

#### Introduction

Alternative ways of evaluating a given query

- Equivalent expressions
- Different algorithms for each operation





# Introduction (Cont.)

An evaluation plan defines exactly what algorithm is used for each operation, and how the execution of the operations is coordinated.



Find out how to view query execution plans on your favorite database



# Introduction (Cont.)

- Cost difference between evaluation plans for a query can be enormous
  - E.g. seconds vs. days in some cases
- Steps in cost-based query optimization
  - Generate logically equivalent expressions using equivalence rules
  - 2. Annotate resultant expressions to get alternative query plans
  - 3. Choose the cheapest plan based on estimated cost
- Estimation of plan cost based on:
  - Statistical information about relations. Examples:
    - number of tuples, number of distinct values for an attribute
  - Statistics estimation for intermediate results
    - to compute cost of complex expressions
  - Cost formulae for algorithms, computed using statistics



## **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

- $t_{\tau}$  time to transfer one block
- $t_s$  time for one seek
- Cost for b block transfers plus Seeks  $b * t_{\tau} + S * t_{s}$
- 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



# **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
  - Number of blocks of student: 100

So rows ber blick





#### **Nested-Loop Join**

To compute the theta join  $r \Join_{\theta} s$ 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.  $\int Nsr would n$ end adulation

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



# **Nested-Loop Join (Cont.)**

- In the worst case, if there is enough memory only to hold one block of each relation, the estimated cost is <u>seeks</u> in <u>relation</u>, the estimated cost is <u>seeks</u> in <u>relation</u>, the seeks is <u>seeks</u> in <u>relation</u>, the seeks <u>nrelation</u>, the seeks <u>nrelation</u> <u>nrelation</u>, the seeks <u>nrelation</u> <u>nrelation</u> <u>nrelation</u>.
  Reduces cost to <u>brever</u> availability cost estimate is <u>nrelation</u> <u>nrelation</u>.
  - 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.
  - Block nested-loops algorithm (next slide) is preferable.



# **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<sub>r</sub> of r do begin
                                                        br + bs + br
   for each block B<sub>s</sub> of s do begin
        for each tuple t_r in B_r do begin
                                                        n, *bs + br
           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.
- Improvements to nested loop and block nested loop algorithms:
  - In block nested-loop, use M 2 disk blocks as blocking unit for outer relations, where M = memory size in blocks; use remaining two blocks to buffer inner relation and output

Cost = 
$$[b_r / (M-2)] * b_s + b_r$$
 block transfers +  $2[b_r / (M-2)]$  seeks

- If equi-join attribute forms a key or inner relation, stop inner loop on first match
- Scan inner loop forward and backward alternately, to make use of the blocks remaining in buffer (with LRU replacement)
- Use index on inner relation if available (next slide)



#### **Indexed Nested-Loop Join**

Index lookups can replace file scans if

- join is an equi-join or natural join and
- 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:  $\vec{b}_r(t_T + t_S) + n_r * \vec{c}$ 
  - 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.



## **Example of Nested-Loop Join Costs**

- Compute student  $\bowtie$  takes, with student as the outer relation.
- Let takes have a primary B<sup>+</sup>-tree index on the attribute ID, which contains 20 entries in each index node.
- Since takes has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data
- *student* has 5000 tuples
- Cost of block nested loops join
  - 400\*100 + 100 = 40,100 block transfers + 2 \* 100 = 200 seeks
    - assuming worst case memory
    - may be significantly less with more memory
  - Cost of indexed nested loops join
    - 100 + 5000 \* 5 = 25,100 block transfers and seeks.
    - CPU cost likely to be less than that for block nested loops join



# **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.
  - 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 only 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:

 $b_r + b_s$  block transfers  $+ [b_r/b_b] + [b_s/b_b]$  seeks

+ the cost of sorting if relations are unsorted.

- hybrid merge-join: If one relation is sorted, and the other has a secondary B<sup>+</sup>-tree index on the join attribute
  - Merge the sorted relation with the leaf entries of the B+-tree.
  - Sort the result on the addresses of the unsorted relation's tuples
  - Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
    - Sequential scan more efficient than random lookup



#### Hash-Join

- 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  $i = h(t_r [JoinAttrs])$ .
  - $r_{0^{\prime}}, r_{1}..., r_{n}$  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<sub>i</sub>* and the *s* tuple in *s<sub>i</sub>*.

## **Transformation of Relational Expressions**

- Two relational algebra expressions are said to be **equivalent** if the two expressions generate the same set of tuples on every *legal* database instance
  - Note: order of tuples is irrelevant
  - we don't care if they generate different results on databases that violate integrity constraints
- In SQL, inputs and outputs are multisets of tuples
  - Two expressions in the multiset version of the relational algebra are said to be equivalent if the two expressions generate the same multiset of tuples on every legal database instance.
- An equivalence rule says that expressions of two forms are equivalent
  - Can replace expression of first form by second, or vice versa



#### **Equivalence Rules**

1. Conjunctive selection operations can be deconstructed into a sequence of individual selections.

$$\sigma_{\theta_1 \land \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$$

2. Selection operations are commutative.

 $\sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$ 

3. Only the last in a sequence of projection operations is needed, the others can be omitted.

 $\Pi_{L_1}(\Pi_{L_2}(\dots(\Pi_{Ln}(E))\dots)) = \Pi_{L_1}(E)$ 

4. Selections can be combined with Cartesian products and theta joins.

$$\alpha. \quad \sigma_{\theta}(\mathsf{E}_{1} \mathsf{X} \mathsf{E}_{2}) = \mathsf{E}_{1} \, \bowtie_{\theta} \mathsf{E}_{2}$$
  
$$\beta. \quad \sigma_{\theta 1}(\mathsf{E}_{1} \,\bowtie_{\theta 2} \mathsf{E}_{2}) = \mathsf{E}_{1} \, \bowtie_{\theta 1^{^{\wedge}} \theta 2} \mathsf{E}_{2}$$



- 5. Theta-join operations (and natural joins) are commutative.  $\bowtie_1 \bowtie_{\theta} E_2 = \bigotimes_{\theta} \bowtie_{\theta} E_1$
- 6. (a) Natural join operations are associative:  $(E_1 \quad E_2) \quad E_3 = E_1 \quad (E_2 \quad E_3)$

(b) Theta joins are associative in the following manner:

$$\bowtie(E_1 \quad \bigotimes_{\theta 1} E_2) \quad \underset{\theta 2^{\circ} \theta 3}{\longrightarrow} E_3 = E_1 \quad \underset{\theta 1^{\circ} \theta 3}{\longrightarrow} (E_2 \quad \underset{\theta 2}{\bowtie} E_3)$$

where  $\theta_2$  involves attributes from only  $E_2$  and  $E_3$ .



#### **Pictorial Depiction of Equivalence Rules**



Database System Concepts - 6th Edition

#### ©Silberschatz, Korth and Sudarshan



- 7. The selection operation distributes over the theta join operation under the following two conditions:
  - (a) When all the attributes in  $\theta_0$  involve only the attributes of one of the expressions ( $E_1$ ) being joined.

$$\sigma_{\theta 0}(\mathsf{E}_1 \boxtimes_{\theta} \mathsf{E}_2) = (\sigma_{\theta 0}(\mathsf{E}_1)) \boxtimes_{\theta} \mathsf{E}_2$$

(b) When  $\theta_1$  involves only the attributes of  $E_1$  and  $\theta_2$  involves only the attributes of  $E_2$ .

$$\sigma_{\theta_1 \theta_2} (\mathsf{E}_1 \bigcup_{\theta} \mathsf{E}_2) = (\sigma_{\theta_1} (\mathsf{E}_1)) \bigcup_{\theta} (\sigma_{\theta_2} (\mathsf{E}_2))$$



8. The projection operation distributes over the theta join operation as follows:

(a) if  $\theta$  involves only attributes from  $L_1 \cup L_2$ :

 $\prod_{L_1\cup L_2} (E_1\boxtimes_{\theta} E_2) = (\prod_{L_1} (E_1))\boxtimes_{\theta} (\prod_{L_2} (E_2))$ 

(b) Consider a join  $E_1 \bowtie_{\theta} E_2$ .

- Let  $L_1$  and  $L_2$  be sets of attributes from  $E_1$  and  $E_2$ , respectively.
- Let  $L_3$  be attributes of  $E_1$  that are involved in join condition  $\theta$ , but are not in  $L_1 \cup L_2$ , and
- let  $L_4$  be attributes of  $E_2$  that are involved in join condition  $\theta$ , but are not in  $L_1 \cup L_2$ .  $\prod_{L_1 \cup L_2} (E_1 \Join_{\theta}^1 E_2) \cong \prod_{L_1 \cup L_2} ((\prod_{L_1 \cup L_3} (E_1)) \Join_{\theta} (\prod_{L_2 \cup L_4} (E_2)))$



- 9. The set operations union and intersection are commutative  $E_1 \cup E_2 = E_2 \cup E_1$  $E_1 \cap E_2 = E_2 \cap E_1$ 
  - 9. (set difference is not commutative).
- 10. Set union and intersection are associative.

$$(E_{1} \cup E_{2}) \cup E_{3} = E_{1} \cup (E_{2} \cup E_{3})$$
$$(E_{1} \cap E_{2}) \cap E_{3} = E_{1} \cap (E_{2} \cap E_{3})$$

9. The selection operation distributes over  $\cup$ ,  $\cap$  and –.

$$\sigma_{\theta} (E_1 - E_2) = \sigma_{\theta} (E_1) - \sigma_{\theta} (E_2)$$
  
and similarly for  $\cup$  and  $\cap$  in place of

Also: 
$$O_{\theta}(E_1 - E_2) = O_{\theta}(E_1) - E_2$$

and similarly for  $\,\cap\,$  in place of  $\,$  –, but not for  $\,\cup\,$ 

12. The projection operation distributes over union

$$\Pi_{L}(E_{1} \cup E_{2}) = (\Pi_{L}(E_{1})) \cup (\Pi_{L}(E_{2}))$$

#### **Transformation Example: Pushing Selections**

- Query: Find the names of all instructors in the Music department, along with the titles of the courses that they teach
  - Π<sub>name, title</sub>(σ<sub>dept\_name= 'Music</sub>")
     (instructor κ (teache)

(instructor  $\bowtie$  (teaches  $\bowtie$   $\Pi_{course_{id, title}}$  (course))))

Transformation using rule 7a.

• 
$$\Pi_{name, title}((\sigma_{dept_name_{inter}}(instructor))) (feaches (  $\Pi_{course_{id, title}}(course)))$$$

Performing the selection as early as possible reduces the size of the relation to be joined.

## **Example with Multiple Transformations**

Query: Find the names of all instructors in the Music department who have taught a course in 2009, along with the titles of the courses that they taught

•  $\Pi_{name, title}(\sigma_{dept_name= 'Music'' year = 2009}))$ (instructor  $\bowtie$  (teaches)  $\bowtie$   $\Pi_{course_id, title}$  (course)))))

- Transformation using join associatively (Rule 6a):
  - $\Pi_{name, title}(\sigma_{dept, name= "Music"^gear = 2009})$ ((instructor teaches)  $\Pi_{course_id, title}$  (course)))

Second form provides an opportunity to apply the "perform selections early" rule, resulting in the subexpression

$$\sigma_{dept\_name = `Music"}(instructor) \quad \sigma_{year = 2009}$$

(teaches)



## **Multiple Transformations (Cont.)**



(a) Initial expression tree

(b) Tree after multiple transformations

#### **Transformation Example: Pushing Projections**

- Consider:  $\Pi_{name, title}(\sigma_{dept_name= "Music"}(instructor) teaches)$  $\bowtie_{course_id, title}(course))))$
- When we compute

( $\sigma_{dept_name = "Music"}$  (instructor teaches)

we obtain a relation whose schema is:

```
(ID, name, dept_name, salary, course_id, sec_id, semester, year)
```

Push projections using equivalence rules 8a and 8b; eliminate unneeded attributes from intermediate results to get:

$$\Pi_{name, title}(\Pi_{name, course_id} ( \\ \bowtie_{dept_name= `Music"} (instructor) teaches)) \\ \Pi_{course_id, title} (course))))$$

Performing the projection as early as possible reduces the size of the relation to be joined.



## **Join Ordering Example**

For all relations 
$$r_1, r_2$$
, and  $r_3$ ,  
 $(\not R_1 \quad \not R_2) \quad r_3 = \not r_1 \quad (\not R_2 \quad r_3)$   
(Join Associativity)  
If  $r_2 \bowtie r_3$  is quite large and  $r_1 \bowtie r_2$  is small, we choose  
 $( \begin{matrix} & & & \\ r_1 & & r_2 \end{matrix}) \quad r_3$ 

so that we compute and store a smaller temporary relation.



# Join Ordering Example (Cont.)

Consider the expression

 $\Pi_{name, title}(\sigma_{dept\_name= `Music"} (instructor) teaches) \bowtie \\ \bowtie \Pi_{course id. title} (course))))$ 

Could compute teaches  $\bowtie \Pi_{course\_id, title}$  (course) first, and join result with

 $\sigma_{dept_name= \text{'Music''}}(instructor)$ but the result of the first join is likely to be a large relation.

- Only a small fraction of the university's instructors are likely to be from the Music department
  - it is better to compute

*σ*<sub>dept\_name= "Music"</sub> (instructor) teaches

first.

## **Enumeration of Equivalent Expressions**

- Query optimizers use equivalence rules to **systematically** generate expressions equivalent to the given expression
- Can generate all equivalent expressions as follows:
  - Repeat
    - apply all applicable equivalence rules on every subexpression of every equivalent expression found so far
    - add newly generated expressions to the set of equivalent expressions

Until no new equivalent expressions are generated above

- The above approach is very expensive in space and time
  - Two approaches
    - Optimized plan generation based on transformation rules
    - Special case approach for queries with only selections, projections and joins



## **Choice of Evaluation Plans**

- Must consider the interaction of evaluation techniques when choosing evaluation plans
  - choosing the cheapest algorithm for each operation independently may not yield best overall algorithm. E.g.
    - merge-join may be costlier than hash-join, but may provide a sorted output which reduces the cost for an outer level aggregation.
    - nested-loop join may provide opportunity for pipelining
- Practical query optimizers incorporate elements of the following two broad approaches:
  - 1. Search all the plans and choose the best plan in a cost-based fashion.
  - 2. Uses heuristics to choose a plan.



## **Choice of Evaluation Plans**

- Must consider the interaction of evaluation techniques when choosing evaluation plans
  - choosing the cheapest algorithm for each operation independently may not yield best overall algorithm. E.g.
    - merge-join may be costlier than hash-join, but may provide a sorted output which reduces the cost for an outer level aggregation.
    - nested-loop join may provide opportunity for pipelining
- Practical query optimizers incorporate elements of the following two broad approaches:
  - 1. Search all the plans and choose the best plan in a cost-based fashion.
  - 2. Uses heuristics to choose a plan.