next()**
  - Operator invokes `get_next()` recursively on its inputs
  - Performs processing and produces an output tuple
- **close()**: clean-up state

CSE 444 - Spring 2015

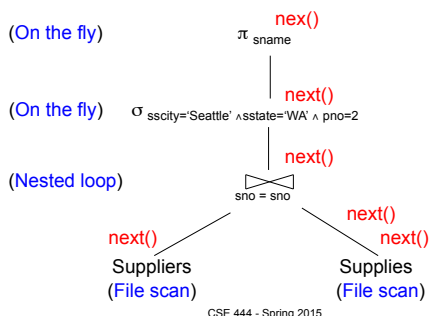
9

## Pipelined Query Execution



10

## Pipelined Query Execution



11

## Pipelined Execution

- Tuples generated by an operator are immediately sent to the parent
- Benefits:
  - No operator synchronization issues
  - Saves cost of writing intermediate data to disk
  - Saves cost of reading intermediate data from disk
- This approach is used whenever possible

CSE 444 - Spring 2015

12

## Intermediate Tuple Materialization

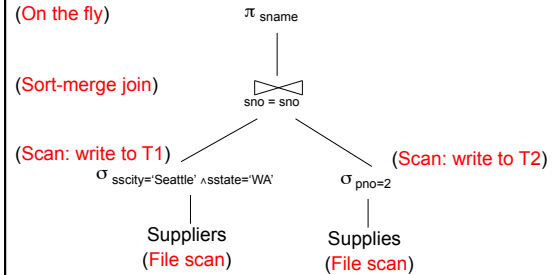
- Tuples generated by an operator are written to disk and in intermediate table

- No direct benefit
- Necessary:
  - For certain operator implementations
  - When we don't have enough memory

CSE 444 - Spring 2015

13

## Intermediate Tuple Materialization



CSE 444 - Spring 2015

14

## Memory Management

Each operator:

- Pre-allocates heap space for tuples
  - Pointers to base data in buffer pool
  - Or new tuples on the heap
- Allocates memory for its internal state
  - Either on heap or buffer pool (depends on system)

DMBS may **limit** how much memory each operator, or each query can use

CSE 444 - Spring 2015

15

## 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
- Execution of the physical plan is pull-based
- Operators given a limited amount of memory

CSE 444 - Spring 2015

16

## Operator Algorithms

CSE 444 - Spring 2015

17

## Operator Algorithms

Design criteria

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

CSE 444 - Spring 2015

18

## 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  $< T(R)$

CSE 444 - Spring 2015

19

## Convention

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

CSE 444 - Spring 2015

20

## Example: Cost of Scanning a Table

- Result may be unsorted:  $B(R)$
- Result needs to be sorted:  $3B(R)$ 
  - We will discuss sorting later

CSE 444 - Spring 2015

21

## 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: read the book

CSE 444 - Spring 2015

22

## Join Algorithms

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

CSE 444 - Spring 2015

23

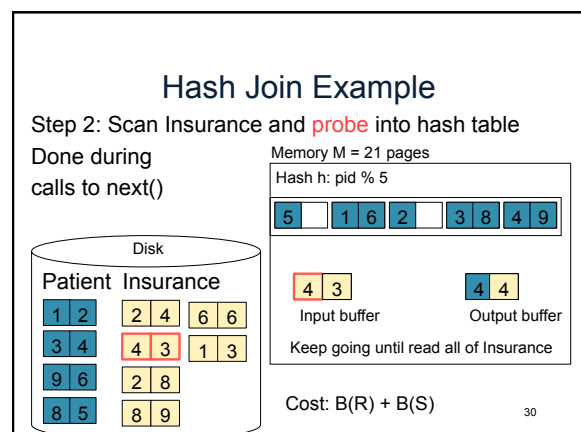
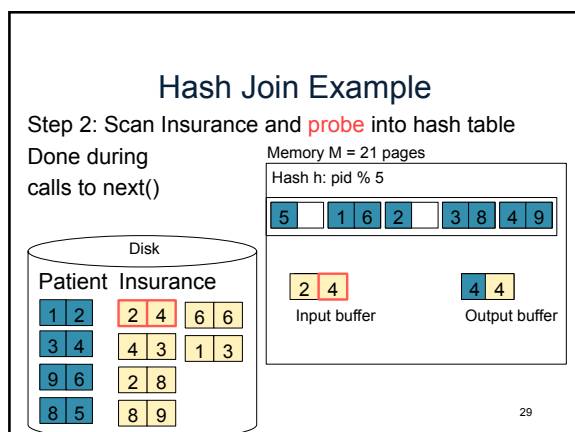
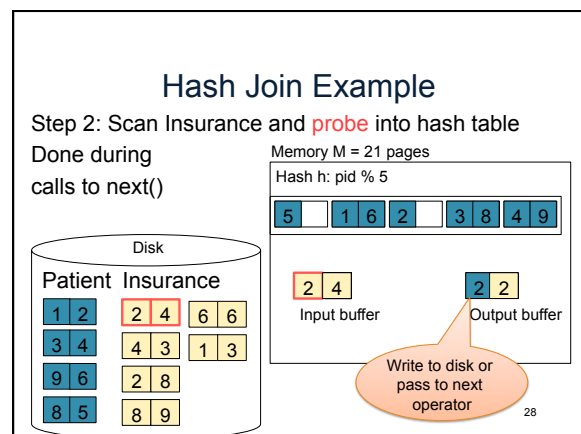
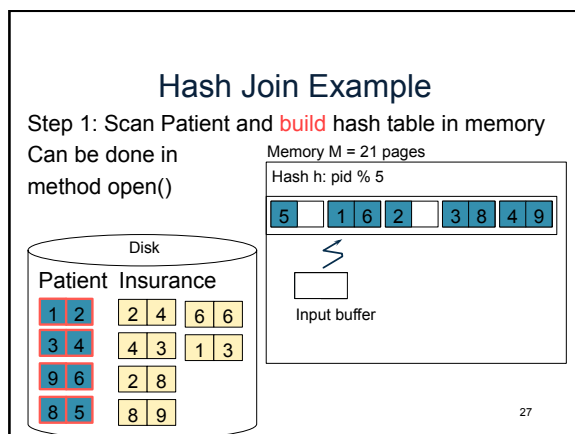
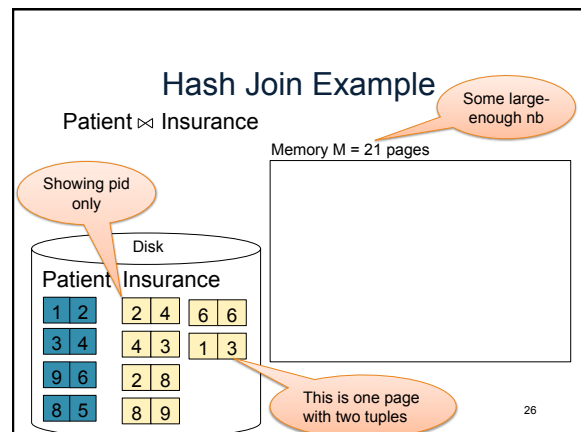
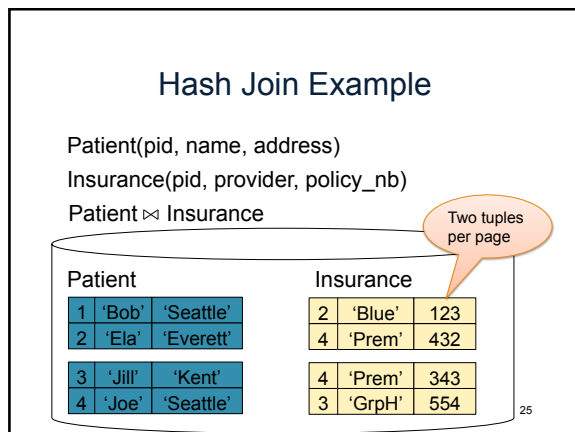
## Hash Join

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  $B(R) \leq M$

CSE 444 - Spring 2015

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?

CSE 444 - Spring 2015

31

## 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

CSE 444 - Spring 2015

32

## 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?

CSE 444 - Spring 2015

33

## 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  $(t_1, t_2)$ 
```

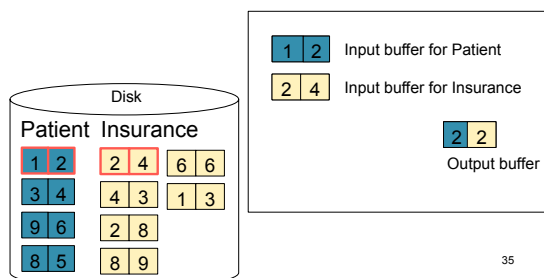
What is the Cost?

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

CSE 444 - Spring 2015

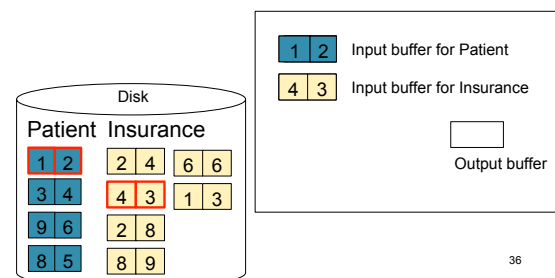
34

## Page-at-a-time Refinement



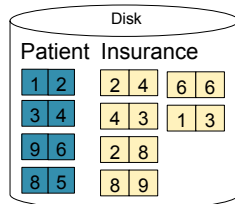
35

## Page-at-a-time Refinement



36

## Page-at-a-time Refinement



1 2 Input buffer for Patient  
 2 8 Input buffer for Insurance  
 Keep going until read all of Insurance  
 Then repeat for next page of Patient... until end of Patient  
 2 2 Output buffer

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

37

## 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)
  
```

What is the Cost?

CSE 444 - Spring 2015

38

## 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 444 - Spring 2015

39

## 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

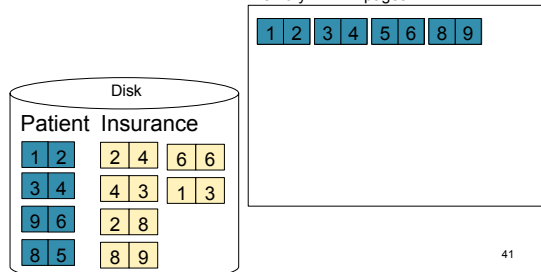
CSE 444 - Spring 2015

40

## Sort-Merge Join Example

Step 1: Scan Patient and sort in memory

Memory M = 21 pages

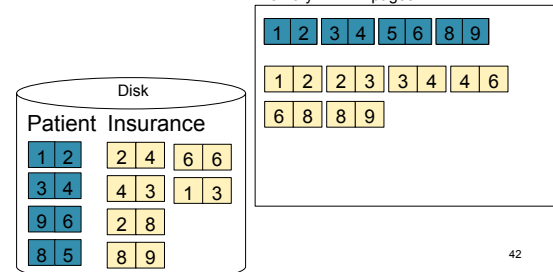


41

## Sort-Merge Join Example

Step 2: Scan Insurance and sort in memory

Memory M = 21 pages



42

