# Chapter 12: Query Processing

**Database System Concepts, 6th Ed.**

Tuesday, April 2, 2013
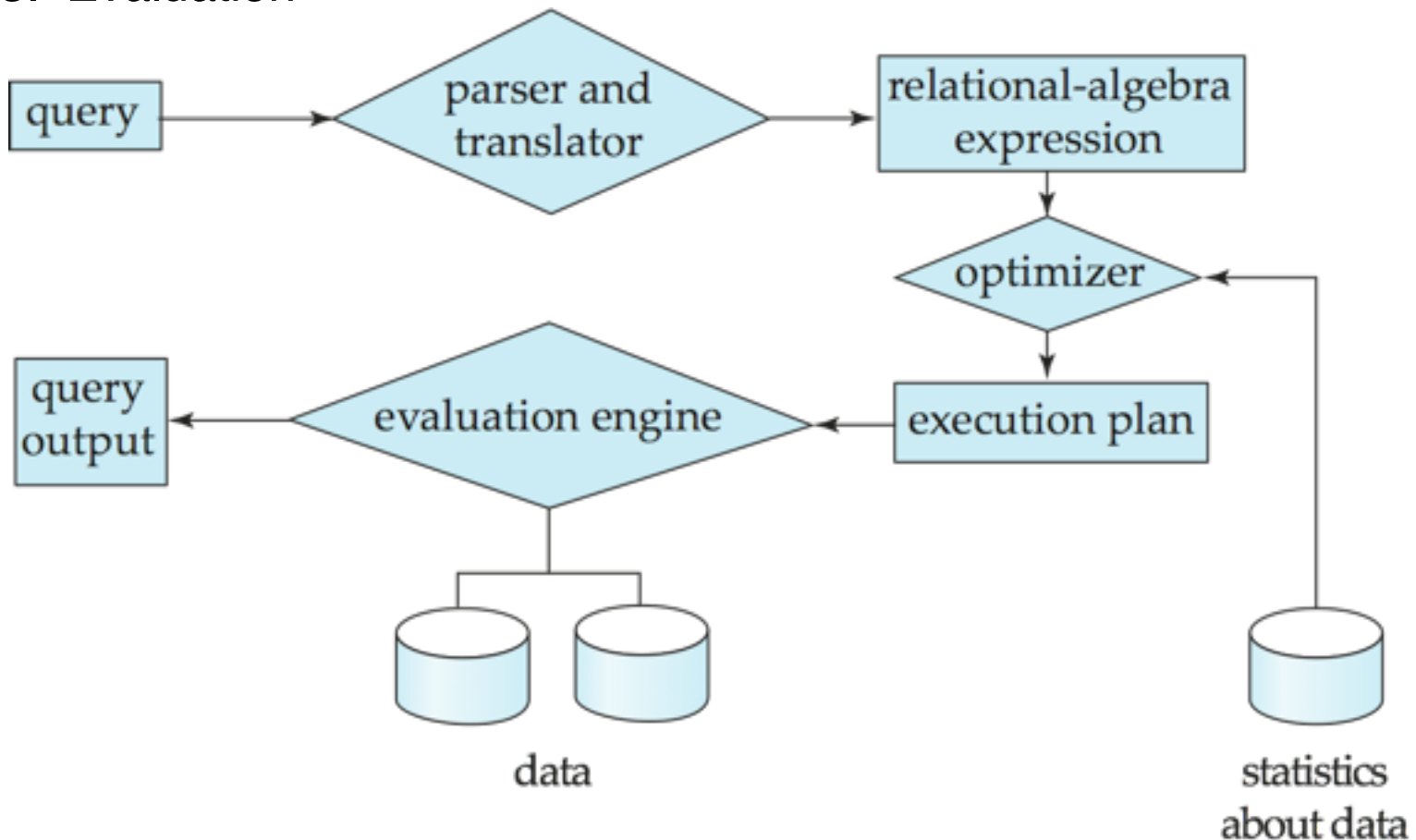
# Chapter 12: Query Processing

■Overview

- ■ Measures of Query Cost

- ■ Selection Operation

- ■ Sorting

- ■ Join Operation

- ■ Other Operations

- ■ Evaluation of Expressions

# Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation

# Query Optimization

■ Amongst all equivalent evaluation plans choose the one with lowest cost (Chap 14).

■ In this chapter we study

- How to measure query costs

- Algorithms for evaluating relational algebra operations

- How to combine algorithms for individual operations in order to evaluate a complete expression

# Chapter 12:  Query Processing

- Overview

- # Measures of Query Cost

- Selection Operation

- Sorting

- Join Operation

- Other Operations

- Evaluation of Expressions

# Measures of Query Cost

- Cost is generally measured as **total elapsed time** for answering query
  - Many factors contribute to time cost
    - *disk accesses, CPU*, or even network *communication*
- Typically disk access is the predominant cost, and is also relatively easy to estimate.   Measured by taking into account
  - **Number of seeks**          * average-seek-cost
  - **Number of blocks read**     * average-block-read-cost
  - **Number of blocks written** * average-block-write-cost
    - Cost to write a block is greater than cost to read a block
      - data is read back after being written to ensure that the write was successful

# Measures of Query Cost (Cont.)

- For simplicity we just use the **number of block transfers** *from disk and the* **number of seeks** as the cost measures

  - $t_T$ – time to transfer one block

  - $t_S$ – time for one seek

  - Cost for b block transfers plus S seeks
    $$b * t_T + S * t_S$$

- We ignore CPU costs for simplicity

  - Real systems do take CPU cost into account

- We do not include cost to writing output to disk in our cost formulae (why?)

# Chapter 12: Query Processing

- Overview

- Measures of Query Cost

- # Selection Operation

- Sorting

- Join Operation

- Other Operations

- Evaluation of Expressions

# Selection Operation

■ **File scan**

■ Algorithm **A1** (**linear search**).  Scan each file block and test all records to see whether they satisfy the selection condition.

- Cost estimate = $b_r$ block transfers + 1 seek

  ‣ $b_r$ denotes number of blocks containing records from relation $r$

- If selection is on a key attribute, can stop on finding record

  ‣ cost = $(b_r/2)$ block transfers + 1 seek

- Linear search can be applied regardless of

  ‣ selection condition or

  ‣ ordering of records in the file, or

  ‣ availability of indices

# Selections Using Indices

# Selections Using Indices

- **Index scan** – search algorithms that use an index
  - selection condition must be on search-key of index.
- **A2** (**primary index, equality on key**).  Retrieve a single record that satisfies the corresponding equality condition
  - $Cost = (h_i + 1) * (t_T + t_S)$
- **A3** (**primary index, equality on nonkey**) Retrieve multiple records.
  - Records will be on consecutive blocks
    - ‣ Let b = number of blocks containing matching records
  - $Cost = h_i * (t_T + t_S) + t_S + t_T * b$

Tuesday, April 2, 2013

# Selections Using Indices

# Selections Using Indices

- **A4** (**secondary index, equality on nonkey**).
  - Retrieve a single record if the search-key is a candidate key
    - *Cost* = $(h_i + 1) * (t_T + t_S)$
  - Retrieve multiple records if search-key is not a candidate key
    - each of *n* matching records may be on a different block
    - Cost = $(h_i + n) * (t_T + t_S)$
      - Can be very expensive!

# Selections Involving Comparisons

# Selections Involving Comparisons

■ Can implement selections of the form $\sigma_{A \leq V}(r)$ or $\sigma_{A \geq V}(r)$ by using

- a linear file scan,
- or by using indices in the following ways:

# Selections Involving Comparisons

■ Can implement selections of the form $\sigma_{A \le V}(r)$ or $\sigma_{A \ge V}(r)$ by using

   ● a linear file scan,

   ● or by using indices in the following ways:

■ **A5** (**primary index, comparison**). (Relation is sorted on A)

   ‣ For $\sigma_{A \ge V}(r)$ use index to find first tuple $\ge v$ and scan relation sequentially from there

   ‣ For $\sigma_{A \le V}(r)$ just scan relation sequentially till first tuple $> v;$ do not use index

# Selections Involving Comparisons

- Can implement selections of the form $\sigma_{A \leq V}(r)$ or $\sigma_{A \geq V}(r)$ by using

  - a linear file scan,

  - or by using indices in the following ways:

- **A5** (**primary index, comparison**)*. (Relation is sorted on A)

  - For $\sigma_{A \geq V}(r)$ use index to find first tuple $\geq v$ and scan relation sequentially from there

  - For $\sigma_{A \leq V}(r)$ just scan relation sequentially till first tuple $> v;$ do not use index

- **A6** (**secondary index, comparison**).

  - For $\sigma_{A \geq V}(r)$ use index to find first index entry $\geq v$ and scan index sequentially from there, to find pointers to records.

  - For $\sigma_{A \leq V}(r)$ just scan leaf pages of index finding pointers to records, till first entry $> v$

# Implementation of Complex Selections

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta_1 \wedge \theta_2 \wedge \ldots \theta_n}(r)$

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta_1 \wedge \theta_2 \wedge \ldots \theta_n}(r)$

- **A7** (**conjunctive selection using one index**).

  - Select a combination of $\theta_i$ and algorithms A1 through A7 that results in the least cost for $\sigma_{\theta_i}(r)$.

  - Test other conditions on tuple after fetching it into memory buffer.

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta 1 \wedge \theta 2 \wedge \ldots \theta n}(r)$

- **A7** (**conjunctive selection using one index**).
  - Select a combination of $\theta_i$ and algorithms A1 through A7 that results in the least cost for $\sigma_{\theta i}(r)$.
  - Test other conditions on tuple after fetching it into memory buffer.

- **A8** (**conjunctive selection using composite index**).
  - Use appropriate composite (multiple-key) index if available.

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta 1} \wedge {}_{\theta 2} \wedge \ldots {}_{\theta n}(r)$

- **A7** (**conjunctive selection using one index**).
  - Select a combination of $\theta_i$ and algorithms A1 through A7 that results in the least cost for $\sigma_{\theta i}(r)$.
  - Test other conditions on tuple after fetching it into memory buffer.

- **A8** (**conjunctive selection using composite index**).
  - Use appropriate composite (multiple-key) index if available.

- **A9** (**conjunctive selection by intersection of identifiers**).
  - Requires indices with record pointers.
  - Use corresponding index for each condition, and take intersection of all the obtained sets of record pointers.
  - Then fetch records from file
  - If some conditions do not have appropriate indices, apply test in memory.

# Algorithms for Complex Selections

Tuesday, April 2, 2013

# Algorithms for Complex Selections

- **Disjunction:** $\sigma_{\theta_1 \vee \theta_2 \vee \ldots \theta_n}(r)$.

- **A10** (**disjunctive selection by union of identifiers**).

  - Applicable if *all* conditions have available indices.

    ‣ Otherwise use linear scan.

  - Use corresponding index for each condition, and take union of all the obtained sets of record pointers.

  - Then fetch records from file

# Chapter 12: Query Processing

- Overview
- Measures of Query Cost
- Selection Operation
- # Sorting
- Join Operation
- Other Operations
- Evaluation of Expressions

Tuesday, April 2, 2013

# Sorting

- We may build an index on the relation, and then use the index to read the relation in sorted order.  May lead to one disk block access for each tuple.

- For relations that fit in memory, techniques like quicksort can be used.  For relations that don't fit in memory, **external sort-merge** is a good choice.

Tuesday, April 2, 2013

# External Sort-Merge

Let $M$ denote memory size (in pages).

1. **Create sorted runs**.  Let $i$ be 0 initially.

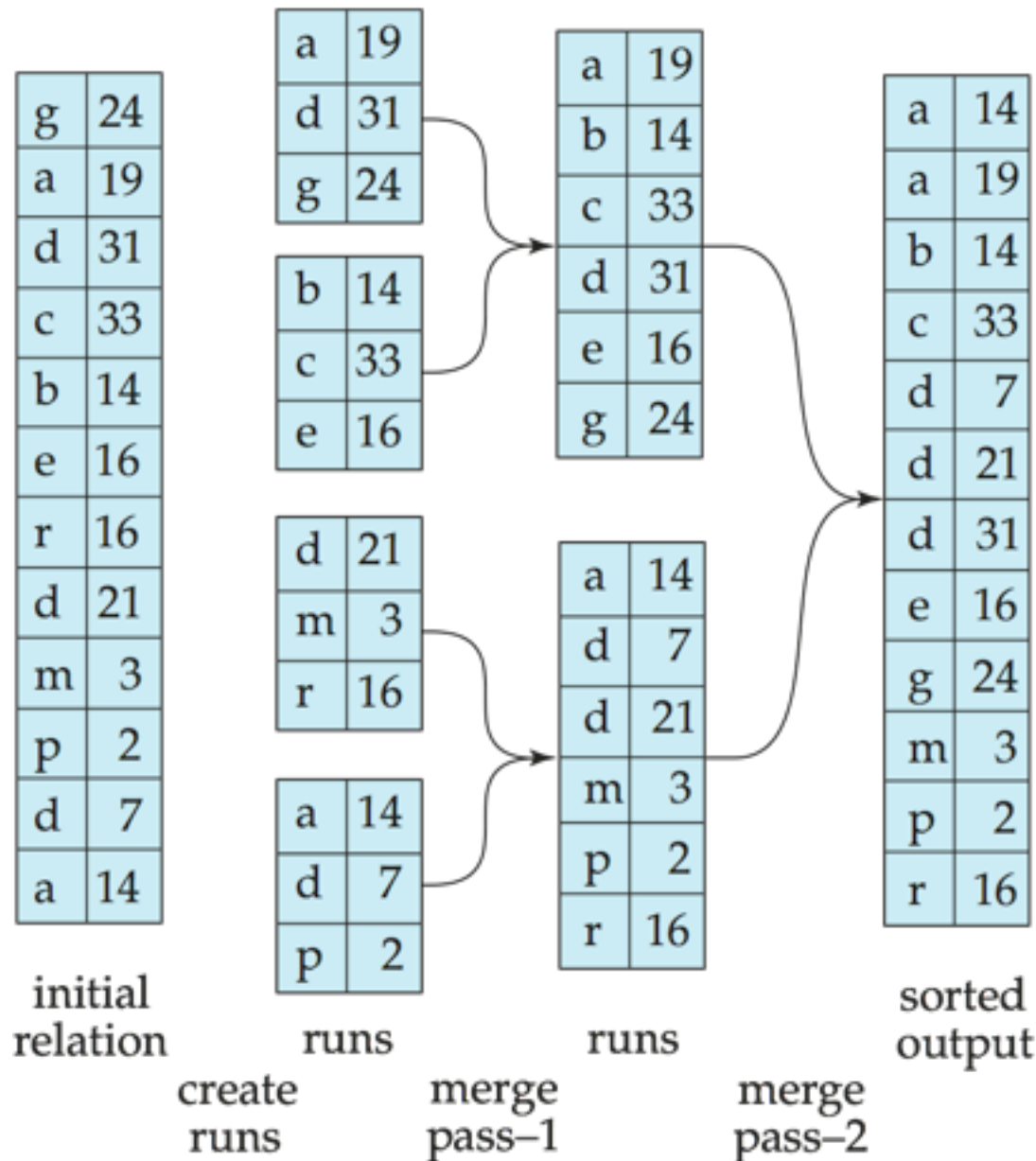   Repeatedly do the following till the end of the relation:
   - (a)  Read $M$ blocks of relation into memory
   - (b)  Sort the in-memory blocks

   - (c)  Write sorted data to run $R_i$; increment $i$.

   Let the final value of $i$ be $N$

2. *Merge the runs (next slide)…..*

| initial relation | runs | runs | sorted output |
|---|---|---|---|

create runs    merge pass–1    merge pass–2

# Chapter 12:  Query Processing

- Overview

- Measures of Query Cost

- Selection Operation

- Sorting

- Join Operation
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join

- Other Operations

- Evaluation of Expressions

# Join Operation

- **Several different algorithms to implement joins**
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- **Choice based on cost estimate**
- **Examples use the following information**
  - Number of records of *student*: 5,000     *takes*: 10,000
  - Number of blocks of    *student*:    100     *takes*:     400

# Nested-Loop Join

■ To compute the theta join $r \bowtie_{\theta} s$

**for each** tuple $t_r$ **in** $r$ **do begin**

  **for each tuple** $t_s$ **in** $s$ **do begin**

   test pair $(t_r, t_s)$ to see if they satisfy the join condition $\theta$

   if they do, add $t_r \cdot t_s$ to the result.

  **end**
**end**

■ $r$ is called the **outer relation** and $s$ the **inner relation** of the join.

■ Requires no indices and can be used with **any kind of join condition (not necessarily equi-join)**.

■ Expensive since it examines every pair of tuples in the two relations.

# Nested-Loop Join (Cont.)

■ In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \qquad \text{seeks}$$

■ If the smaller relation fits entirely in memory, use that as the inner relation.

r [x x] [x x] [x x]

- ● Reduces cost to $b_r + b_s$ block transfers and 2 seeks

■ Assuming worst case memory availability cost estimate is

s [x x] [x x]

- ● with *student* as outer relation:
  - ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,
  - ‣ $5000 + 100 = 5100$ seeks
- ● with *takes* as the outer relation
  - ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

■ If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$
$$n_r + b_r \quad\quad\quad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

r [x x] [x x] [x x] [1]

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

s [x x] [x x]

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

  $n_r * b_s + b_r$   block transfers, plus

  $n_r + b_r$         seeks

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    - $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    - $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    - $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

[1]

r [x x] [x x] [x x]

[2]

s [x x] [x x]

Tuesday, April 2, 2013

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

  $$n_r * b_s + b_r \quad \text{block transfers, plus}$$
  $$n_r + b_r \qquad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    - $5000 * 400 + 100 = 2,000,100$ block transfers,

    - $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    - $10000 * 100 + 400 = 1,000,400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

$$r \overset{1}{[x\ x]} [x\ x] [x\ x]$$

$$s \overset{2}{[x\ x]} \overset{3}{[x\ x]}$$

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \qquad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

1

r [x x] [x x] [x x]

4
2          3

s [x x] [x x]

# Nested-Loop Join (Cont.)

■ In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \quad\quad\quad \text{seeks}$$

■ If the smaller relation fits entirely in memory, use that as the inner relation.

   ● Reduces cost to $b_r + b_s$ block transfers and 2 seeks

■ Assuming worst case memory availability cost estimate is

   ● with *student* as outer relation:

     ‣ 5000 * 400 + 100 = 2,000,100 block transfers,

     ‣ 5000 + 100 = 5100 seeks

   ● with *takes* as the outer relation

     ‣ 10000 * 100 + 400 = 1,000,400 block transfers and 10,400 seeks

■ If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

1

r [x x] [x x] [x x]

4     5
2     3

s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \quad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:
    - ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,
    - ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation
    - ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

$$r \begin{array}{cc} 1 & 6 \\ [x\ x] & [x\ x] \end{array} [x\ x]$$

$$s \begin{array}{cc} 4 & 5 \\ 2 & 3 \\ [x\ x] & [x\ x] \end{array}$$

# Nested-Loop Join (Cont.)

■ In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$
$$n_r + b_r \qquad \text{seeks}$$

■ If the smaller relation fits entirely in memory, use that as the inner relation.

- Reduces cost to $b_r + b_s$ block transfers and 2 seeks

■ Assuming worst case memory availability cost estimate is

- with *student* as outer relation:
  ‣ $5000 * 400 + 100 = 2,000,100$ block transfers,
  ‣ $5000 + 100 = 5100$ seeks

- with *takes* as the outer relation
  ‣ $10000 * 100 + 400 = 1,000,400$ block transfers and 10,400 seeks

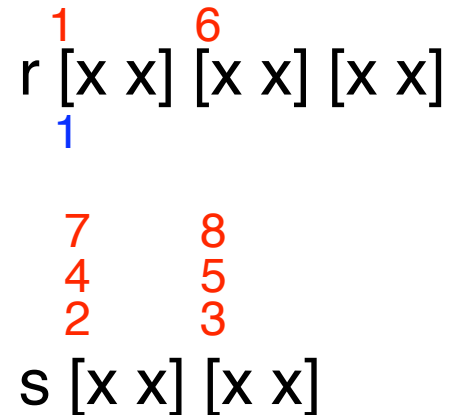■ If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]
1        6

7
4        5
2        3
s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \quad\quad\quad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

1     6

r [x x] [x x] [x x]

7   8
4   5
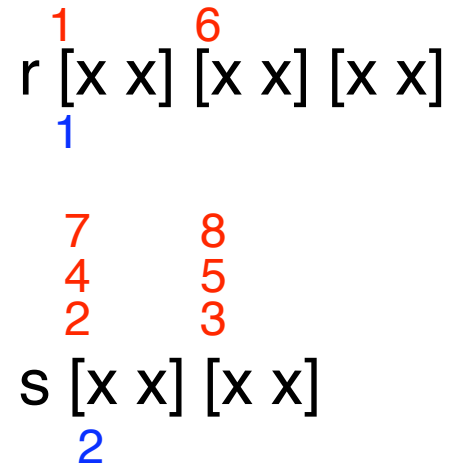2   3

s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \quad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:
    - ‣ 5000 ∗ 400 + 100 = 2,000,100 block transfers,
    - ‣ 5000 + 100 = 5100 seeks

  - with *takes* as the outer relation
    - ‣ 10000 ∗ 100 + 400 = 1,000,400 block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

1          6
r [x x] [x x] [x x]
1

7      8
4      5
2      3
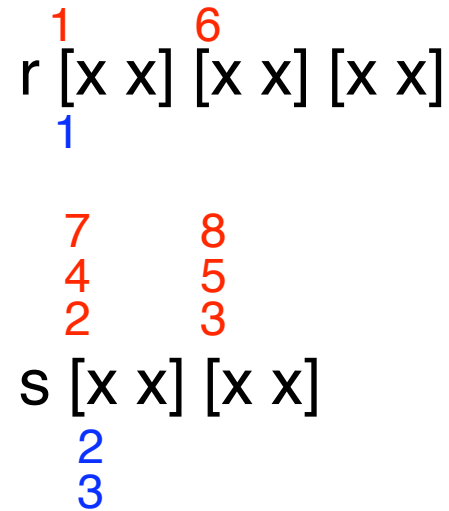s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

  $$n_r * b_s + b_r \quad \text{block transfers, plus}$$
  $$n_r + b_r \qquad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.
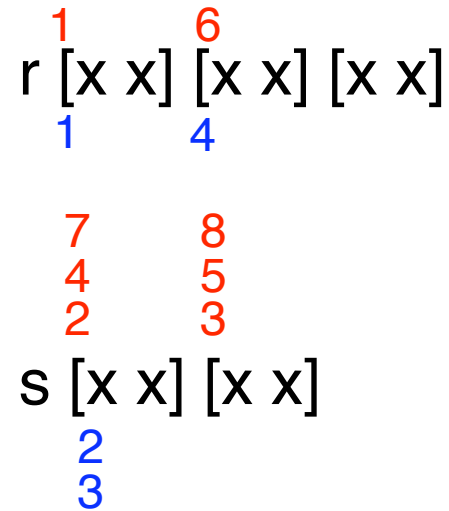
1     6

r [x x] [x x] [x x]

1

7    8
4    5
2    3

s [x x] [x x]

2

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

  $n_r * b_s + b_r$   block transfers, plus

  $n_r + b_r$     seeks

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ 5000 * 400 + 100 = 2,000,100 block transfers,

    ‣ 5000 + 100 = 5100 seeks

  - with *takes*  as the outer relation

    ‣ 10000 * 100 + 400 = 1,000,400 block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

1          6
r [x x] [x x] [x x]
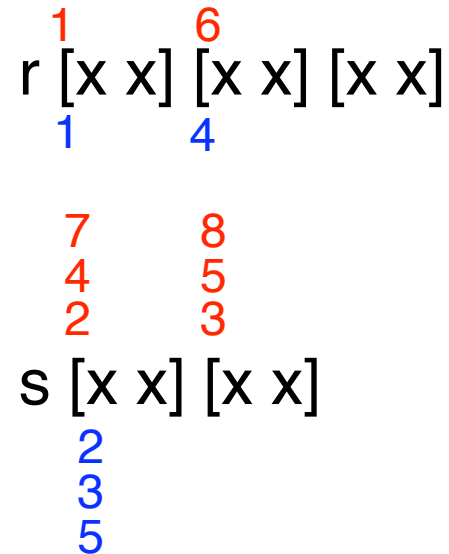1

7          8
4          5
2          3
s [x x] [x x]
2
3

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

    $n_r * b_s + b_r$   block transfers, plus

    $n_r + b_r$          seeks

- If the smaller relation fits entirely in memory, use that as the inner relation.

    - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

    - with *student* as outer relation:

        ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

        ‣ $5000 + 100 = 5100$ seeks

    - with *takes* as the outer relation

        ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.
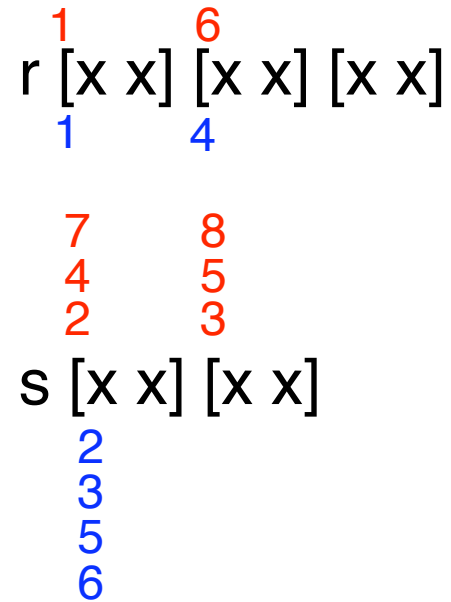
1   6
r [x x] [x x] [x x]
1       4

7       8
4       5
2       3
s [x x] [x x]
  2
  3

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \quad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]

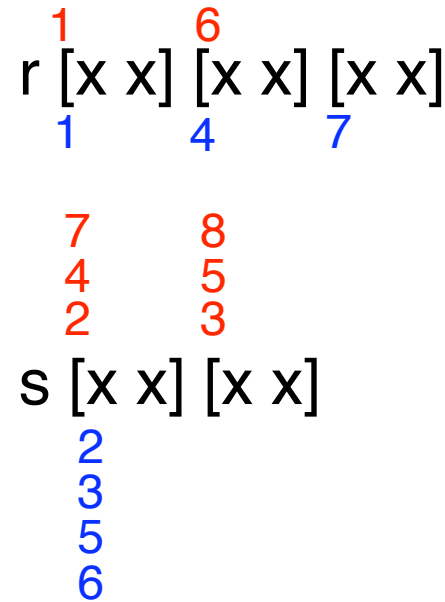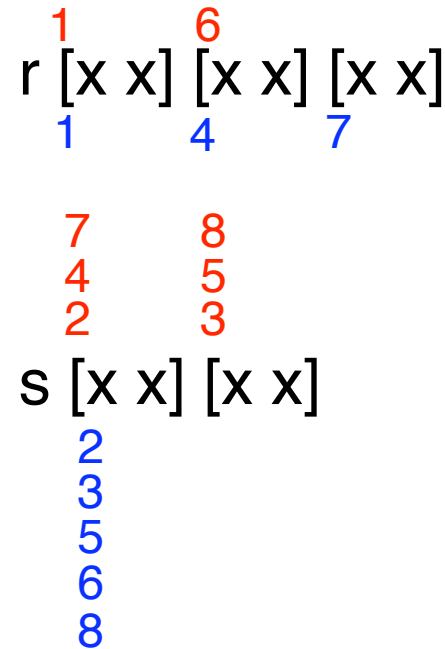s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \qquad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ $5000 * 400 + 100 = 2{,}000{,}100$ block transfers,

    ‣ $5000 + 100 = 5100$ seeks

  - with *takes* as the outer relation

    ‣ $10000 * 100 + 400 = 1{,}000{,}400$ block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]

s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

  $n_r * b_s + b_r$   block transfers, plus

  $n_r + b_r$           seeks

- If the smaller relation fits entirely in memory, use that as the inner relation.

  - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

  - with *student* as outer relation:

    ‣ 5000 * 400 + 100 = 2,000,100 block transfers,

    ‣ 5000 + 100 = 5100 seeks

  - with *takes*  as the outer relation

    ‣ 10000 * 100 + 400 = 1,000,400 block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]

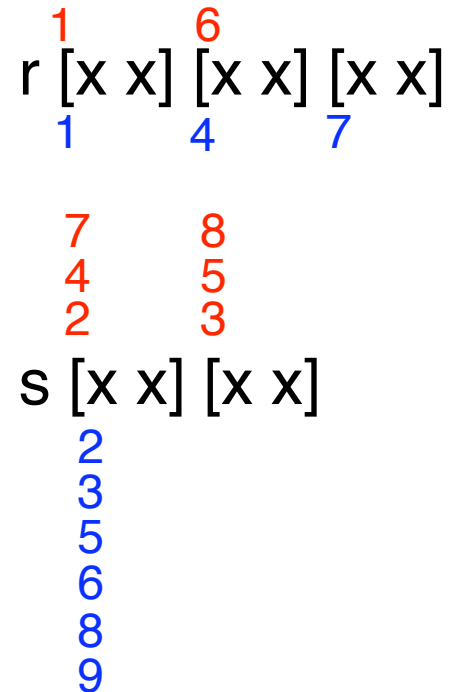s [x x] [x x]

# Nested-Loop Join (Cont.)

- In the worst case, the estimated cost is

    $$n_r * b_s + b_r \quad \text{block transfers, plus}$$

    $$n_r + b_r \qquad \text{seeks}$$

- If the smaller relation fits entirely in memory, use that as the inner relation.

    - Reduces cost to $b_r + b_s$ block transfers and 2 seeks

- Assuming worst case memory availability cost estimate is

    - with *student* as outer relation:

        - 5000 $*$ 400 + 100 = 2,000,100 block transfers,

        - 5000 + 100 = 5100 seeks

    - with *takes* as the outer relation

        - 10000 $*$ 100 + 400 = 1,000,400 block transfers and 10,400 seeks

- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]

1   6
1   4   7

7   8
4   5
2   3

s [x x] [x x]

2
3
5
6
8

# Nested-Loop Join (Cont.)

■ In the worst case, the estimated cost is

$$n_r * b_s + b_r \quad \text{block transfers, plus}$$

$$n_r + b_r \qquad \text{seeks}$$

■ If the smaller relation fits entirely in memory, use that as the inner relation.

● Reduces cost to $b_r + b_s$ block transfers and 2 seeks

■ Assuming worst case memory availability cost estimate is

● with *student* as outer relation:

▸ 5000 * 400 + 100 = 2,000,100 block transfers,

▸ 5000 + 100 = 5100 seeks

● with *takes* as the outer relation

▸ 10000 * 100 + 400 = 1,000,400 block transfers and 10,400 seeks

■ If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.

r [x x] [x x] [x x]

s [x x] [x x]

# Block Nested-Loop Join

■ Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.

**for each** block $B_r$ **of** $r$ **do begin**

    **for each** block $B_s$ **of** *s* **do begin**

        **for each** tuple $t_r$ **in** $B_r$ **do begin**

            **for each** tuple $t_s$ **in** $B_s$ **do begin**

                Check if $(t_r, t_s)$ satisfy the join condition

                if they do, add $t_r \cdot t_s$ to the result.

            **end**

        **end**

    **end**

**end**

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation *s* is read once for each *block* in the outer relation

- Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]

s [x x] [x x]

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

  - Each block in the inner relation $s$ is read once for each *block* in the outer relation

- Best case: $b_r + b_s$ block transfers + 2 seeks.

$$r^1 [x\ x]\ [x\ x]\ [x\ x]$$

$$s\ [x\ x]\ [x\ x]$$

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

  - Each block in the inner relation $s$ is read once for each *block* in the outer relation

- Best case: $b_r + b_s$ block transfers + 2 seeks.

1

r [x x] [x x] [x x]

2

s [x x] [x x]

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation $s$ is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

1

r [x x] [x x] [x x]

2      3

s [x x] [x x]

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation $s$ is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

$$\overset{1}{r} [x\ x]\ \overset{4}{[x\ x]}\ [x\ x]$$

$$\overset{2}{s} [x\ x]\ \overset{3}{[x\ x]}$$

# Block Nested-Loop Join (Cont.)

■ Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

● Each block in the inner relation *s* is read once for each *block* in the outer relation

■ Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
1       4

s [x x] [x x]
5       3
2

Tuesday, April 2, 2013

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
    - Each block in the inner relation *s* is read once for each *block* in the outer relation

- Best case: $b_r + b_s$ block transfers + 2 seeks.

1       4
r [x x] [x x] [x x]

5       6
2       3
s [x x] [x x]

# Block Nested-Loop Join (Cont.)

■ Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

● Each block in the inner relation *s* is read once for each *block* in the outer relation

■ Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]

s [x x] [x x]

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
  - Each block in the inner relation $s$ is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
1        4        7

s [x x] [x x]
  8
  5        6
  2        3

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation $s$ is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

$$
r \begin{array}{ccc} 1 & 4 & 7 \\ [x\ x] & [x\ x] & [x\ x] \end{array}
$$

$$
s \begin{array}{cc} 8 & 9 \\ 5 & 6 \\ 2 & 3 \\ [x\ x] & [x\ x] \end{array}
$$

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation *s* is read once for each *block* in the outer relation
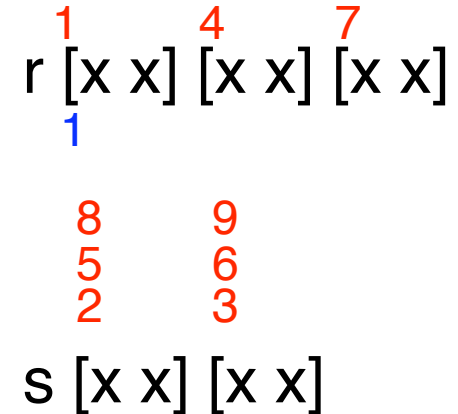- Best case: $b_r + b_s$ block transfers + 2 seeks.

$$r \ [x \ x] \ [x \ x] \ [x \ x]$$

1      4      7

1

8      9
5      6
2      3

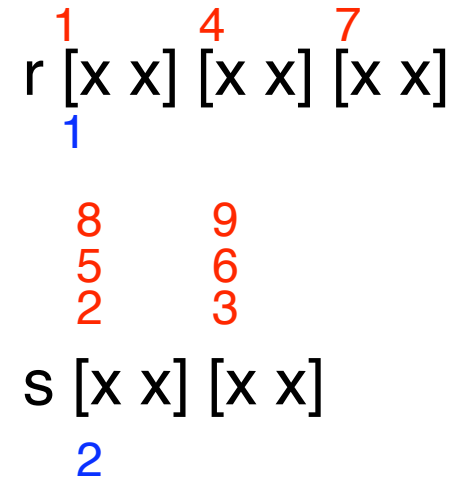$$s \ [x \ x] \ [x \ x]$$

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
  - Each block in the inner relation *s* is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
1      4      7
1

8    9
5    6
2    3

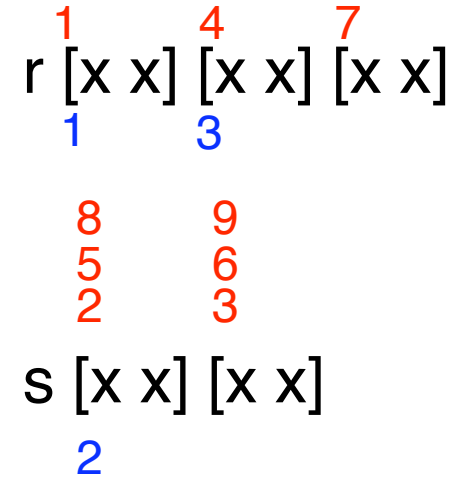s [x x] [x x]
2

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks
  - Each block in the inner relation *s* is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
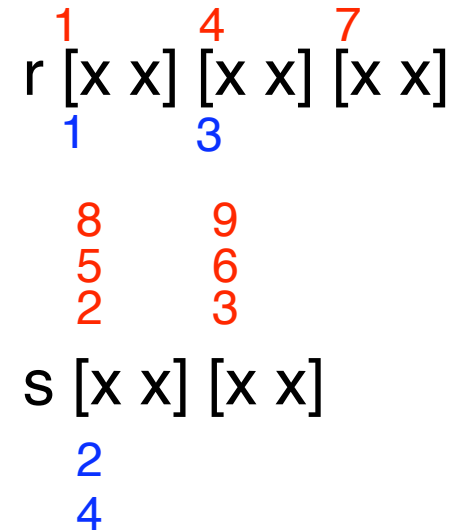
s [x x] [x x]

# Block Nested-Loop Join (Cont.)

■ Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

  ● Each block in the inner relation $s$ is read once for each *block* in the outer relation

■ Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
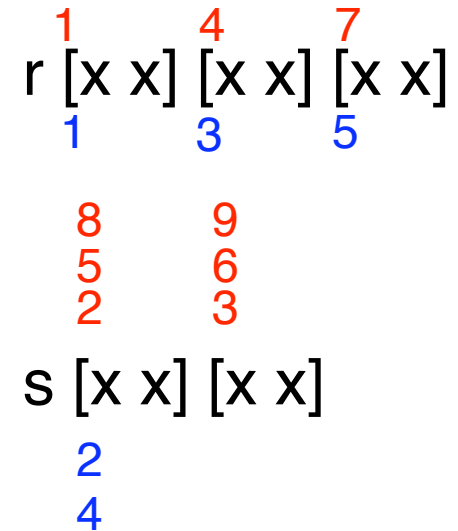
s [x x] [x x]

# Block Nested-Loop Join (Cont.)

■ Worst case estimate: $b_r * b_s + b_r$ block transfers + 2 * $b_r$ seeks

- Each block in the inner relation $s$ is read once for each *block* in the outer relation

■ Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]
1        4        7
   1        3        5

8     9
5     6
2     3

s [x x] [x x]
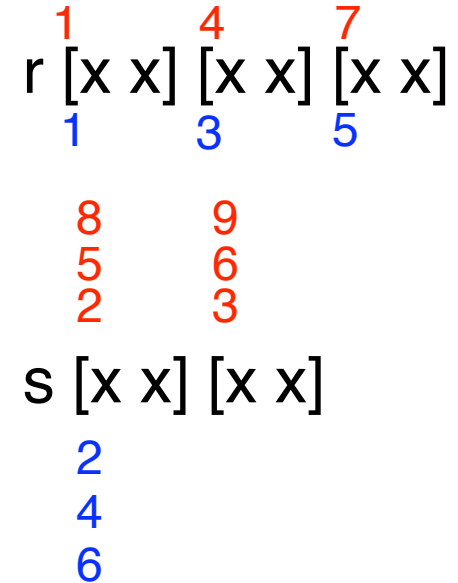   2
   4

# Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
  - Each block in the inner relation *s* is read once for each *block* in the outer relation
- Best case: $b_r + b_s$ block transfers + 2 seeks.

r [x x] [x x] [x x]

<table>
<tr><td>1</td><td>4</td><td>7</td></tr>
<tr><td>1</td><td>3</td><td>5</td></tr>
</table>

<table>
<tr><td>8</td><td>9</td></tr>
<tr><td>5</td><td>6</td></tr>
<tr><td>2</td><td>3</td></tr>
</table>

s [x x] [x x]

2
4
6

# Indexed Nested-Loop Join

r [x x] [x x] [x x]

- **Index lookups can replace file scans if**
  - join is an equi-join or natural join and        s [x x] [x x]
  - an index is available on the inner relation's join attribute
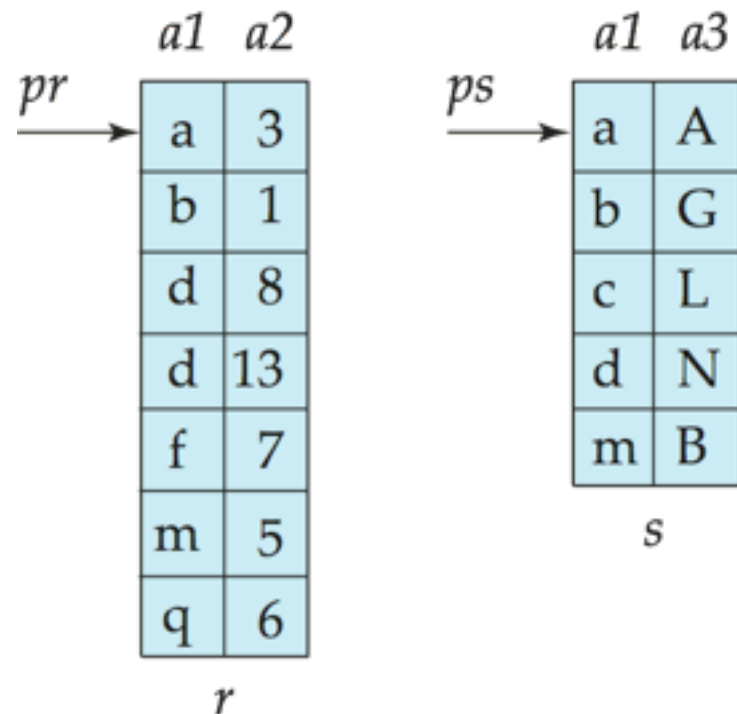    - ‣ Can construct an index just to compute a join.
- **For each tuple $t_r$ in the outer relation $r$, use the index to look up tuples in $s$ that satisfy the join condition with tuple $t_r$.**
- **Worst case:** buffer has space for only one page of $r$, and, for each tuple in $r$, we perform an index lookup on $s$.
- **Cost of the join:** $b_r (t_T + t_S) + n_r * c$
  - Where $c$ is the cost of traversing index and fetching all matching $s$ tuples for one tuple or $r$
  - $c$ can be estimated as cost of a single selection on $s$ using the join condition.
- **If indices are available on join attributes of both $r$ and $s$, use the relation with fewer tuples as the outer relation.**

# Merge-Join

1. Sort both relations on their join attribute (if not already sorted on the join attributes).

2. Merge the sorted relations to join them
   1. Join step is similar to the merge stage of the sort-merge algorithm.
   2. Main difference is handling of duplicate values in join attribute — every pair with same value on join attribute must be matched
   3. Detailed algorithm in book

# Merge-Join (Cont.)

- Can be used for equi-joins and natural joins

- Each block needs to be read only once (assuming all tuples for any given value of the join attributes fit in memory

- Thus the cost of merge join is:

$$b_r + b_s \text{ block transfers } + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$

  - \+ the cost of sorting if relations are unsorted.

  - $b_b$: # of buffer blocks allocated for each relation (how many blocks we read each time)
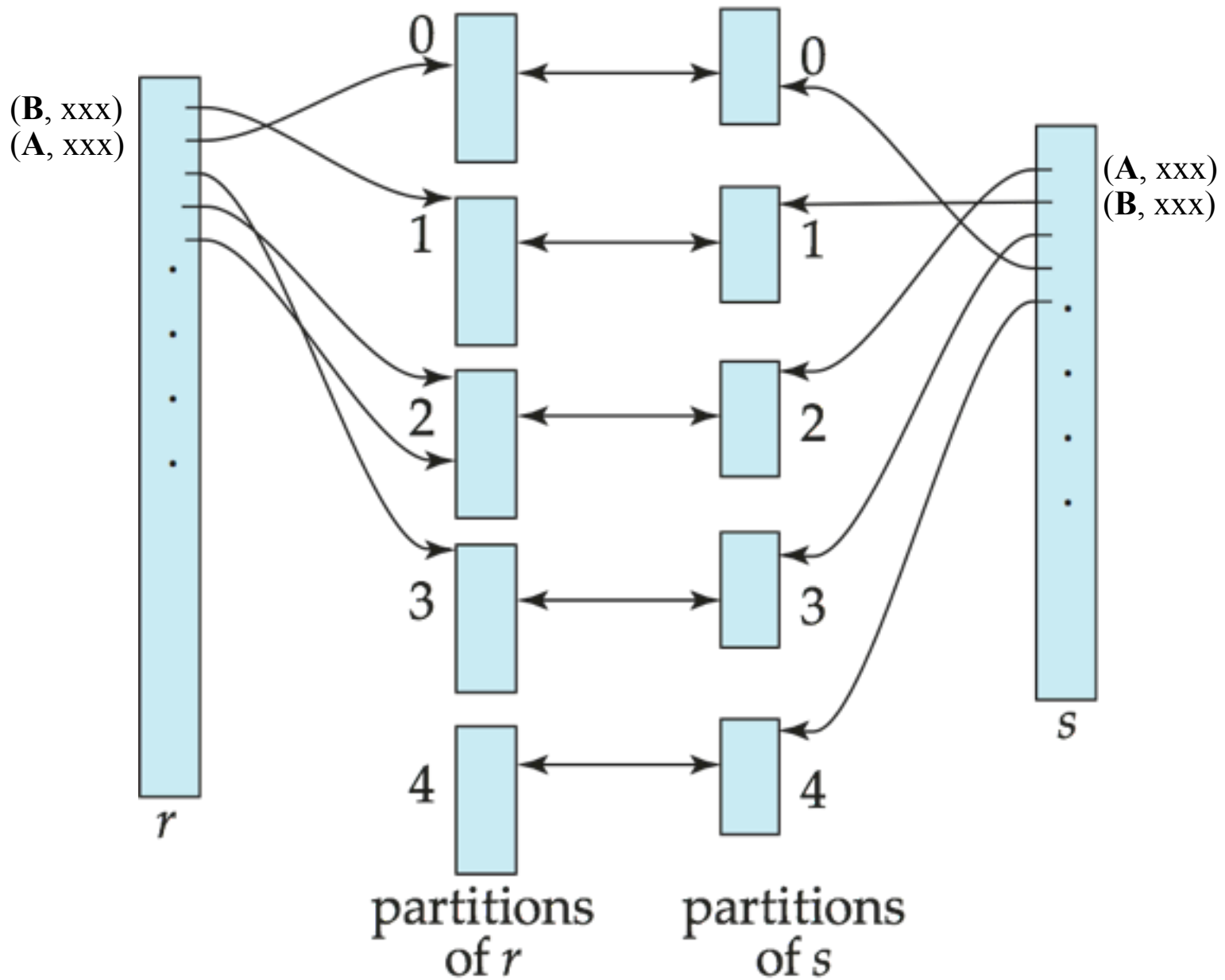
# Hash-Join

r={1, 2, 3, 4, 5, 6, 7}                 r1={1, 2, 3, 4}, r2={5, 6, 7}
s={2, 4, 6, 8}                           s1={2, 4}, s2={6, 8}

- Applicable for equi-joins and natural joins.

- A hash function $h$ is used to partition tuples of both relations

- $h$ maps *JoinAttrs* values to {0, 1, ..., $n$}, where *JoinAttrs* denotes the common attributes of $r$ and $s$ used in the natural join.

  - $r_0, r_1, \ldots, r_n$ denote partitions of $r$ tuples

    ‣ Each tuple $t_r \in r$ is put in partition $r_i$ where $i = h(t_r[JoinAttrs])$.

  - $r_0, r_1 \ldots, r_n$ denotes partitions of $s$ tuples

    ‣ Each tuple $t_s \in s$ is put in partition $s_i$, where $i = h(t_s[JoinAttrs])$.

- *Note:* In book, $r_i$ is denoted as $H_{ri}$, $s_i$ is denoted as $H_{si}$ and

  $n$ is denoted as $n_h$.

# Hash-Join (Cont.)

- *r* tuples in $r_i$ need only to be compared with *s* tuples in $s_i$

  Need not be compared with *s* tuples in any other partition, since:

  - an *r* tuple and an *s* tuple that satisfy the join condition will have the same value for the join attributes.

  - If that value is hashed to some value *i*, the *r* tuple has to be in $r_i$ and the *s* tuple in $s_i$.

Tuesday, April 2, 2013

# Hash-Join Algorithm

The hash-join of $r$ and $s$ is computed as follows.

1. Partition the relation $s$ using hashing function $h$. When partitioning a relation, one block of memory is reserved as the output buffer for each partition.

2. Partition $r$ similarly.

3. For each $i$:

   (a) Load $s_i$ into memory and build an in-memory hash index on it using the join attribute.

   (b) Read the tuples in $r_i$ from the disk one by one. For each tuple $t_r$ locate each matching tuple $t_s$ in $s_i$ using the in-memory hash index. Output the concatenation of their attributes.

Relation $s$ is called the **build input** and $r$ is called the **probe input**.

# Hash-Join algorithm (Cont.)

■ The value *n* and the hash function *h* is chosen such that each $s_i$ should fit in memory.

- Typically n is chosen as $\lceil b_s/M \rceil * f$ where f is a "**fudge factor**", typically around 1.2

- The probe relation partitions $s_i$ need not fit in memory

- $b_s$: # of disk blocks for relation S.

- M: # of memory pages

# Cost of Hash-Join

- Cost of hash join:

$$3(b_r + b_s) + 4 * n_h \text{ block transfers}$$

$$2(\lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil) + 2 * n_h \text{ seeks}$$

  - partitioning
    - read: $b\_r + b\_s$ blocks
    - write: $b\_r + b\_s$ blocks + 2 $n_h$ blocks (last block of each partition)
  - matching
    - build: $b\_s + n_h$ blocks
    - probe: $b\_r + n_h$ blocks

- If the entire build input can be kept in main memory no partitioning is required
  - Cost estimate goes down to $b_r + b_s$.

# Chapter 12:  Query Processing

- Overview

- Measures of Query Cost

- Selection Operation

- Sorting

- Join Operation

# Other Operations

- duplicate elimination / projection

- aggregation / set operations

- outer join

- Evaluation of Expressions

# Other Operations

- **Duplicate elimination** can be implemented via hashing or sorting.

  - On sorting duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.

  - *Optimization:* duplicates can be deleted during run generation as well as at intermediate merge steps in external sort-merge.

  - Hashing is similar – duplicates will come into the same bucket.

- **Projection:**

  - perform projection on each tuple

  - followed by duplicate elimination.

# Other Operations : Aggregation

Tuesday, April 2, 2013

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.
  - Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.

  - Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.

  - *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values

# Other Operations : Aggregation

■ **Aggregation** can be implemented in a manner similar to duplicate elimination.

- Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.

- *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values

  ‣ For count, min, max, sum: keep aggregate values on tuples found so far in the group.

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.

  - Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.

  - *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values

    ▸ For count, min, max, sum: keep aggregate values on tuples found so far in the group.

      – When combining partial aggregate for count, add up the aggregates

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.
  - Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.
  - *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values
    - ‣ For count, min, max, sum: keep aggregate values on tuples found so far in the group.
      - – When combining partial aggregate for count, add up the aggregates
    - ‣ For avg, keep sum and count, and divide sum by count at the end

# Other Operations : Set Operations

- **Set operations** ($\cup$, $\cap$ and —): can either use variant of merge-join after sorting, or variant of hash-join.

- E.g., Set operations using hashing:
  1. Partition both relations using the same hash function
  2. Process each partition $i$ as follows.
     1. Using a different hashing function, build an in-memory hash index on $r_i$.
     2. Process $s_i$ as follows
        - $r \cup s$:
          1. Add tuples in $s_i$ to the hash index if they are not already in it.
          2. At end of $s_i$ add the tuples in the hash index to the result.

# Other Operations : Set Operations

■ E.g., Set operations using hashing:

1. as before partition $r$ and $s,$

2. as before, process each partition $i$ as follows

   1. build a hash index on $r_i$

   2. Process $s_i$ as follows

      ● $r \cap s$:

         1. output tuples in $s_i$ to the result if they are already there in the hash index

      ● $r - s:$

         1. for each tuple in $s_i$, if it is there in the hash index, delete it from the index.

         2. At end of $s_i$ add remaining tuples in the hash index to the result.

# Chapter 12:  Query Processing

- Overview

- Measures of Query Cost

- Selection Operation

- Sorting

- Join Operation

- Other Operations

- # Evaluation of Expressions

  - materialization

  - pipelining

    - pull-based (demand-driven; lazy)

    - push-based (produce-driven; eager)

# Evaluation of Expressions

- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
  - **Materialization**:  generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk.  Repeat.
  - **Pipelining**:  pass on tuples to parent operations even as an operation is being executed
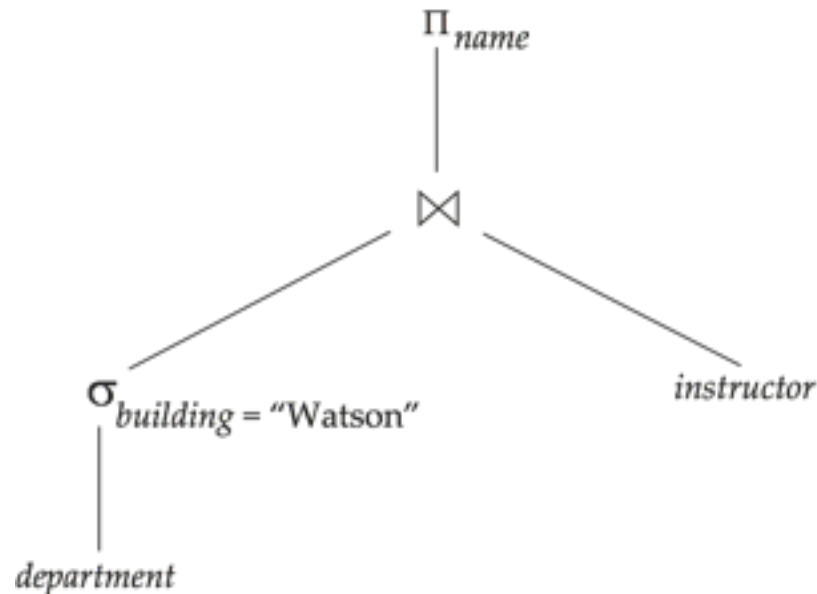- We study above alternatives in more detail

# Materialization

- **Materialized evaluation**: evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.

- E.g., in figure below, compute and store

$$\sigma_{building="Watson"}(department)$$

then compute the store its join with *instructor,* and finally compute the projection on *name.*

# Materialization (Cont.)

- Materialized evaluation is always applicable

- Cost of writing results to disk and reading them back can be quite high

- **Double buffering**: use two output buffers for each operation, when one is full write it to disk while the other is getting filled
  - Allows overlap of disk writes with computation and reduces execution time

# Pipelining

- **Pipelined evaluation** :  evaluate several operations simultaneously, passing the results of one operation on to the next.
- E.g., in previous expression tree, don't store result of

$$\sigma_{building="Watson"}(department)$$

  - instead, pass tuples directly to the join..  Similarly, don't store result of join, pass tuples directly to projection.
- Much cheaper than materialization: no need to store a temporary relation to disk.
- Pipelining may not always be possible – e.g., sort, hash-join.
- For pipelining to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.
- Pipelines can be executed in two ways:  **demand driven** and **producer driven**

# Pipelining (Cont.)

- In **demand driven** (or **lazy** or **pull-based)** evaluation

  - system repeatedly requests next tuple from top level operation

  - Each operation requests next tuple from children operations as required, in order to output its next tuple

- In **producer-driven** (or **eager** or **push-based)** pipelining

  - Operators produce tuples eagerly and pass them up to their parents

    ‣ Buffer maintained between operators, child puts tuples in buffer, parent removes tuples from buffer

    ‣ if buffer is full, child waits till there is space in the buffer, and then generates more tuples

  - System schedules operations that have space in output buffer and can process more input tuples

# Pipelining (Cont.)

- Implementation of demand-driven pipelining
  - Each operation is implemented as an **iterator** implementing the following operations
    - open()
      - E.g. file scan: initialize file scan
        - » state: pointer to beginning of file
      - E.g.merge join: sort relations;
        - » state: pointers to beginning of sorted relations
    - next()
      - E.g. for file scan: Output next tuple, and advance and store file pointer
      - E.g. for merge join: continue with merge from earlier state till
        next output tuple is found. Save pointers as iterator state.
    - close()

# Chapter 12: Query Processing

- Overview
- Measures of Query Cost
- Selection Operation
- Sorting
- Join Operation
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Other Operations
  - duplicate elimination / projection
  - aggregation / set operations
  - outer join
- Evaluation of Expressions
  - materialization
  - pipelining
    ‣ pull-based (demand-driven; lazy)
    ‣ push-based (produce-driven; eager)

Tuesday, April 2, 2013

# End of Chapter

**Database System Concepts, 6th Ed.**

©Silberschatz, Korth and Sudarshan
See www.db-book.com for conditions on re-use

47

Tuesday, April 2, 2013