

BITS Pilani
Pilani Campus

Dr. Mukesh Kumar Rohil
Department of Computer Science & Information Systems



Database Systems: Query Processing and Query Optimization

Database Systems: Query Processing and Query Optimization



Dr. Mukesh Kumar Rohil
Instructor of the course on
Database Systems

**Computer Science &
Information Systems
Department**

**Birla Institute of
Technology & Science
Pilani – 333031
(Rajasthan)**

Learning Objectives



Query Processing

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

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

Resources



The subsequent slides are adopted and/or prepared from the material available in/with the following prescribed textbook (TB) for the course:

1. Silberschatz A, Korth H F, & Sudarshan S, ***Database System Concepts***, 7e, TMH, 2019..



Introduction

- **Query processing** refers to the range of activities involved in extracting data from a database.
- The activities include:
 1. translation of queries in high-level database languages into expressions that can be used at the physical level of the file system,
 2. a variety of query-optimizing transformations, and
 3. actual evaluation of queries.

Introduction ...2



- **Query optimization** is the process of selecting the most efficient query-evaluation plan from among the many strategies usually possible for processing a given query, especially if the query is complex.
- Users may not write their queries ready to be processed efficiently.
- Expect the system to construct a query-evaluation plan that minimizes the cost of query evaluation.
- **One aspect of optimization** occurs at the relational-algebra level, where the system attempts to find an expression that is equivalent to the given expression, but more efficient to execute.

Introduction ...3

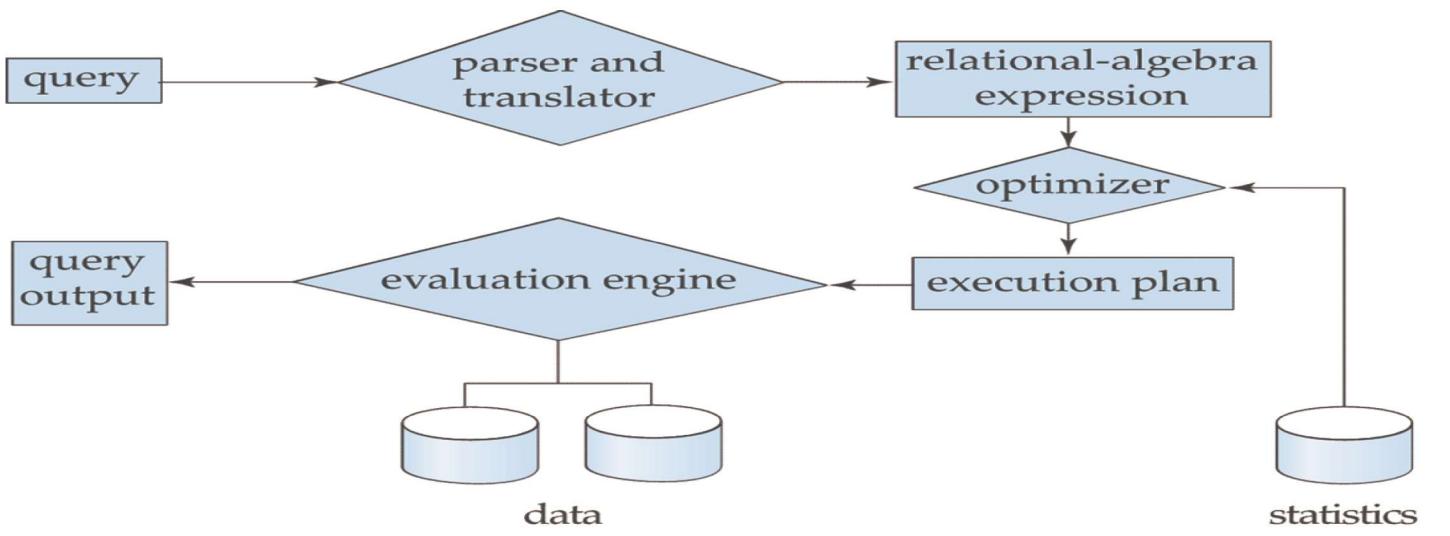


- **Another aspect of optimization** is selecting a detailed strategy for processing the query, such as choosing the algorithm to use for executing an operation, choosing the specific indices to use, and so on.
- The **difference in cost** (in terms of evaluation time) between a good strategy and a bad strategy is often substantial and may be several orders of magnitude.
- Hence, it is worthwhile for the system to spend a substantial amount of time on the selection of a good strategy for processing a query, even if the query is executed only once.

Basic Steps in Query Processing



1. Parsing and translation
2. Optimization
3. Evaluation



Basic Steps in Query Processing (Cont.)



❑ Parsing and translation

- ❑ translate the query into its internal form. This is then translated into relational algebra.
- ❑ Parser checks syntax, verifies relations

❑ Evaluation

- ❑ The query-execution engine takes a query-evaluation plan, executes that plan, and returns the answers to the query.

Basic Steps in Query Processing: Optimisation



- A relational algebra expression may have many equivalent expressions
 - E.g., $\sigma_{\text{salary} < 75000}(\Pi_{\text{salary}}(\text{instructor}))$ is equivalent to $\Pi_{\text{salary}}(\sigma_{\text{salary} < 75000}(\text{instructor}))$
- Each relational algebra operation can be evaluated using one of several different algorithms
 - Correspondingly, a relational-algebra expression can be evaluated in many ways.
- Annotated expression specifying detailed evaluation strategy is called an **evaluation-plan**.
 - E.g., can use an index on *salary* to find instructors with $\text{salary} < 75000$,
 - or can perform complete relation scan and discard instructors with $\text{salary} \geq 75000$

Basic Steps in Query Processing: Optimization (Cont.)



- **Query Optimization:** Amongst all equivalent evaluation plans choose the one with lowest cost.
 - Cost is estimated using statistical information from the database catalog
 - e.g. number of tuples in each relation, size of tuples, etc.
- We will discuss:
 - 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?
 - Later, we will discuss how to optimize queries, that is, how to find an evaluation plan with lowest estimated cost.



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 ...2

- 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

Measures of Query Cost ...3



- Several algorithms can reduce disk I/O by using extra buffer space
 - Amount of real memory available to buffer depends on other concurrent queries and OS processes, known only during execution
 - We often use worst case estimates, assuming only the minimum amount of memory needed for the operation is available
- Required data may be buffer resident already, avoiding disk I/O
 - But hard to take into account for cost estimation

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
- Note: binary search generally does not make sense since data is not stored consecutively
 - except when there is an index available,
 - and binary search requires more seeks than index search



Selections Using Indices

- **Index scan** – search algorithms that use an index
 - selection condition must be on search-key of index.

Assumptions:

1. h_i = height of the B+-tree, and
 2. a random I/O operation is required for each node in the path from the root to a leaf.
- **A2 (primary index, equality on key)**. Retrieve a single record that satisfies the corresponding equality condition
 - $\text{Cost} = (h_i + 1) * (t_T + t_S)$
 - **A3 (primary index, equality on non-key)** Retrieve multiple records.
 - Records will be on consecutive blocks
 - Let b = number of blocks containing matching records
 - $\text{Cost} = h_i * (t_T + t_S) + t_S + t_T * b$

Selections Using Indices ...2



- **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



- 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 $> v$
 - In either case, retrieve records that are pointed to
 - requires an I/O for each record
 - Linear file scan may be cheaper

Implementation of Complex Selections



- **Conjunction:** $\sigma_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n}(r)$
- **A7 (conjunctive selection using one index).**
 - Select a combination of θ_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



- **Disjunction:** $\sigma_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n}(r)$.
- **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
- **Negation:** $\sigma_{\neg\theta}(r)$
 - Use linear scan on file
 - If very few records satisfy $\neg\theta$, and an index is applicable to θ
 - Find satisfying records using index and fetch from file

Sorting & External Sort-Merge



- 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.

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 .
 - (d) Let the final value of i be N
2. **Merge the runs (*next slide*).....**



External Sort-Merge (Cont.)

2. Merge the runs (N-way merge). We assume (for now) that $N < M$.

1. Use N blocks of memory to buffer input runs, and 1 block to buffer output. Read the first block of each run into its buffer page
2. **repeat**
 1. Select the first record (in sort order) among all buffer pages
 2. Write the record to the output buffer. If the output buffer is full write it to disk.
 3. Delete the record from its input buffer page. If the buffer page becomes empty then read the next block (if any) of the run into the buffer.
3. **until** all input buffer pages are empty:

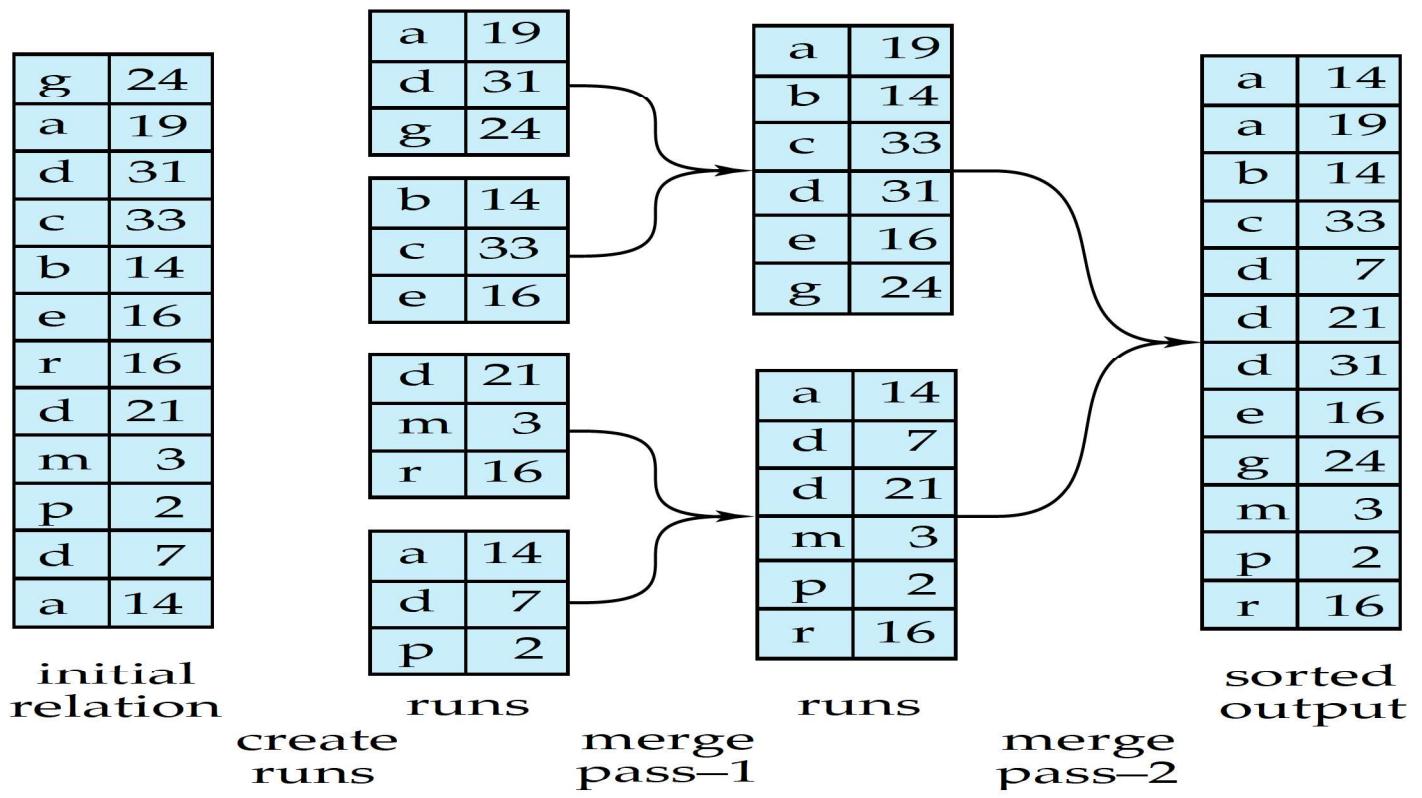
External Sort-Merge (Cont.)



If $N \geq M$, several merge passes are required.

- In each pass, contiguous groups of $M - 1$ runs are merged.
- A pass reduces the number of runs by a factor of $M - 1$, and creates runs longer by the same factor.
 - E.g. If $M=11$, and there are 90 runs, one pass reduces the number of runs to 9, each 10 times the size of the initial runs
- Repeated passes are performed till all runs have been merged into one.

Example: External Sorting Using Sort-Merge





External Merge Sort (Cont.)

Cost analysis:

- 1 block per run leads to too many seeks during merge
 - Instead use b_b buffer blocks per run
 - ➔ read/write b_b blocks at a time
 - Can merge $\lfloor M/b_b \rfloor - 1$ runs in one pass
- Total number of merge passes required: $\lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil$.
- Block transfers for initial run creation as well as in each pass is $2b_r$,
 - for final pass, we don't count write cost
 - we ignore final write cost for all operations since the output of an operation may be sent to the parent operation without being written to disk
 - Thus total number of block transfers for external sorting:
$$b_r (2 \lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil + 1)$$
- Seek: discussed in the next slide

External Merge Sort (Cont.)



- Cost of seeks
 - During run generation: one seek to read each run and one seek to write each run
 - $2\lceil b_r/M \rceil$
 - During the merge phase
 - Need $2\lceil b_r/b_b \rceil$ seeks for each merge pass
 - except the final one which does not require a write
 - Total number of seeks:
$$2\lceil b_r/M \rceil + \lceil b_r/b_b \rceil (2\lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil - 1)$$

Join Operation



- Several different algorithms to implement joins
 - **Nested-loop join:** For each row on the left-hand side, the right-hand side is scanned for a match based on the join condition.
 - Block nested-loop join
 - Indexed nested-loop join
 - **Merge-join:** A Merge Join relies on its two inputs being sorted on the join attributes.
 - **Hash-join:** the Hash Join algorithm builds an in-memory hash table of the smaller of its two inputs, and then reads the larger input and probes the in-memory hash table to find matches.
- 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 θ
 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.
- 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
$$n_r * b_s + b_r \text{ block transfers, plus}$$
$$n_r + b_r \text{ seeks}$$
- 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.
- 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  $B_r$  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.
- 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 = $\lceil b_r / (M-2) \rceil * b_s + b_r$ block transfers + $2 \lceil b_r / (M-2) \rceil$ seeks
 - If equi-join attribute forms a key on 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 Least Recently Used (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: $b_r(t_T + t_S) + n_r * c$
 - Where c is the cost of traversing index and fetching all matching s tuples for one tuple of 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

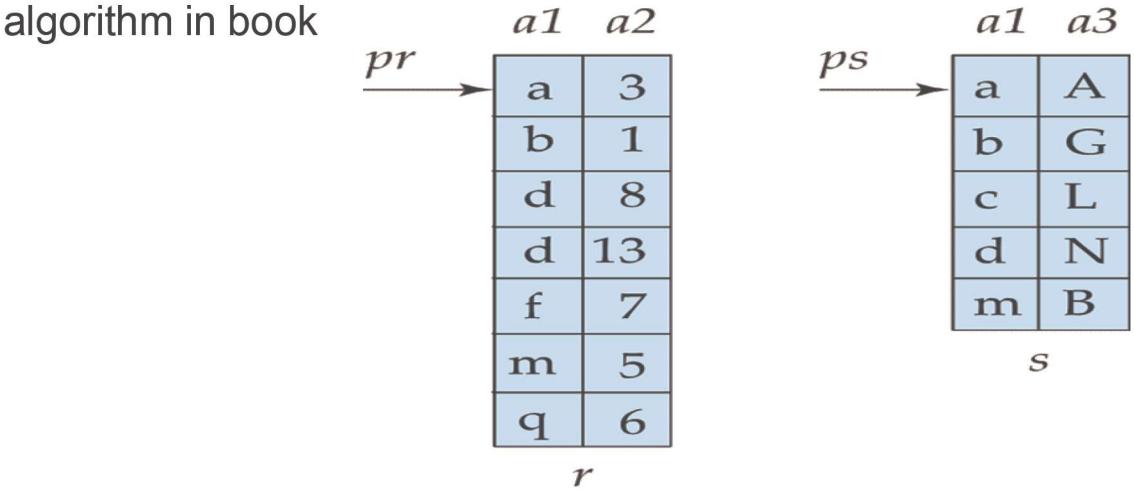


- Compute $student \bowtie takes$, with $student$ as the outer relation.
- Let $takes$ have a primary B⁺-tree index on the attribute ID , which contains 20 entries in each index node.
- Since $takes$ has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data
- $student$ 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.
 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 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:

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

- + the cost of sorting if relations are unsorted.

hybrid merge-join: If one relation is sorted, and the other has a secondary B⁺-tree index on the join attribute

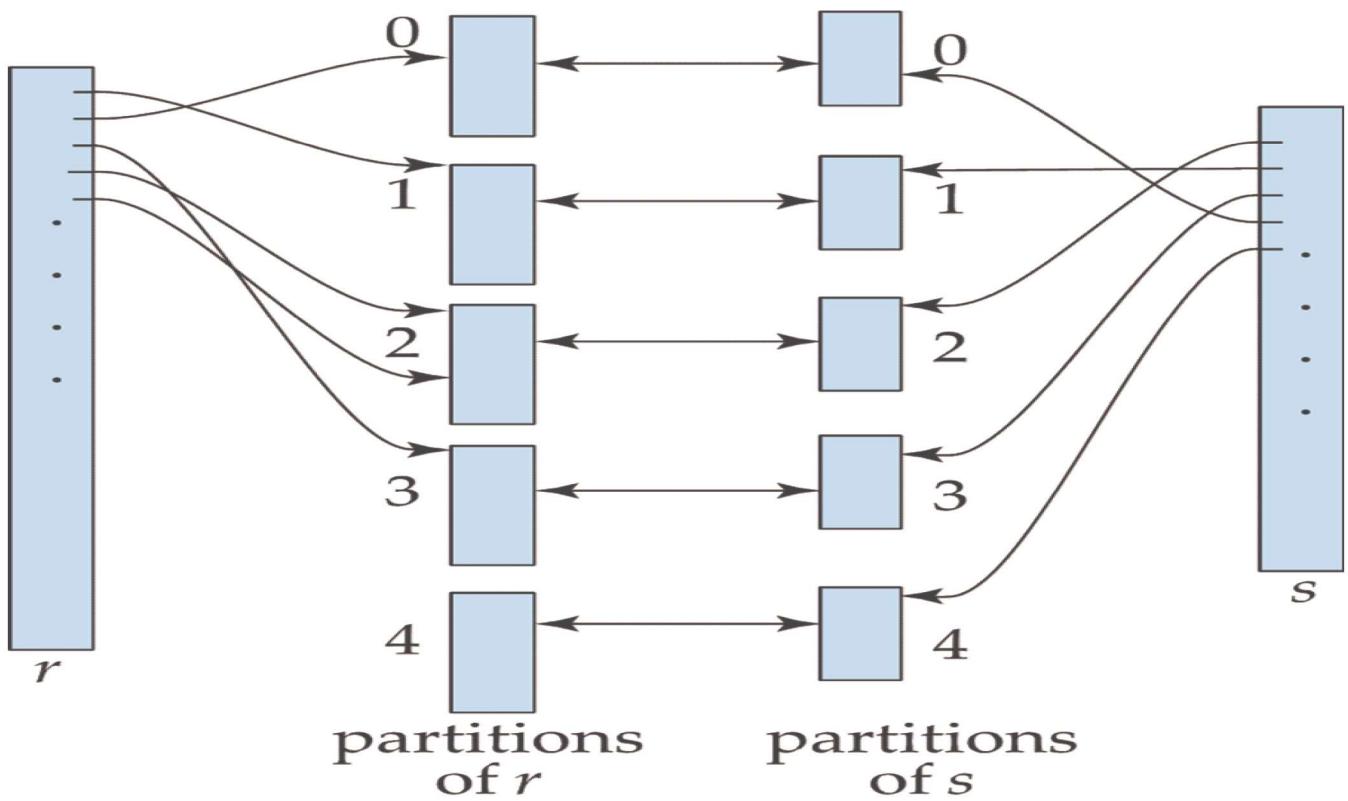
- Merge the sorted relation with the leaf entries of the B⁺-tree .
- Sort the result on the addresses of the unsorted relation's tuples
- Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
 - Sequential scan more efficient than 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 $\{0, 1, \dots, n\}$, where $JoinAttrs$ denotes the common attributes of r and s used in the natural join.
 - r_0, r_1, \dots, r_n denote partitions of r tuples
 - Each tuple $t_r \in r$ is put in partition r_i where $i = h(t_r[JoinAttrs])$.
 - s_0, s_1, \dots, s_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.)



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 .

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. This hash index uses a different hash function than the earlier one h .
 - (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
- **Recursive partitioning** required if number of partitions n is greater than number of pages M of memory.
 - instead of partitioning n ways, use $M - 1$ partitions for s
 - Further partition the $M - 1$ partitions using a different hash function
 - Use same partitioning method on r
 - Rarely required: e.g., with block size of 4 KB, recursive partitioning not needed for relations of < 1GB with memory size of 2MB, or relations of < 36 GB with memory of 12 MB

Handling of Overflows



- Partitioning is said to be **skewed** if some partitions have significantly more tuples than some others
- **Hash-table overflow** occurs in partition s_i if s_i does not fit in memory. Reasons could be
 - Many tuples in s_i with same value for join attributes
 - Bad hash function
- **Overflow resolution** can be done in build phase
 - Partition s_i is further partitioned using different hash function.
 - Partition r_j must be similarly partitioned.
- **Overflow avoidance** performs partitioning carefully to avoid overflows during build phase
 - E.g. partition build relation into many partitions, then combine them
- Both approaches fail with large numbers of duplicates
 - Fallback option: use block nested loops join on overflowed partitions

Cost of Hash-Join



- If recursive partitioning is not required: cost of hash join is
$$3(b_r + b_s) + 4 * n_h$$
 block transfers +
$$2(\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil)$$
 seeks
- If recursive partitioning required:
 - number of passes required for partitioning build relation s to less than M blocks per partition is $\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_s/M) \rceil$
 - best to choose the smaller relation as the build relation.
 - Total cost estimate is:
$$2(b_r + b_s)\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_s/M) \rceil + b_r + b_s$$
 block transfers +
$$2(\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil)\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_s/M) \rceil$$
 seeks
- If the entire build input can be kept in main memory no partitioning is required
 - Cost estimate goes down to $b_r + b_s$.

Example of Cost of Hash-Join



instructor \bowtie *teaches*

- Assume that memory size is 20 blocks
- $b_{instructor} = 100$ and $b_{teaches} = 400$.
- *instructor* is to be used as build input. Partition it into five partitions, each of size 20 blocks. This partitioning can be done in one pass.
- Similarly, partition *teaches* into five partitions, each of size 80. This is also done in one pass.
- Therefore total cost, ignoring cost of writing partially filled blocks:
 - $3(100 + 400) = 1500$ block transfers +
 - $2(\lceil 100/3 \rceil + \lceil 400/3 \rceil) = 336$ seeks

Hybrid Hash–Join



- Useful when memory sizes are relatively large, and the build input is bigger than memory.

- **Main feature of hybrid hash join:**

Keep the first partition of the build relation in memory.

E.g. With memory size of 25 blocks, *instructor* can be partitioned into five partitions, each of size 20 blocks.

- Division of memory:
 - The first partition occupies 20 blocks of memory
 - 1 block is used for input, and 1 block each for buffering the other 4 partitions.
- *teaches* is similarly partitioned into five partitions each of size 80
 - the first is used right away for probing, instead of being written out
- Cost of $3(80 + 320) + 20 + 80 = 1300$ block transfers for hybrid hash join, instead of 1500 with plain hash-join.
- Hybrid hash-join most useful if $M \gg \sqrt{b_s}$

Complex Joins



- Join with a conjunctive condition:
 - $r \bowtie_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n} s$
 - Either use nested loops/block nested loops, or
 - Compute the result of one of the simpler joins $r \bowtie_{\theta_i} s$
 - final result comprises those tuples in the intermediate result that satisfy the remaining conditions
 - $\theta_1 \wedge \dots \wedge \theta_{i-1} \wedge \theta_{i+1} \wedge \dots \wedge \theta_n$
 - Join with a disjunctive condition
 - $r \bowtie_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n} s$
 - Either use nested loops/block nested loops, or
 - Compute as the union of the records in individual joins $r \bowtie_{\theta_i} s$:
 - $(r \bowtie_{\theta_1} s) \cup (r \bowtie_{\theta_2} s) \cup \dots \cup (r \bowtie_{\theta_n} s)$

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

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:
 - Partition both relations using the same hash function
 - 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.

Other Operations : Outer Join



- **Outer join** can be computed either as
 - A join followed by addition of null-padded non-participating tuples.
 - by modifying the join algorithms.
- Modifying merge join to compute $r \bowtie s$ ($\bowtie = \text{Left Outer Join}$)
 - In $r \bowtie s$, non participating tuples are those in $r - \Pi_R(r \bowtie s)$
 - Modify merge-join to compute $r \bowtie s$:
 - During merging, for every tuple t_r from r that do not match any tuple in s , output t_r padded with nulls.
 - Right outer-join and full outer-join can be computed similarly.

Other Operations : Outer Join



- Modifying hash join to compute $r \bowtie s$
 - If r is probe relation, output non-matching r tuples padded with nulls
 - If r is build relation, when probing keep track of which r tuples matched s tuples. At end of s , output non-matched r tuples padded with nulls

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.

Materialization

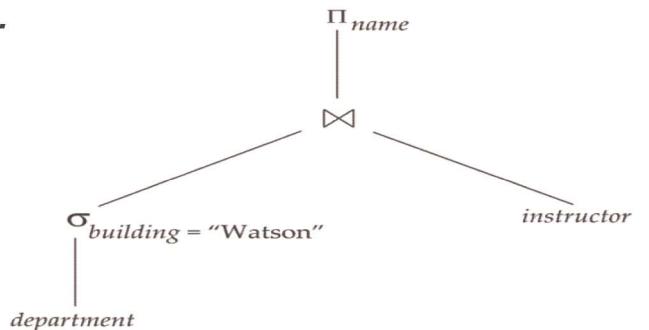


- **Materialized evaluation:** evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.

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

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

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



Materialization (Cont.)



- Materialized evaluation is always applicable (i.e. at any level of the expression tree)
- Cost of writing results to disk and reading them back can be quite high
 - Our cost formulas for operations ignore cost of writing results to disk, so
 - Overall cost = Sum of costs of individual operations + cost of writing intermediate results to disk
- **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 so as to reduce 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 **call-by-need** or **pull** evaluation
 - system repeatedly requests next tuple from top level (root) operation
 - Each operation requests next tuple from children operations as required, in order to output its next tuple
 - In between calls, operation has to maintain “**state**” so it knows what to return next
- In **producer-driven** or **eager** or **push** 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
- Alternative name: **pull** and **push** models of pipelining

Pipelining (Cont.)



- Implementation of demand-driven pipelining
 - Each operation is implemented as an **iterator** (an object to loop through collections like arrays etc.) 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()**

Evaluation Algorithms for Pipelining



- Some algorithms are not able to output results even as they get input tuples
 - E.g. merge join, or hash join
 - intermediate results written to disk and then read back
- Algorithm variants to generate (at least some) results on the fly, as input tuples are read in
 - E.g. hybrid hash join generates output tuples even as probe relation tuples in the in-memory partition (partition 0) are read in
 - **Double-pipelined join technique:** Hybrid hash join, modified to buffer partition 0 tuples of both relations in-memory, reading them as they become available, and output results of any matches between partition 0 tuples
 - When a new r_0 tuple is found, match it with existing s_0 tuples, output matches, and save it in r_0
 - Symmetrically for s_0 tuples

Query Optimization



Query Optimization

Query Optimization: Topics



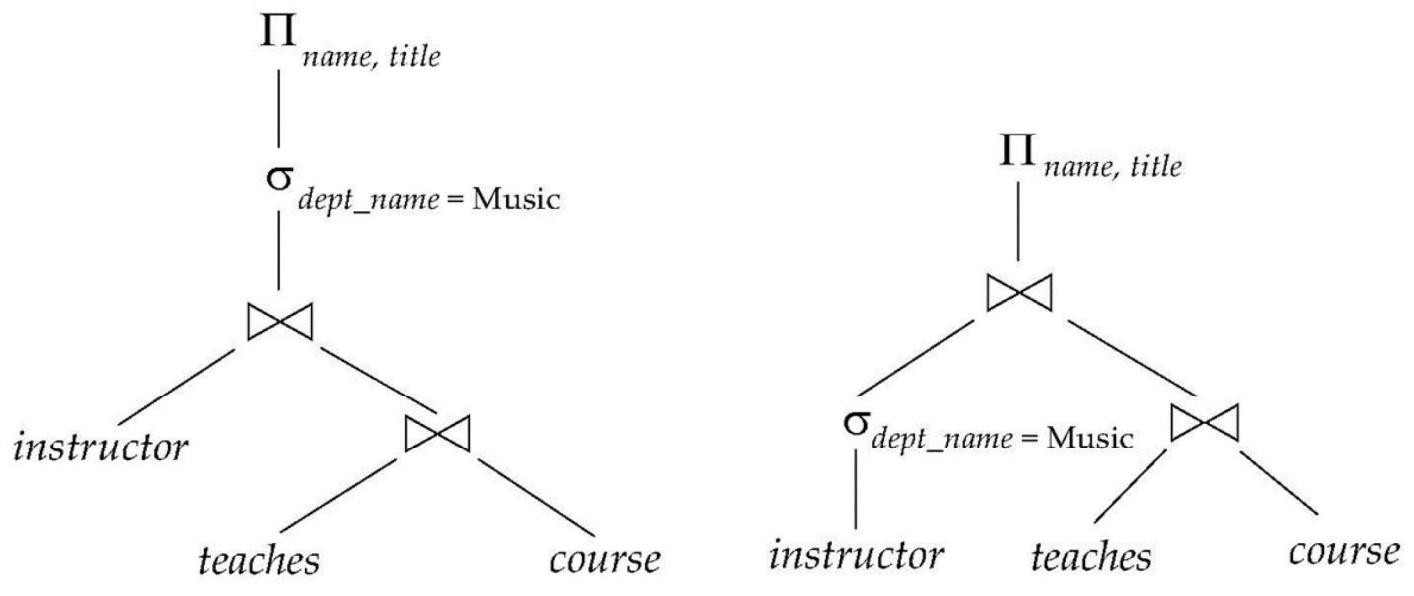
- Introduction
 - Transformation of Relational Expressions
 - Catalog Information for Cost Estimation
 - Statistical Information for Cost Estimation
 - Cost-based optimization
 - Dynamic Programming for Choosing Evaluation Plans
 - Materialized views
-

Introduction



Alternative ways of evaluating a given query

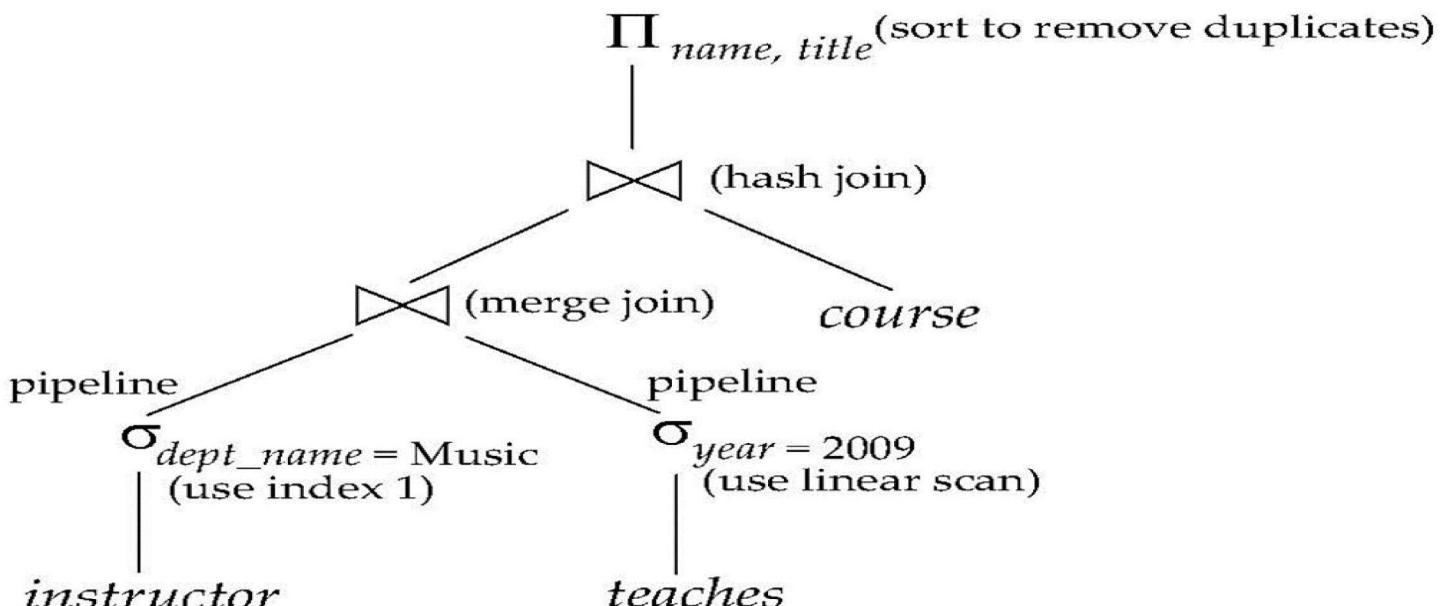
- Equivalent expressions
- Different algorithms for each operation



Introduction (cont.)



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





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**
 - Annotate resultant expressions to get alternative query plans
 - 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

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 (i.e. allows for multiple instances for each of its elements) 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.

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

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 θ_2 involves attributes from only E_2 and E_3 .



Equivalence Rules (cont.)

7. The selection operation distributes over the theta join operation under the following two conditions:
 - (a) When all the attributes in θ_0 involve only the attributes of one of the expressions (E_1) being joined.

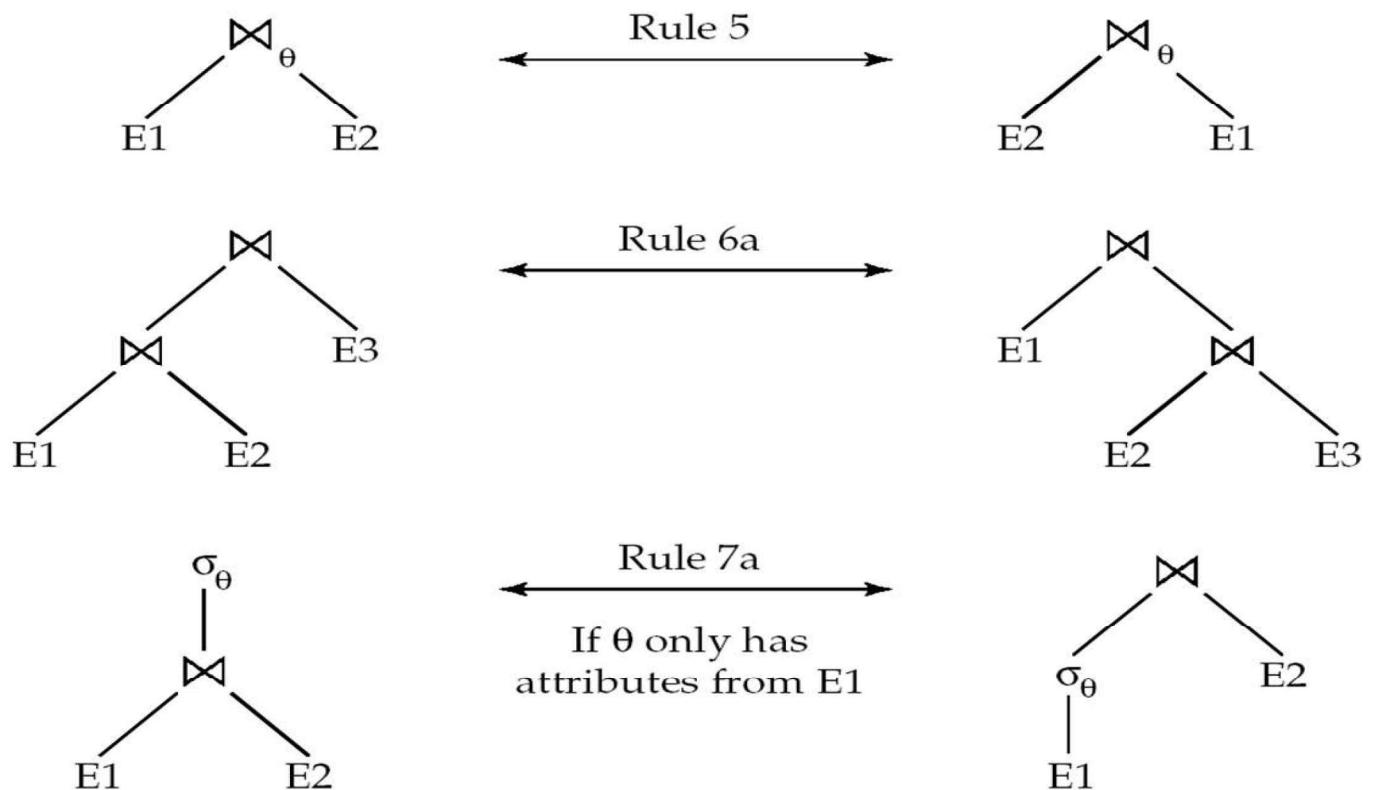
$$\sigma_{\theta_0}(E_1 \bowtie_{\theta} E_2) = (\sigma_{\theta_0}(E_1)) \bowtie_{\theta} E_2$$

- (b) When θ_1 involves only the attributes of E_1 and θ_2 involves only the attributes of E_2 .

$$\sigma_{\theta_1 \wedge \theta_2}(E_1 \bowtie_{\theta} E_2) = (\sigma_{\theta_1}(E_1)) \bowtie_{\theta} (\sigma_{\theta_2}(E_2))$$



Pictorial Depiction of Equivalence Rules



Equivalence Rules (cont.)



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

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

$$\Pi_{L_1 \cup L_2} (E_1 \bowtie_{\theta} E_2) = (\Pi_{L_1}(E_1)) \bowtie_{\theta} (\Pi_{L_2}(E_2))$$

(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 θ , but are not in $L_1 \cup L_2$, and
- let L_4 be attributes of E_2 that are involved in join condition θ , but are not in $L_1 \cup L_2$.

$$\Pi_{L_1 \cup L_2} (E_1 \bowtie_{\theta} E_2) = \Pi_{L_1 \cup L_2} ((\Pi_{L_1 \cup L_3}(E_1)) \bowtie_{\theta} (\Pi_{L_2 \cup L_4}(E_2)))$$



Equivalence Rules (cont.)

9. The set operations union and intersection are commutative

$$\begin{aligned} E_1 \cup E_2 &= E_2 \cup E_1 \\ E_1 \cap E_2 &= E_2 \cap E_1 \end{aligned}$$

■ (set difference is not commutative).

10. Set union and intersection are associative.

$$\begin{aligned} (E_1 \cup E_2) \cup E_3 &= E_1 \cup (E_2 \cup E_3) \\ (E_1 \cap E_2) \cap E_3 &= E_1 \cap (E_2 \cap E_3) \end{aligned}$$

11. The selection operation distributes over \cup , \cap and $-$.

$$\begin{aligned} \sigma_0(E_1 - E_2) &= \sigma_0(E_1) - \sigma_0(E_2) \\ \text{and similarly for } \cup \text{ and } \cap \text{ in place of } - \end{aligned}$$

Also: $\sigma_0(E_1 - E_2) = \sigma_0(E_1) - E_2$

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
 - $\Pi_{name, title}(\sigma_{dept_name= \text{`Music'}}(instructor \bowtie (teaches \bowtie \Pi_{course_id, title} (course))))$
- **Transformation using rule 7a.**
 - $\Pi_{name, title}((\sigma_{dept_name= \text{`Music'}}(instructor)) \bowtie (teaches \bowtie \Pi_{course_id, title} (course)))$
- 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

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

Transformation using join associatively (Rule 6a):

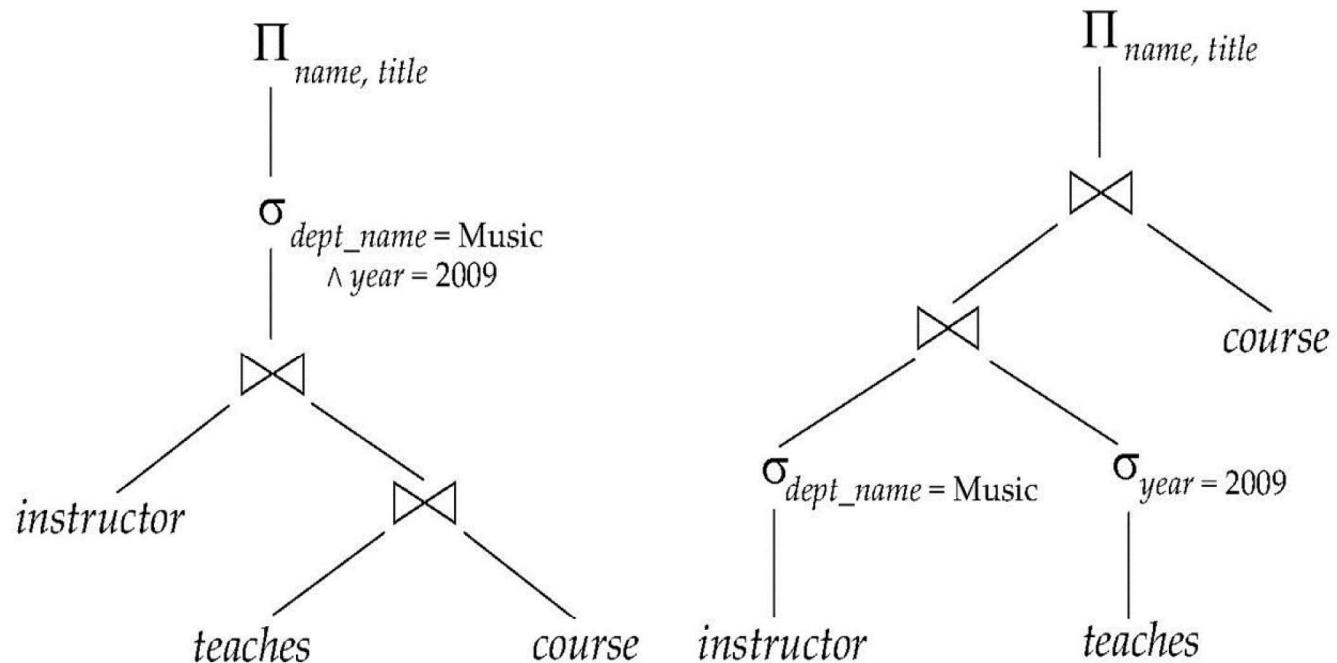
- $\Pi_{name, title}(\sigma_{dept_name = "Music"} \wedge gear = 2009 ((instructor \bowtie teaches) \bowtie \Pi_{course_id, title} (course)))$

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

$$\sigma_{dept_name = "Music"} (instructor) \bowtie \sigma_{year = 2009} (teaches)$$



Multiple Transformations (Cont.)



(a) Initial expression tree

(b) Tree after multiple transformations



Transformation Example: Pushing Projections

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

When we compute

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

- 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:
$$\Pi_{name, title}(\Pi_{name, course_id}(\sigma_{dept_name = "Music"}(instructor \bowtie teaches)) \bowtie \Pi_{course_id, title}(course)))$$
- 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 ,
$$(r_1 \bowtie r_2) \bowtie r_3 = r_1 \bowtie (r_2 \bowtie r_3)$$
(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

$$\begin{aligned} & \Pi_{name, title}(\sigma_{dept_name= ‘Music’}(instructor) \bowtie teaches) \\ & \bowtie \Pi_{course_id, title}(course))) \end{aligned}$$

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

$$\sigma_{dept_name= ‘Music’}(instructor)$$

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

$$\sigma_{dept_name= ‘Music’}(instructor) \bowtie teaches$$

first.

Enumeration of Equivalent Expressions

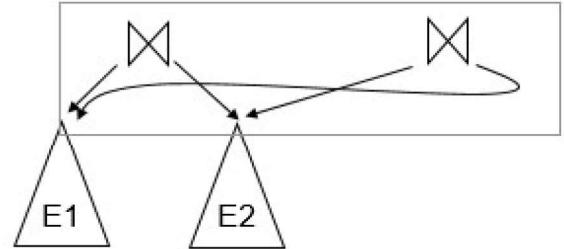


- Query optimizers use equivalence rules to **systematically** generate expressions equivalent to the given expression
- Can generate all equivalent expressions as follows:
 - Repeat
 - apply all applicable equivalence rules on every subexpression of every equivalent expression found so far
 - add newly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

Implementing Transformation Based Optimization



- Space requirements reduced by sharing common sub-expressions:
 - when E1 is generated from E2 by an equivalence rule, usually only the top level of the two are different, subtrees below are the same and can be shared using pointers
 - E.g. when applying join commutativity



- Same sub-expression may get generated multiple times
 - Detect duplicate sub-expressions and share one copy

Time requirements are reduced by not generating all expressions

- Dynamic programming
 - We will discuss only the special case of dynamic programming for join order optimization

Cost Estimation



- Cost of each operator computed as described in earlier
 - 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, histogram etc.
- More on cost estimation later

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.

Cost-Based Optimization



- Consider finding the best join-order for $r_1 \bowtie r_2 \bowtie \dots r_n$.
- For bushy binary tree plane, 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! [For bushy plan: $(2(n - 1))!/(n - 1)!$, for deep plans may be less say $n!$ for left deep plan i.e. join operators on left side]
- No need to generate all the join orders. Using dynamic programming, the least-cost join order for any subset of $\{r_1, r_2, \dots r_n\}$ is computed only once and stored for future use.



Dynamic Programming

- Save the result for future reference, so as to avoid solving the same problem again.
- The optimal solutions to the sub-problems contribute to the optimal solution of the given problem.

There are two ways of doing this.

- 1.) **Top-Down** : Start solving the given problem by breaking it down. If you see that the problem has been solved already, then just return the saved answer. If it has not been solved, solve it and save the answer. This is usually easy to think of and very intuitive. This is referred to as *Memoization*.
- 2.) **Bottom-Up** : Analyze the problem and see the order in which the sub-problems are solved and start solving from the trivial sub-problem, up towards the given problem. In this process, it is guaranteed that the sub-problems are solved before solving the problem. This is referred to as *Dynamic Programming*.

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



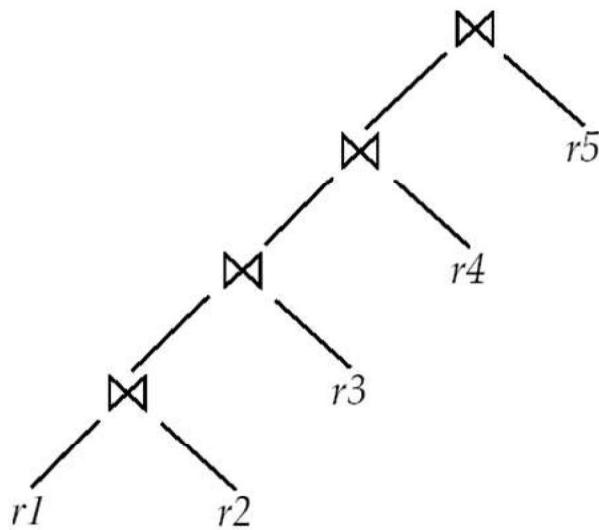
Join Order Optimization Algorithm

```
procedure findbestplan(S)
    if (bestplan[S].cost ≠ ∞)
        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 ≠ 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]
```

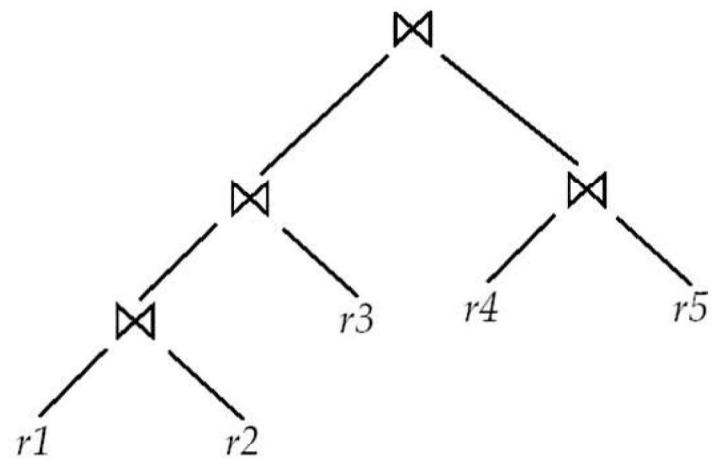
Left Deep Join Trees



- 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 S_1 of S such that $S_1 \neq S$ ”
 - By: **for each** relation r in S
let $S_1 = S - r$.
- 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 does not affect time/space complexity significantly

Cost Based Optimization with Equivalence Rules



- **Physical equivalence rules** allow logical query plan to be converted to physical query plan specifying what algorithms are used for each operation.
- Efficient optimizer based on equivalent rules depends on
 - A space efficient representation of expressions which avoids making multiple copies of subexpressions
 - Efficient techniques for detecting duplicate derivations of expressions
 - A form of dynamic programming based on **memoization**, which stores the best plan for a subexpression the first time it is optimized, and reuses it on repeated optimization calls on same subexpression
 - Cost-based pruning techniques that avoid generating all plans
- Pioneered by the Volcano project and implemented in the SQL Server optimizer

Heuristic Optimization



- Cost-based optimization is expensive, even with dynamic programming.
- Systems may use *heuristics* to reduce the number of choices that must be made in a cost-based fashion.
- Heuristic optimization transforms the query-tree by using a set of rules that typically (but not in all cases) improve execution performance:
 - Perform selection early (reduces the number of tuples)
 - Perform projection early (reduces the number of attributes)
 - Perform most restrictive selection and join operations (i.e. with smallest result size) before other similar operations.
 - Some systems use only heuristics, others combine heuristics with partial cost-based optimization.



Structure of Query Optimizers

- Many optimizers consider only left-deep join orders.
 - Plus heuristics to push selections and projections down the query tree
 - Reduces optimization complexity and generates plans amenable to pipelined evaluation.
- Heuristic optimization used in some versions of Oracle:
 - Repeatedly pick “best” relation to join next
 - Starting from each of n starting points. Pick best among these
- Intricacies (elaborately complex detail) of SQL complicate query optimization
 - E.g. nested subqueries

Structure of Query Optimizers (Cont.)



- Some query optimizers integrate heuristic selection and the generation of alternative access plans.
 - Frequently used approach
 - heuristic rewriting of nested block structure and aggregation
 - followed by cost-based join-order optimization for each block
 - Some optimizers (e.g. SQL Server) apply transformations to entire query and do not depend on block structure
 - **Optimization cost budget** to stop optimization early (if cost of plan is less than cost of optimization)
 - **Plan caching** to reuse previously computed plan if query is resubmitted
 - Even with different constants in query
- Even with the use of heuristics, cost-based query optimization imposes a substantial overhead.
 - But is worth it for expensive queries
 - Optimizers often use simple heuristics for very cheap queries, and perform exhaustive enumeration for more expensive queries



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 $\Pi_A(r)$.

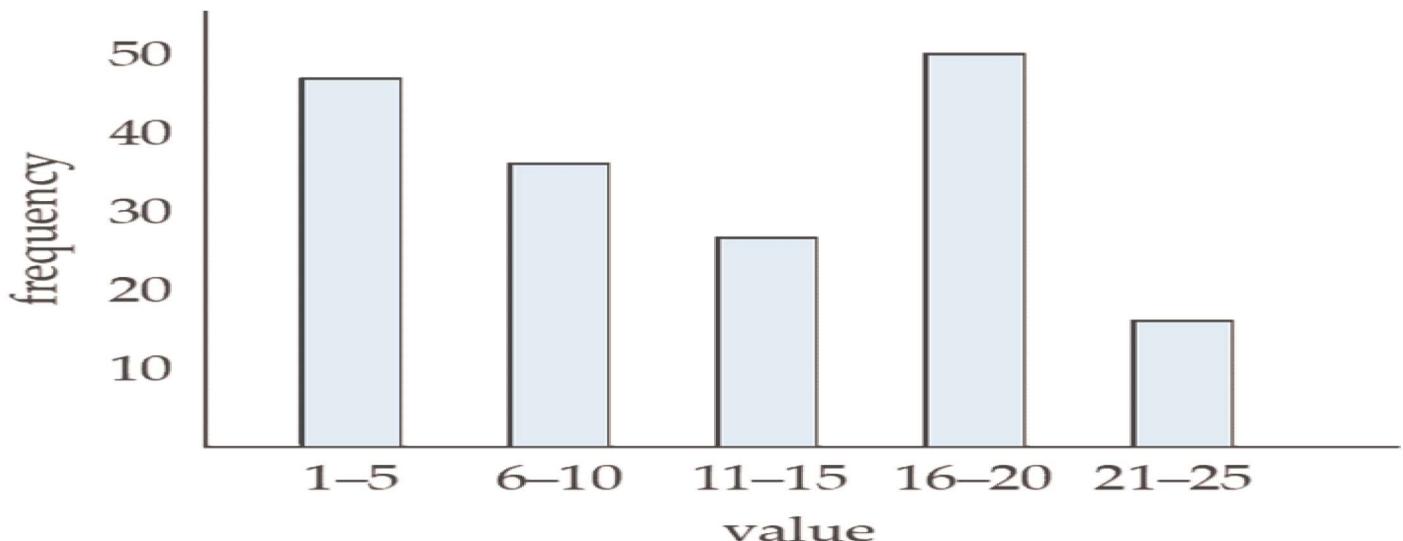
If tuples of r are stored together physically in a file, then:

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



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 $v < \min(A, r)$
 - $c = n_r \cdot \frac{v - \min(A, r)}{\max(A, r) - \min(A, r)}$
 - 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 θ_i is the probability that a tuple in the relation r satisfies θ_i .

- If s_i is the number of satisfying tuples in r , the selectivity of θ_i is given by s_i/n_r

Conjunction: $\sigma_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n}(r)$. Assuming independence, estimate of tuples in the result is:

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

Disjunction: $\sigma_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n}(r)$. Estimated number of tuples:

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

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

$$n_r - \text{size}(\sigma_\theta(r))$$



Join Operation: Running Example

- Running example:
 $student \quad takes$

Catalog information for join examples:

$$n_{student} = 5,000.$$

$$f_{student} = 50, \text{ which implies that } b_{student} = 5000/50 = 100.$$

$$n_{takes} = 10000.$$

$$f_{takes} = 25, \text{ which implies that } b_{takes} = 10000/25 = 400.$$

$V(ID, takes) = 2500$, which implies that on average, each student who has taken a course has taken 4 courses.

- Attribute ID in $takes$ is a foreign key referencing $student$.
- $V(ID, student) = 5000$ (*primary key!*)

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.
-

Size Estimation for Other Operations



Projection: estimated size of $\prod_A(r) = V(A,r)$

Aggregation : estimated size of ${}_A g_F(r) = V(A,r)$

Set operations

- For unions/intersections of selections on the same relation: rewrite and use size estimate for selections
 - E.g. $\sigma_{\theta_1}(r) \cup \sigma_{\theta_2}(r)$ can be rewritten as $\sigma_{\theta_1 \vee \theta_2}(r)$
- For operations on different relations:
 - estimated size of $r \cup s = \text{size of } r + \text{size of } s.$
 - estimated size of $r \cap s = \text{minimum size of } r \text{ and size of } s.$
 - estimated size of $r - s = r.$
 - All the three estimates may be quite inaccurate, but provide upper bounds on the sizes.

Size Estimation (Cont.)



Outer join:

- Estimated size of $r \bowtie s = \text{size of } r \bowtie s + \text{size of } r$
 - Case of right outer join is symmetric
- Estimated size of $r \bowtie s = \text{size of } r \bowtie s + \text{size of } r + \text{size of } s$

Estimation of Number of Distinct Values



Selections: $\sigma_{\theta}(r)$

- If θ forces A to take a specified value: $V(A, \sigma_{\theta}(r)) = 1$.
 - e.g., $A = 3$
- If θ forces A to take on one of a specified set of values:
 $V(A, \sigma_{\theta}(r)) = \text{number of specified values.}$
 - (e.g., $(A = 1 \vee A = 3 \vee A = 4)$),
- If the selection condition θ is of the form $A \text{ op } r$
estimated $V(A, \sigma_{\theta}(r)) = V(A.r) * s$
 - where s is the selectivity of the selection.
- In all the other cases: use approximate estimate of
 $\min(V(A, r), n_{\sigma\theta(r)})$
 - More accurate estimate can be got using probability theory, but this one works fine generally

Estimation of Number of Distinct Values (cont.)



Joins: $r \bowtie s$

- If all attributes in A are from r
estimated $V(A, r \bowtie s) = \min(V(A, r), n_{r \bowtie s})$
- If A contains attributes A_1 from r and A_2 from s , then
estimated
 $V(A, r \bowtie s) =$
 $\min(V(A_1, r) * V(A_2 - A_1, s), V(A_1 - A_2, r) * V(A_2, s), n_{r \bowtie s})$
 - More accurate estimate can be got using probability theory, but this one works fine generally

Estimation of Distinct Values (Cont.)



- Estimation of distinct values are straightforward for projections.
 - They are the same in $\Pi_{A(r)}$ as in r .
- The same holds for grouping attributes of aggregation.
- For aggregated values
 - For $\min(A)$ and $\max(A)$, the number of distinct values can be estimated as $\min(V(A,r), V(G,r))$ where G denotes grouping attributes
 - For other aggregates, assume all values are distinct, and use $V(G,r)$



Additional Optimization Techniques

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
    - teaches can have duplicate IDs
    - Can be modified to handle duplicates correctly
  - 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 general, SQL queries of the form below can be rewritten as shown

Rewrite: **select ... from  $L_1$**

**where  $P_1$  and exists (select \***  
**from  $L_2$  where  $P_2$ )**

To: **create table  $t_1$  as**

**select distinct  $V$**   
**from  $L_2$**   
**where  $P_2^1$**

**select ...**  
**from  $L_1, t_1$**   
**where  $P_1$  and  $P_2^2$**

- $P_2^1$  contains predicates in  $P_2$  that do not involve any correlation variables
- $P_2^2$  reintroduces predicates involving correlation variables, with relations renamed appropriately
- $V$  contains all attributes used in predicates with correlation variables

# Optimizing Nested Subqueries (Cont.)



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

```
create table t1 as
 select distinct ID
 from teaches
 where year = 2007
```

```
select name
from instructor, t1
where t1.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.

# Materialized Views



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

```
create 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

View maintenance can be done by

- Manually defining triggers on insert, delete, and update of each relation in the view definition
- Manually written code to update the view whenever database relations are updated
- Periodic recomputation (e.g. nightly)
- Above methods are directly supported by many database systems

- Avoids manual effort/correctness issues

# Incremental View Maintenance

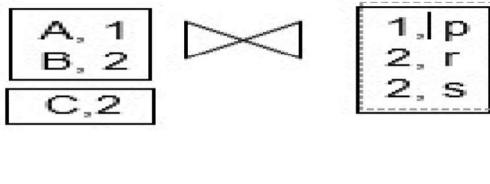


- The changes (inserts and deletes) to a relation or expressions are referred to as its **differential**
  - Set of tuples inserted to and deleted from  $r$  are denoted  $i_r$  and  $d_r$
- To simplify our description, we only consider inserts and deletes
  - We replace updates to a tuple by deletion of the tuple followed by insertion of the update tuple
- We discuss how to compute the change to the result of each relational operation, given changes to its inputs
- We then discuss how to handle relational algebra expressions

# Join Operation



- Consider the materialized view  $v = r \bowtie s$  and an update to  $r$ .
- Let  $r^{old}$  and  $r^{new}$  denote the old and new states of relation  $r$
- Consider the case of an insert to  $r$ :
  - We can write  $r^{new} \bowtie s$  as  $(r^{old} \cup i_r) \bowtie s$
  - And rewrite the above to  $(r^{old} \bowtie s) \cup (i_r \bowtie s)$
  - But  $(r^{old} \bowtie s)$  is simply the old value of the materialized view, so the incremental change to the view is just  $i_r \bowtie s$
- Thus, for inserts  $v^{new} = v^{old} \cup (i_r \bowtie s)$
- Similarly for deletes  $v^{new} = v^{old} - (d_r \bowtie s)$



# Selection and Projection Operations



Selection: Consider a view  $v = \sigma_{\theta}(r)$ .

- $v^{new} = v^{old} \cup \sigma_{\theta}(i_r)$
- $v^{new} = v^{old} - \sigma_{\theta}(d_r)$

Projection is a more difficult operation

- $R = (A, B)$ , and  $r(R) = \{(a, 2), (a, 3)\}$
- $\Pi_A(r)$  has a single tuple  $(a)$ .
- If we delete the tuple  $(a, 2)$  from  $r$ , we should not delete the tuple  $(a)$  from  $\Pi_A(r)$ , but if we then delete  $(a, 3)$  as well, we should delete the tuple

For each tuple in a projection  $\Pi_A(r)$ , we will keep a count of how many times it was derived

- On insert of a tuple to  $r$ , if the resultant tuple is already in  $\Pi_A(r)$  we increment its count, else we add a new tuple with count = 1
- On delete of a tuple from  $r$ , we decrement the count of the corresponding tuple in  $\Pi_A(r)$ 
  - if the count becomes 0, we delete the tuple from  $\Pi_A(r)$



# Aggregation Operations

- count :  $v = {}_A g_{count}(B)^{(r)}$ .
  - When a set of tuples  $i_r$  is inserted
    - For each tuple  $r$  in  $i_r$ , if the corresponding group is already present in  $v$ , we increment its count, else we add a new tuple with count = 1
  - When a set of tuples  $d_r$  is deleted
    - for each tuple  $t$  in  $i_r$ , we look for the group  $t.A$  in  $v$ , and subtract 1 from the count for the group.
      - If the count becomes 0, we delete from  $v$  the tuple for the group  $t.A$
- sum:  $v = {}_A g_{sum}(B)^{(r)}$ 
  - We maintain the sum in a manner similar to count, except we add/subtract the B value instead of adding/subtracting 1 for the count
  - Additionally we maintain the count in order to detect groups with no tuples. Such groups are deleted from  $v$ 
    - Cannot simply test for sum = 0 (why?)
- To handle the case of **avg**, we maintain the **sum** and **count** aggregate values separately, and divide at the end



# Aggregate Operations (Cont.)

- **min, max:**  $v = \text{Ag}_{\min(B)}(r)$ .
  - Handling insertions on  $r$  is straightforward.
  - Maintaining the aggregate values **min** and **max** on deletions may be more expensive. We have to look at the other tuples of  $r$  that are in the same group to find the new minimum

# Other Operations



- Set intersection:  $v = r \cap s$ 
  - when a tuple is inserted in  $r$  we check if it is present in  $s$ , and if so we add it to  $v$ .
  - If the tuple is deleted from  $r$ , we delete it from the intersection if it is present.
  - Updates to  $s$  are symmetric
  - The other set operations, *union* and *set difference* are handled in a similar fashion.
- Outer joins are handled in much the same way as joins but with some extra work

# Handling Expressions



- To handle an entire expression, we derive expressions for computing the incremental change to the result of each sub-expressions, starting from the smallest sub-expressions.
- E.g. consider  $E_1 \bowtie E_2$  where each of  $E_1$  and  $E_2$  may be a complex expression
  - Suppose the set of tuples to be inserted into  $E_1$  is given by  $D_1$ 
    - Computed earlier, since smaller sub-expressions are handled first
  - Then the set of tuples to be inserted into  $E_1 \bowtie E_2$  is given by  $D_1 \bowtie E_2$ 
    - This is just the usual way of maintaining joins

# Query Optimization and Materialized Views



- Rewriting queries to use materialized views:
  - A materialized view  $v = r \bowtie s$  is available
  - A user submits a query  $r \bowtie s \bowtie t$
  - We can rewrite the query as  $v \bowtie t$ 
    - Whether to do so depends on cost estimates for the two alternative
- Replacing a use of a materialized view by the view definition:
  - A materialized view  $v = r \bowtie s$  is available, but without any index on it
  - User submits a query  $\sigma_{A=10}(v)$ .
  - Suppose also that  $s$  has an index on the common attribute B, and  $r$  has an index on attribute A.
  - The best plan for this query may be to replace  $v$  by  $r \bowtie s$ , which can lead to the query plan  $\sigma_{A=10}(r) \bowtie s$
- Query optimizer should be extended to consider all above alternatives and choose the best overall plan



# 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 can be found in the manuals of the commercial software.
- 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**

---



## **Additional Optimization Techniques**

---



# Top-K Queries

## Top-K queries

```
select *
from r, s
where r.B = s.B
order by r.A ascending
limit 10
```

- Alternative 1: Indexed nested loops join with r as outer
- Alternative 2: estimate highest r.A value in result and add selection (**and**  $r.A \leq H$ ) to where clause
  - If < 10 results, retry with larger H

# Optimization of Updates



## Halloween problem

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

- If index on A is used to find tuples satisfying  $A > 10$ , and tuples updated immediately, same tuple may be found (and updated) multiple times
- Solution 1: *Always defer updates*
  - collect the updates (old and new values of tuples) and update relation and indices in second pass
  - Drawback: extra overhead even if e.g. update is only on R.B, not on attributes in selection condition
- Solution 2: *Defer only if required*
  - Perform immediate update if update does not affect attributes in where clause, and deferred updates otherwise.



# Join Minimization

## Join minimization

**select r.A, r.B from r, s where r.B = s.B**

- Check if join with s is redundant, drop it
  - E.g. join condition is on foreign key from r to s, r.B is declared as not null, and no selection on s
  - Other sufficient conditions possible  
**select r.A, s2.B  
from r, s as s1, s as s2  
where r.B=s1.B and r.B = s2.B and s1.A < 20 and s2.A < 10**
    - join with s1 is redundant and can be dropped (along with selection on s1)
    - Lots of research in this area since 70s/80s!

# Multiquery Optimization



- Example

Q1: **select \* from (r natural join t) natural join s**

Q2: **select \* from (r natural join u) natural join s**

- Both queries share common subexpression ( $r$  natural join  $s$ )
- May be useful to compute ( $r$  natural join  $s$ ) once and use it in both queries
  - But this may be more expensive in some situations
    - e.g. ( $r$  natural join  $s$ ) may be expensive, plans as shown in queries may be cheaper
- **Multiquery optimization:** find best overall plan for a set of queries, exploiting sharing of common subexpressions between queries where it is useful

# Multiquery Optimization (Cont.)



- Simple heuristic used in some database systems:
  - optimize each query separately
  - detect and exploiting common subexpressions in the individual optimal query plans
    - May not always give best plan, but is cheap to implement
  - **Shared scans:** widely used special case of multiquery optimization
- Set of materialized views may share common subexpressions
  - As a result, view maintenance plans may share subexpressions
  - Multiquery optimization can be useful in such situations



# Parametric Query Optimization

- Example
  - select \* from r natural join s where r.a < \$1
    - value of parameter \$1 not known at compile time
      - known only at run time
      - different plans may be optimal for different values of \$1
    - Solution 1: optimize at run time, each time query is submitted
      - can be expensive
  - Solution 2: **Parametric Query Optimization:**
    - optimizer generates a set of plans, optimal for different values of \$1
      - Set of optimal plans usually small for 1 to 3 parameters
      - Key issue: how to do find set of optimal plans efficiently
    - best one from this set is chosen at run time when \$1 is known
  - Solution 3: **Query Plan Caching**
    - If optimizer decides that same plan is likely to be optimal for all parameter values, it caches plan and reuses it, else re-optimize each time
    - Implemented in many database systems

# Plan Stability Across Optimizer Changes



- What if 95% of plans are faster on database/optimizer version N+1 than on N, but 5% are slower?
  - Why should plans be slower on new improved optimizer?
    - Answer: Two wrongs can make a right, fixing one wrong can make things worse!
- Approaches:
  - Allow hints for tuning queries
    - Not practical for migrating large systems with no access to source code
  - Set optimization level, default to version N (Oracle)
    - And migrate one query at a time after testing both plans on new optimizer
  - Save plan from version N, and give it to optimizer version N+1
    - Sybase, XML representation of plans (SQL Server)



**BITS** Pilani, Pilani Campus



# Thanks ....

---

Any questions please .....

# Thanking you