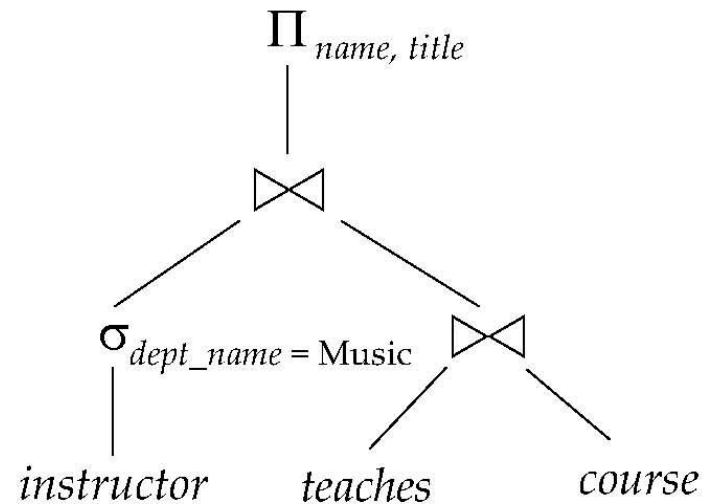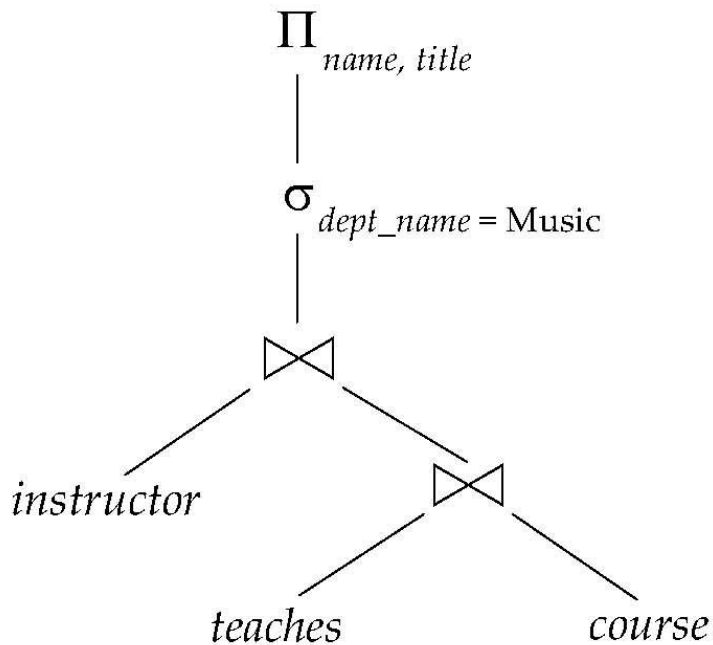# Chapter 13: Query Optimization

# Chapter 13:  Query Optimization

- Introduction

- Transformation of Relational Expressions

- Catalog Information for Cost Estimation

- Statistical Information for Cost Estimation

- Cost-based optimization

- Dynamic Programming for Choosing Evaluation Plans
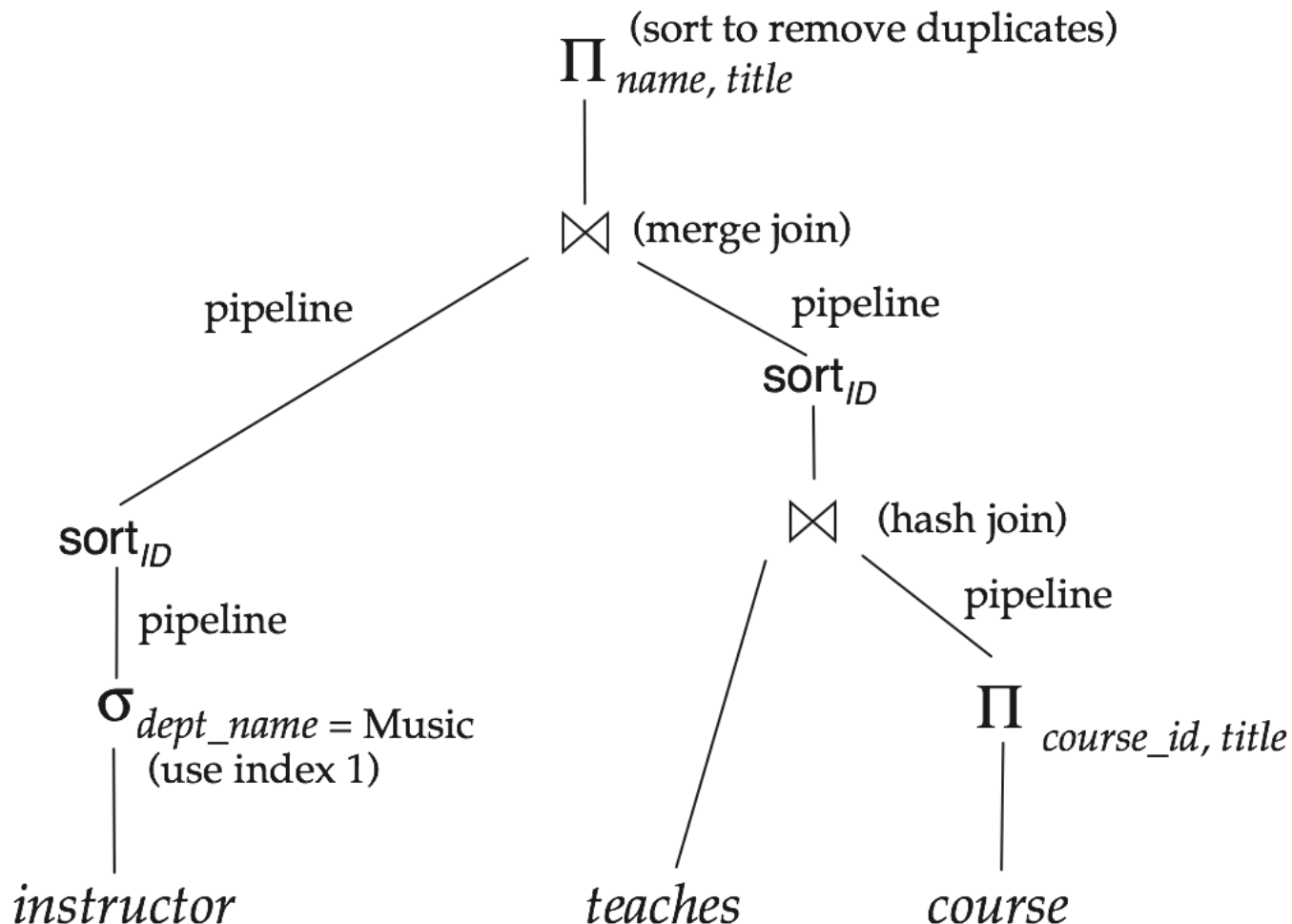
- Materialized views

# 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

# Viewing Query Execution Plans

- All database provide ways to view query execution plans

- E.g. in PostgreSQL, prefix an SQL query with the keyword **explain** to see the plan that is chosen.

- In SQL Server, execute  **set showplan_text on**

  - Any query submitted after this will show the plan instead of executing the query

    - use **set showplan_text off**  to stop showing plans

# 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**
  1. 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

# Generating Equivalent Expressions

# 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.

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

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.

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

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

   a. $\sigma_{\theta}(E_1 \times E_2) = E_1 \bowtie_{\theta} E_2$

   b. $\sigma_{\theta 1}(E_1 \bowtie_{\theta 2} E_2) = E_1 \bowtie_{\theta 1 \wedge \theta 2} E_2$

# Equivalence Rules (Cont.)

5. Theta-join operations (and natural joins) are commutative.

$$E_1 \bowtie_\theta E_2 = E_2 \bowtie_\theta E_1$$

6. (a) Natural join operations are associative:

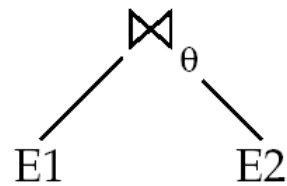$$(E_1 \bowtie E_2) \bowtie E_3 = E_1 \bowtie (E_2 \bowtie E_3)$$

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

$$(E_1 \bowtie_{\theta 1} E_2) \bowtie_{\theta 2 \wedge \theta 3} E_3 = E_1 \bowtie_{\theta 1 \wedge \theta 3} (E_2 \bowtie_{\theta 2} E_3)$$
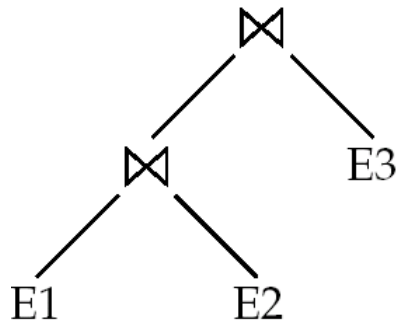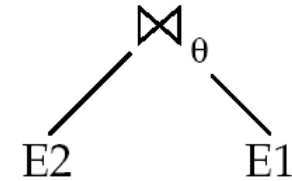
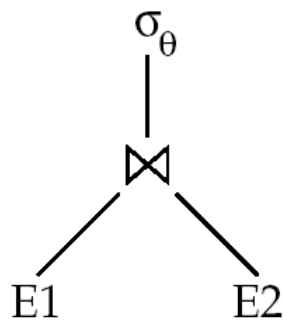where $\theta_2$ involves attributes from only $E_2$ and $E_3$.

# Pictorial Depiction of Equivalence Rules

# Equivalence Rules (Cont.)

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.

   $$\sigma_{\theta 0}(E_1 \bowtie_\theta E_2) = (\sigma_{\theta 0}(E_1)) \bowtie_\theta E_2$$

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

   $$\sigma_{\theta 1 \wedge \theta 2}(E_1 \bowtie_\theta E_2) = (\sigma_{\theta 1}(E_1)) \bowtie_\theta (\sigma_{\theta 2}(E_2))$$

# Equivalence Rules (Cont.)

- Other rules in the book include
    - Pushing projection through joins
    - Set operations
        - associativity/commutativity, except for set difference
        - selection and projection distribute over set operations

# Multiple Transformations



(a) Initial expression tree

(b) Tree after multiple transformations

# Quiz Time

**Quiz Q1**: The expression $\sigma_{r.A=5}$ $(r \bowtie s)$ is equivalent to which of these expressions, given relations r(A,B) and s(B,C)?

(1) $\sigma_{r.A=5}$ $(r) \bowtie s$

(2) $\sigma_{r.A=5}$ $(s) \bowtie r$

(3) neither

(4) both

# Join Ordering Example

- For all relations $r_1$, $r_2$, and $r_3$,

$$(r_1 \bowtie r_2) \bowtie r_3 = r_1 \bowtie (r_2 \bowtie r_3)$$

(Join Associativity)

- If $r_2 \bowtie r_3$ is quite large and $r_1 \bowtie r_2$ is small, we choose

$$(r_1 \bowtie r_2) \bowtie r_3$$

so that we compute and store a smaller temporary relation.

# Join Ordering Example (Cont.)

- Consider the expression

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

- 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

    $$\sigma_{dept\_name = \text{"Music"}} (instructor) \bowtie teaches$$

    first.

# Enumeration of Equivalent Expressions

- Some query optimizers use equivalence rules to **systematically** generate expressions equivalent to the given expression

- Others are special cased for join order optimization, along with heuristics for pushing selections, and other heuristics for other operations such as aggregation

# Cost Estimation

- Cost of each operator computer as described in Chapter 12
  - Need statistics of input relations
    - E.g. number of tuples, sizes of tuples
- Inputs can be results of sub-expressions
  - Need to estimate statistics of expression results
  - To do so, we require additional statistics
    - E.g. number of distinct values for an attribute
- More details on cost estimation are in the book

# Cost-Based Optimization

- Consider finding the best join-order for $r_1 \bowtie r_2 \bowtie \ldots r_n$.

- There are $(2(n-1))!/(n-1)!$ different join orders for above expression. With $n = 7$, the number is 665280, with $n = 10$, the number is greater than 176 billion!

- No need to generate all the join orders. Using dynamic programming, the least-cost join order for any subset of $\{r_1, r_2, \ldots r_n\}$ is computed only once and stored for future use.

# Dynamic Programming in Optimization

- To find best join tree for a set of $n$ relations:

  - To find best plan for a set $S$ of $n$ relations, consider all possible plans of the form: $S_1 \bowtie (S - S_1)$ where $S_1$ is any non-empty subset of $S$.

  - Recursively compute costs for joining subsets of $S$ to find the cost of each plan. Choose the cheapest of the $2^n - 2$ alternatives.

  - Base case for recursion: single relation access plan
    - Apply all selections on $R_i$ using best choice of indices on $R_i$

  - When plan for any subset is computed, store it and reuse it when it is required again, instead of recomputing it
    - Dynamic programming

# Join Order Optimization Algorithm

procedure findbestplan($S$)
    if ($bestplan[S].cost \neq \infty$)
        **return** $bestplan[S]$
    // else $bestplan[S]$ has not been computed earlier, compute it now
    **if** ($S$ contains only 1 relation)
        set $bestplan[S].plan$ and $bestplan[S].cost$ based on the best way
        of accessing $S$  /* Using selections on S and indices on S */

    **else for each** non-empty subset $S1$ of $S$ such that $S1 \neq S$
        P1= findbestplan($S1$)
        P2= findbestplan($S - S1$)
        A = best algorithm for joining results of $P1$ and $P2$
        cost = $P1.cost$ + $P2.cost$ + cost of $A$
        **if** $cost < bestplan[S].cost$
            $bestplan[S].cost$ = cost
            $bestplan[S].plan$ = "execute $P1.plan$; execute $P2.plan$;
                join results of $P1$ and $P2$ using $A$"

    **return** $bestplan[S]$

\* Some modifications to allow indexed nested loops joins on relations that have
  selections (see book)

# Additional Optimization Techniques

- Nested Subqueries

- Materialized Views

# Optimizing Nested Subqueries**

- Nested query example:
  **select** *name*
  **from** *instructor*
  **where exists** (**select** *
          **from** *teaches*
          **where** *instructor.ID = teaches.ID* **and** *teaches.year = 2007*)

- SQL conceptually treats nested subqueries in the where clause as functions that take parameters and return a single value or set of values

  - Parameters are variables from outer level query that are used in the nested subquery; such variables are called **correlation variables**

- Conceptually, nested subquery is executed once for each tuple in the cross-product generated by the outer level **from** clause

  - Such evaluation is called **correlated evaluation**

  - Note: other conditions in where clause may be used to compute a join (instead of a cross-product) before executing the nested subquery

# Optimizing Nested Subqueries (Cont.)

- Correlated evaluation may be quite inefficient since
    - a large number of calls may be made to the nested query
    - there may be unnecessary random I/O as a result
- SQL optimizers attempt to transform nested subqueries to joins where possible, enabling use of efficient join techniques
- E.g.: earlier nested query can be rewritten as
  **select** *name*
  **from** *instructor, teaches*
  **where** *instructor.ID = teaches.ID* **and** *teaches.year = 2007*
    - Note: the two queries generate different numbers of duplicates (why?)
        - ‣ teaches can have duplicate IDs
        - ‣ Can be modified to handle duplicates correctly as we will see
- In general, it is not possible/straightforward to move the entire nested subquery from clause into the outer level query from clause
    - A temporary relation is created instead, and used in body of outer level query

# Optimizing Nested Subqueries (Cont.)

- In our example, the original nested query would be transformed to

  **create table** $t_1$ **as**
  > **select distinct** *ID*
  > **from** *teaches*
  > **where** *year* = *2007*

  **select** *name*
  **from** *instructor*, $t_1$
  **where** $t_1$.*ID* = *instructor.ID*

- The process of replacing a nested query by a query with a join (possibly with a temporary relation) is called **decorrelation**.

- Decorrelation is more complicated when

  - the nested subquery uses aggregation, or

  - when the result of the nested subquery is used to test for equality, or

  - when the condition linking the nested subquery to the other query is **not exists**,

  - and so on.

# Quiz Time

**Quiz Q2**: Given an option of writing a query using a join versus using a correlated subquery
(1) it is always better to write it using a subquery
(2) some optimizers are likely to get a better plan if the query is written using a join than if written using a subquery
(3) some optimizers are likely to get a better plan if the query is written using a subquery
(4) none of the above.

# Materialized Views**

- A **materialized view** is a view whose contents are computed and stored.

- Consider the view
  c**reate view** *department_total_salary*(*dept_name, total_salary*) **as**
  **select** *dept_name*, **sum**(*salary*)
  **from** *instructor*
  **group by** *dept_name*

- Materializing the above view would be very useful if the total salary by department is required frequently

  - Saves the effort of finding multiple tuples and adding up their amounts

# Materialized View Maintenance

- The task of keeping a materialized view up-to-date with the underlying data is known as **materialized view maintenance**

- Materialized views can be maintained by recomputation on every update

- A better option is to use **incremental view maintenance**

  - **Changes to database relations are used to compute changes to the materialized view, which is then updated**

- See book for details on incremental view maintenance

# Materialized View Selection

- **Materialized view selection**: "What is the best set of views to materialize?".

- **Index selection:** "what is the best set of indices to create"
    - closely related, to materialized view selection
        ▸ but simpler
- Materialized view selection and index selection based on typical system **workload** (queries and updates)
    - Typical goal: minimize time to execute workload , subject to constraints on space and time taken for some critical queries/updates
    - One of the steps in database tuning
        ▸ more on tuning in later chapters
- Commercial database systems provide tools (called "tuning assistants" or "wizards") to help the database administrator choose what indices and materialized views to create

# Additional Optimization Techniques

■ See book for details on the following advanced optimization techniques

- Top-K queries

- Halloween problem

  ▸ **update** R **set** A = 5 * A
    **where** A > 10

- Join minimization

- Multiquery optimization

- Parametric query optimization

# Quiz Time

**Quiz Q3**: If all data is stored in main memory
(1) query optimization will no longer be required
(2) query optimization will still be required, but queries will run faster
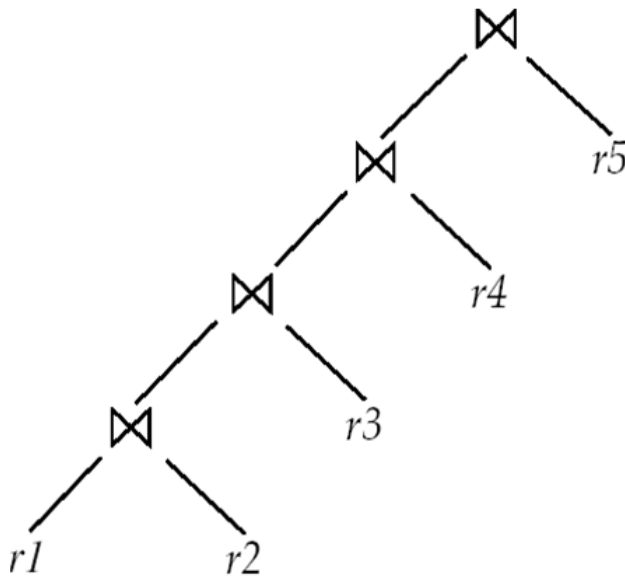(3) query optimization will still be required, but queries will run slower
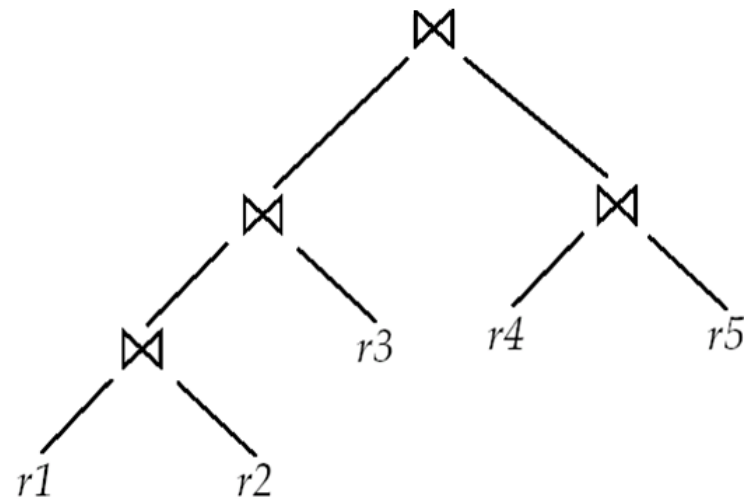(4) none of the above

# End of Chapter

# Left Deep Join Trees

■ In **left-deep join trees,** the right-hand-side input for each join is a relation, not the result of an intermediate join.



(a) Left-deep join tree

(b) Non-left-deep join tree

# Cost of Optimization

- With dynamic programming time complexity of optimization with bushy trees is $O(3^n)$.
    - With $n = 10$, this number is 59000 instead of 176 billion!
- Space complexity is $O(2^n)$
- To find best left-deep join tree for a set of $n$ relations:
    - Consider $n$ alternatives with one relation as right-hand side input and the other relations as left-hand side input.
    - Modify optimization algorithm:
        - Replace "**for each** non-empty subset $S1$ of $S$ such that $S1 \neq S$"
        - By: **for each** relation r in S
                    let S1 = S − r .
- If only left-deep trees are considered, time complexity of finding best join order is $O(n\, 2^n)$
    - Space complexity remains at $O(2^n)$
- Cost-based optimization is expensive, but worthwhile for queries on large datasets (typical queries have small n, generally < 10)

# Interesting Sort Orders

- Consider the expression $(r_1 \bowtie r_2) \bowtie r_3$  (with A as common attribute)

- An **interesting sort order**  is a particular sort order of tuples that could be useful for a later operation

  - Using merge-join to compute $r_1 \bowtie r_2$  may be costlier than hash join but generates result sorted on A

  - Which in turn may make merge-join with $r_3$ cheaper, which may reduce cost of join with $r_3$ and minimizing overall cost

  - Sort order may also be useful for order by and for grouping

- Not sufficient to find the best join order for each subset of the set of *n* given relations

  - must find the best join order for each subset, **for each interesting sort order**

  - Simple extension of earlier dynamic programming algorithms

  - Usually, number of interesting orders is quite small and doesn't affect time/space complexity significantly

# Statistics for Cost Estimation
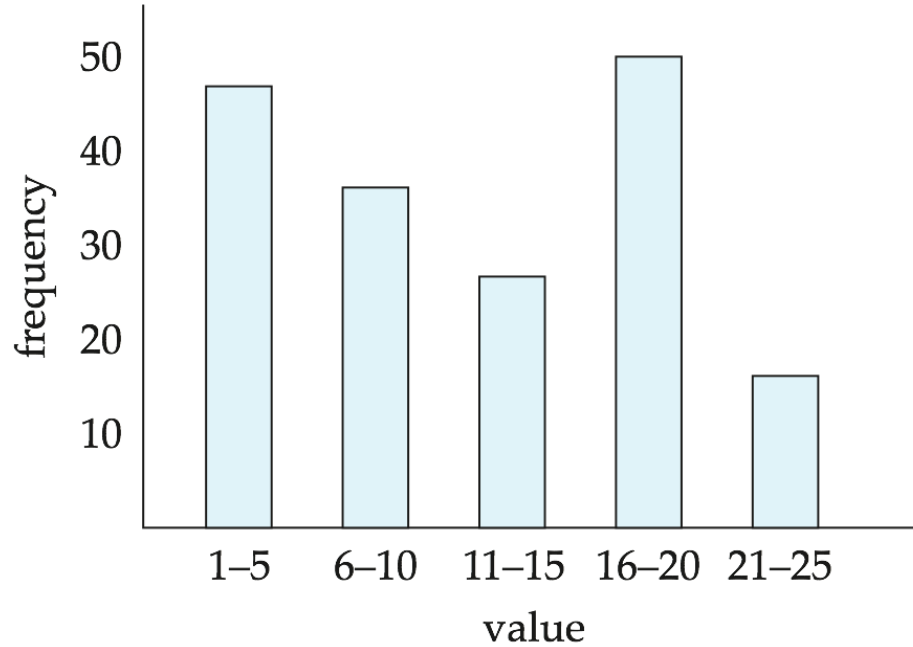
# Statistical Information for Cost Estimation

- $n_r$: number of tuples in a relation $r$.

- $b_r$: number of blocks containing tuples of $r$.

- $l_r$: size of a tuple of $r$.

- $f_r$: blocking factor of $r$ — i.e., the number of tuples of $r$ that fit into one block.

- $V(A, r)$: number of distinct values that appear in $r$ for attribute $A$; same as the size of $\prod_A(r)$.

- If tuples of $r$ are stored together physically in a file, then:

$$b_r = \left\lceil \frac{n_r}{f_r} \right\rceil$$

# Histograms

- Histogram on attribute *age* of relation *person*



- **Equi-width** histograms
- **Equi-depth** histograms

# Selection Size Estimation

- $\sigma_{A=v}(r)$

  - ▸ $n_r / V(A,r)$ : number of records that will satisfy the selection
  - ▸ Equality condition on a key attribute: *size estimate* $= 1$

- $\sigma_{A \leq V}(r)$ (case of $\sigma_{A \geq V}(r)$ is symmetric)

  - Let c denote the estimated number of tuples satisfying the condition.
  - If min(A,r) and max(A,r) are available in catalog
    - ▸ c = 0 if v < min(A,r)
    - ▸ c = $n_r \cdot \dfrac{v - \min(A,r)}{\max(A,r) - \min(A,r)}$

  - If histograms available, can refine above estimate
  - In absence of statistical information $c$ is assumed to be $n_r / 2$.

# Size Estimation of Complex Selections

- The **selectivity** of a condition $\theta_i$ is the probability that a tuple in the relation $r$ satisfies $\theta_i$ .

  - If $s_i$ is the number of satisfying tuples in $r$, the selectivity of $\theta_i$ is given by $s_i/n_r$.

- **Conjunction:** $\sigma_{\theta1 \wedge \theta2 \wedge \dots \wedge \theta n}(r)$. *Assuming indepdence,* estimate of tuples in the result is:

$$n_r * \frac{s_1 * s_2 * \dots * s_n}{n_r^n}$$

- **Disjunction:** $\sigma_{\theta1 \vee \theta2 \vee \dots \vee \theta n}(r)$. Estimated number of tuples:

$$n_r * \left( 1 - (1 - \frac{s_1}{n_r}) * (1 - \frac{s_2}{n_r}) * \dots * (1 - \frac{s_n}{n_r}) \right)$$

- **Negation:** $\sigma_{\neg\theta}(r)$. Estimated number of tuples:

$$n_r - size(\sigma_\theta(r))$$

# Estimation of the Size of Joins

- The Cartesian product $r \times s$ contains $n_r \cdot n_s$ tuples; each tuple occupies $s_r + s_s$ bytes.

- If $R \cap S = \emptyset$, then $r \bowtie s$ is the same as $r \times s$.

- If $R \cap S$ is a key for $R$, then a tuple of $s$ will join with at most one tuple from $r$

  - therefore, the number of tuples in $r \bowtie s$ is no greater than the number of tuples in $s$.

- If $R \cap S$ in S is a foreign key in $S$ referencing $R$, then the number of tuples in $r \bowtie s$ is exactly the same as the number of tuples in $s$.

    - The case for $R \cap S$ being a foreign key referencing $S$ is symmetric.

- In the example query $student \bowtie takes$, ID in $takes$ is a foreign key referencing $student$

  - hence, the result has exactly $n_{takes}$ tuples, which is 10000

# Estimation of the Size of Joins (Cont.)

- If $R \cap S = \{A\}$ is not a key for $R$ or $S$.
  If we assume that every tuple $t$ in $R$ produces tuples in $R \bowtie S$, the number of tuples in $R \bowtie S$ is estimated to be:

$$\frac{n_r * n_s}{V(A, s)}$$

  If the reverse is true, the estimate obtained will be:

$$\frac{n_r * n_s}{V(A, r)}$$

  The lower of these two estimates is probably the more accurate one.

- Can improve on above if histograms are available

  - Use formula similar to above, for each cell of histograms on the two relations

# Estimation of the Size of Joins (Cont.)

- Compute the size estimates for *depositor* $\bowtie$ *customer* without using information about foreign keys:

  - *V(ID, takes)* = 2500, and
    *V(ID, student)* = 5000

  - The two estimates are 5000 * 10000/2500 = 20,000 and 5000 * 10000/5000 = 10000

  - We choose the lower estimate, which in this case, is the same as our earlier computation using foreign keys.

# More on Estimation

- See book for
    - Size estimation details for other operations
    - Details on how to estimate number of distinct values in result of an operation