# CSE 444: Database Internals

Lecture 7
Query Execution and
Operator Algorithms (part 1)

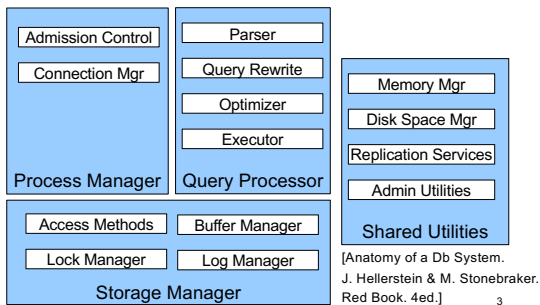CSE 444 - Winter 2018                1

---

## 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 - Winter 2018                2

---

## DBMS Architecture

| Process Manager | Query Processor | Shared Utilities |
|---|---|---|
| Admission Control | Parser | |
| Connection Mgr | Query Rewrite | Memory Mgr |
| | Optimizer | Disk Space Mgr |
| | Executor | Replication Services |
| | | Admin Utilities |

**Storage Manager**
Access Methods | Buffer Manager
Lock Manager | Log Manager

[Anatomy of a Db System.
J. Hellerstein & M. Stonebraker.
Red Book. 4ed.]          3

---

## 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
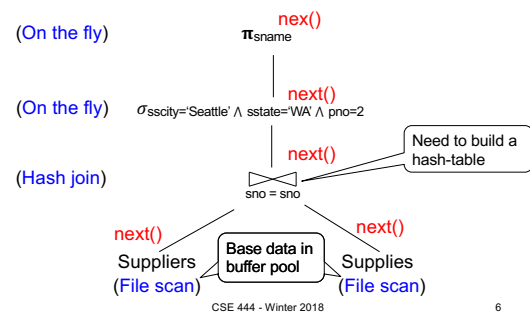
CSE 444 - Winter 2018                4

---

## 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*

CSE 444 - Winter 2018                5

---

## Pipelined Query Execution

(On the fly)                nex()
                            $\pi_{sname}$

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

(Hash join)                 next()          Need to build a
                            ⋈               hash-table
                          sno = sno

        next()                              next()
      Suppliers      Base data in      Supplies
      (File scan)    buffer pool       (File scan)

CSE 444 - Winter 2018                6

---

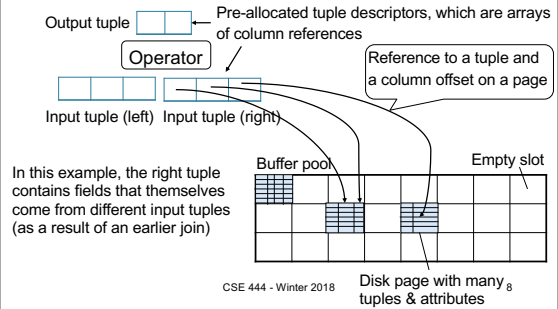## Memory Management

Each operator:
- Pre-allocates heap space for input/output tuples
  - Option 1: Array of pointers to base data in buffer pool
  - Option 2: New tuples on the heap
- Allocates memory for its internal state
  - Either on heap or in buffer pool (depends on system)

DMBS **limits** how much memory each operator, or each query can use
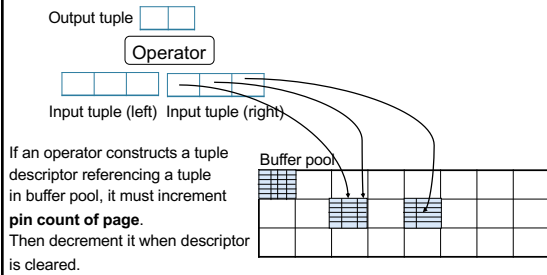
---

## In Flight Tuples (option 1)

Output tuple
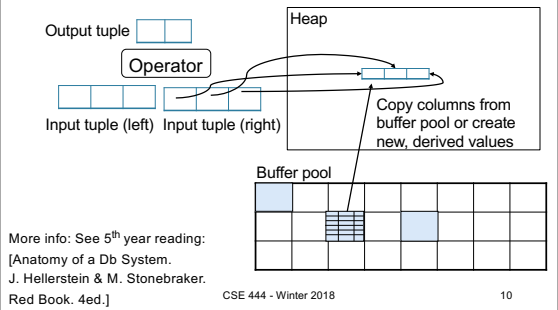
Pre-allocated tuple descriptors, which are arrays of column references

Operator

Reference to a tuple and a column offset on a page

Input tuple (left)   Input tuple (right)

In this example, the right tuple contains fields that themselves come from different input tuples (as a result of an earlier join)

Buffer pool

Empty slot

Disk page with many $_8$ tuples & attributes

---

## In Flight Tuples (option 1)

Output tuple

Operator

Input tuple (left)   Input tuple (right)

If an operator constructs a tuple descriptor referencing a tuple in buffer pool, it must increment **pin count of page**. Then decrement it when descriptor is cleared.

Buffer pool

---

## In Flight Tuples (option 2)

Output tuple

Heap

Operator

Input tuple (left)   Input tuple (right)

Copy columns from buffer pool or create new, derived values

Buffer pool

More info: See 5[th] year reading:
[Anatomy of a Db System.
J. Hellerstein & M. Stonebraker.
Red Book. 4ed.]

---

## Operator Algorithms
## (Quick review from 344 today
## & new algorithms next time)

---

## Operator Algorithms

Design criteria

- Cost: IO, CPU, Network

- Memory utilization

- Load balance (for parallel operators)

## 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)**

---

## 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

---

## Outline

- **Join operator algorithms**

  Review
  - One-pass algorithms (Sec. 15.2 and 15.3)
  - Index-based algorithms (Sec 15.6)

  New
  - 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

---

## Join Algorithms

- Hash join

- Nested loop join

- Sort-merge join

---

## Hash Join

Hash join:  R ⋈ S
- Scan R, build buckets in main memory
- Then scan S and join
- Cost: B(R) + B(S)

- One-pass algorithm when B(R) ≤ M

---

## Hash Join Example

Patient(pid, name, address)
Insurance(pid, provider, policy_nb)
Patient ⋈ 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' | | 3 | 'GrpH' | 554 |

18

# Hash Join Example

Patient ⋈ Insurance

Some large-enough nb

Showing pid only

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 |

This is one page with two tuples

19

---

# 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

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 |

20

---

# 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

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

Disk

Patient  Insurance

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

2 | 4    Input buffer

2 | 2    Output buffer

Write to disk or pass to next operator

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

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

Disk

Patient  Insurance

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

2 | 4    Input buffer

4 | 4    Output buffer

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

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

Disk

Patient  Insurance

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

4 | 3    Input buffer

4 | 4    Output buffer

Keep going until read all of Insurance

Cost: B(R) + B(S)

23

---

# Nested Loop Joins

- Tuple-based nested loop R ⋈ 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?

## Nested Loop Joins

- Tuple-based nested loop R ⋈ 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 ($t_1$, $t_2$)

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

What is the Cost?

CSE 444 - Winter 2018     25

---

## 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 - Winter 2018     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
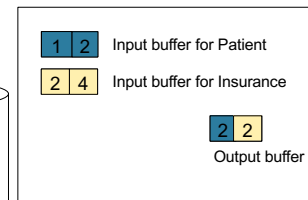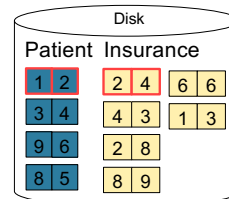>       if $t_1$ and $t_2$ join then output ($t_1$, $t_2$)

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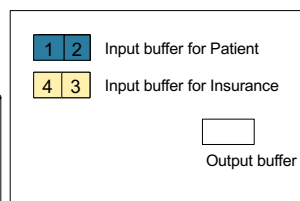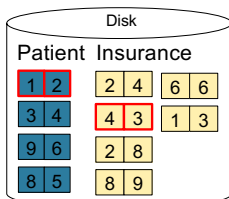
What is the Cost?

CSE 444 - Winter 2018     27

---

## Page-at-a-time Refinement



Disk

Patient  Insurance

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

1 2   Input buffer for Patient
2 4   Input buffer for Insurance

2 2
Output buffer

28

---

## Page-at-a-time Refinement

Disk

Patient  Insurance

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

1 2   Input buffer for Patient
4 3   Input buffer for Insurance

Output buffer

29

---

## Page-at-a-time Refinement

Disk

Patient  Insurance

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

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
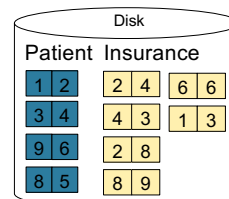
2 2
Output buffer

Cost: $B(R) + B(R)B(S)$    30

## 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 $t_1$ 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 - Winter 2018     31

---

## 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 $t_1$ in r, $t_2$ in s
         if $t_1$ and $t_2$ join then output $(t_1,t_2)$

- Cost: $B(R) + B(R)B(S)/(M-1)$

What is the Cost?

CSE 444 - Winter 2018     32
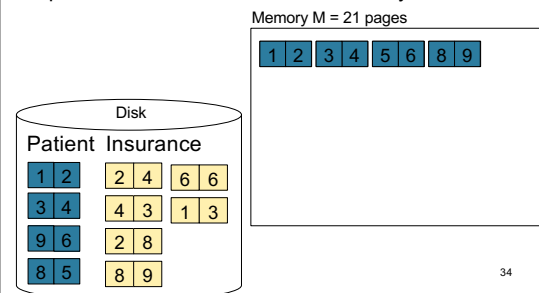
---

## 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 $B(S) + B(R) <= M$
- Typically, this is NOT a one pass algorithm

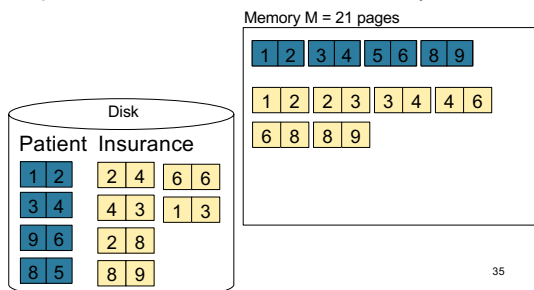CSE 444 - Winter 2018     33

---

## Sort-Merge Join Example

Step 1: Scan Patient and sort in memory



34

---

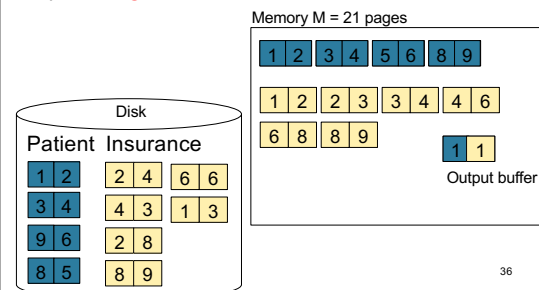## Sort-Merge Join Example

Step 2: Scan Insurance and sort in memory



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

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

| 1 | 2 | 2 | 3 | 3 | 4 | 4 | 6 |

| 6 | 8 | 8 | 9 |

| 2 | 2 |

Output buffer

Keep going until end of first relation

Disk

Patient  Insurance

| 1 | 2 |  | 2 | 4 |  | 6 | 6 |

| 3 | 4 |  | 4 | 3 |  | 1 | 3 |

| 9 | 6 |  | 2 | 8 |

| 8 | 5 |  | 8 | 9 |

37