

8. Query Processing

Goals: Understand the basic concepts underlying the steps in query processing and optimization and estimating query processing cost; apply query optimization techniques;

Contents:

- Overview
- Catalog Information for Cost Estimation
- Measures of Query Cost
- Selection
- Join Operations
- Other Operations
- Evaluation and Transformation of Expressions

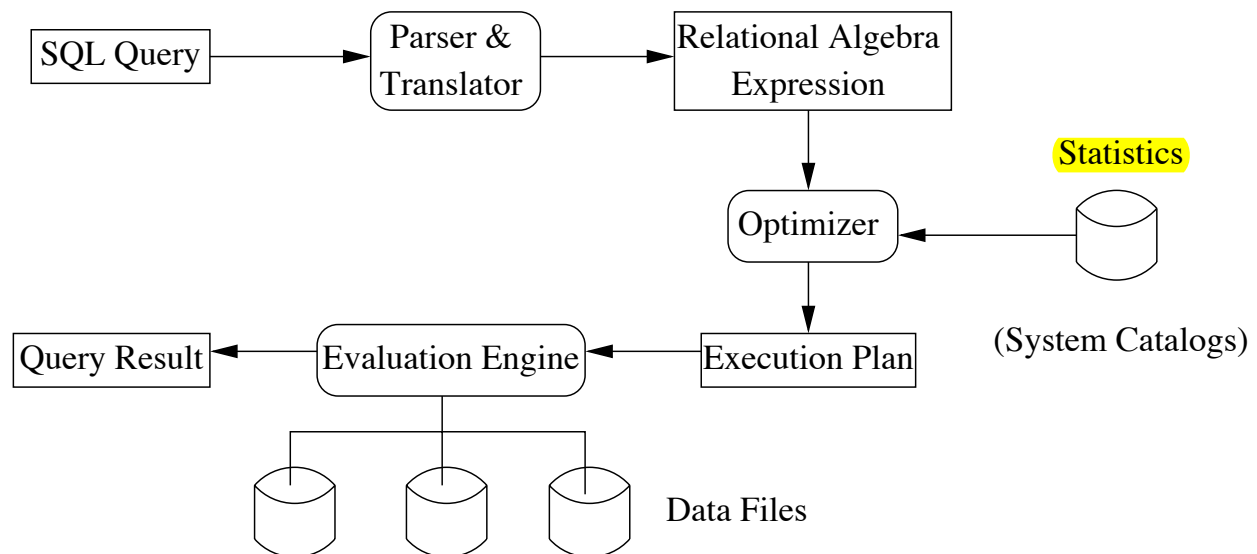
Query Processing & Optimization

Task: Find an efficient *physical query plan* (aka *execution plan*) for an SQL query

Goal: Minimize the *evaluation time* for the query, i.e., compute query result as fast as possible

Cost Factors: Disk accesses, read/write operations, [I/O, page transfer] (CPU time is typically ignored)

Basic Steps in Processing an SQL Query



- **Parsing and Translating**

- Translate the query into its internal form (parse tree). This is then translated into an expression of the relational algebra.
- Parser checks syntax, validates relations, attributes and access permissions

- **Evaluation**

- The query execution engine takes a physical query plan (aka execution plan), executes the plan, and returns the result.

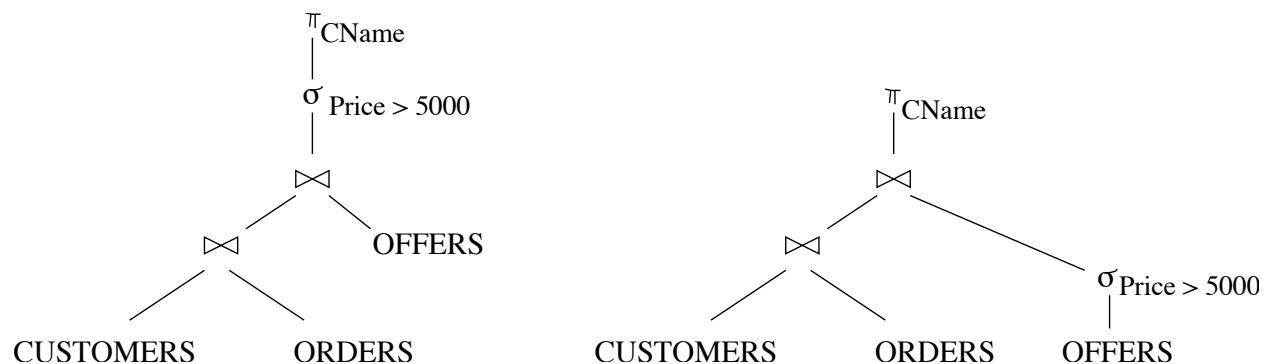
- **Optimization:** Find the “cheapest” execution plan for a query

- A relational algebra expression may have many equivalent expressions, e.g.,

$$\pi_{CName}(\sigma_{Price > 5000}((CUSTOMERS \bowtie ORDERS) \bowtie OFFERS))$$

$$\pi_{CName}((CUSTOMERS \bowtie ORDERS) \bowtie (\sigma_{Price > 5000}(OFFERS)))$$

Representation as *logical query plan* (a tree):



Non-leaf nodes \equiv operations of relational algebra (with parameters); Leaf nodes \equiv relations

- A relational algebra expression can be evaluated in many ways. An annotated expression specifying detailed evaluation strategy is called the **execution plan** (includes, e.g., whether index is used, join algorithms, . . .)
- Among all semantically equivalent expressions, the one with the least costly evaluation plan is chosen. Cost estimate of a plan is based on **statistical information** in the system catalogs.

Catalog Information for Cost Estimation

Information about relations and attributes:

- N_R : number of tuples in the relation R.
- B_R : number of blocks that contain tuples of the relation R.
- S_R : size of a tuple of R.
- F_R : blocking factor; number of tuples from R that fit into one block ($F_R = \lceil N_R/B_R \rceil$)
- $V(A, R)$: number of distinct values for attribute A in R.
- $SC(A, R)$: selectivity of attribute A
 \equiv average number of tuples of R that satisfy an equality condition on A.
 $SC(A, R) = N_R/V(A, R)$.

Information about indexes:

- HT_I : number of levels in index I (B^+ -tree). *height*
- LB_I : number of blocks occupied by leaf nodes in index I (first-level blocks).
- Val_I : number of distinct values for the search key.

Some relevant tables in the Oracle system catalogs:

USER_TABLES	USER_TAB_COLUMNS	USER_INDEXES
NUM_ROWS	NUM_DISTINCT	BLEVEL
BLOCKS	LOW_VALUE	LEAF_BLOCKS
EMPTY_BLOCKS	HIGH_VALUE	DISTINCT_KEYS
AVG_SPACE	DENSITY	AVG_LEAF_BLOCKS_PER_KEY
CHAIN_CNT	NUM_BUCKETS	
AVG_ROW_LEN	LAST_ANALYZED	

Measures of Query Cost

- There are many possible ways to estimate cost, e.g., based on disk accesses, CPU time, or communication overhead.
- Disk access is the predominant cost (in terms of time); relatively easy to estimate; therefore, number of block transfers from/to disk is typically used as measure.
 - Simplifying assumption: each **block** transfer has the same cost.
- Cost of algorithm (e.g., for join or selection) depends on database **buffer size**; more memory for DB buffer reduces disk accesses. Thus DB buffer size is a parameter for estimating cost.
- We refer to the cost estimate of algorithm S as **cost(S)**. We do not consider cost of writing output to disk.

Selection Operation

$\sigma_{A=a}(R)$ where a is a constant value, A an attribute of R

- *File Scan* – search algorithms that locate and retrieve records that satisfy a selection condition
- **S1 – Linear search**

$$\text{cost}(S1) = B_R$$

Selection Operation (cont.)

- **S2 – Binary search**, i.e., the file ordered based on attribute A (primary index)

$$\text{cost}(S2) = \lceil \log_2(B_R) \rceil + \left\lceil \frac{SC(A, R)}{F_R} \right\rceil - 1$$

- $\lceil \log_2(B_R) \rceil \equiv$ cost to locate the first tuple using binary search
- Second term \equiv blocks that contain records satisfying the selection.
- If A is primary key, then $SC(A, R) = 1$, hence $\text{cost}(S2) = \lceil \log_2(B_R) \rceil$.

- **Example** (for Employee DB)

- $F_{\text{Employee}} = 10$;
 $V(\text{Deptno}, \text{Employee}) = 50$ (different departments)
 - $N_{\text{Employee}} = 10,000$ (Relation Employee has 10,000 tuples)
 - Assume selection $\sigma_{\text{Deptno}=20}(\text{Employee})$ and Employee is sorted on search key Deptno :
 - $\implies 10,000/50 = 200$ tuples in Employee belong to Deptno 20;
(assuming an equal distribution)
 $200/10 = 20$ blocks for these tuples
 - \implies A binary search finding the first block would require $\lceil \log_2(1,000) \rceil = 10$ block accesses
- Total cost of binary search is 10+20 block accesses (versus 1,000 for linear search and Employee not sorted by Deptno).

- **Index scan** – search algorithms that use an index (here, a B^+ -tree); selection condition is on search key of index
- **S3** – **Primary index** I for A , **A primary key**, equality $A = a$
 $\text{cost}(S3) = HT_I + 1$ (only 1 tuple satisfies condition)

- **S4** – **Primary index** I on **non-key** A equality $A = a$

$$\text{cost}(S4) = HT_I + \left\lceil \frac{SC(A, R)}{F_R} \right\rceil$$

- **S5** – Non-primary (**non-clustered**) index on **non-key** A , equality $A = a$

$$\text{cost}(S5) = HT_I + SC(A, R)$$

Worst case: each matching record resides in a different block.

- **Example (Cont.):**

- Assume primary (B^+ -tree) index for attribute Deptno
- $200/10=20$ blocks accesses are required to read Employee tuples
- If B^+ -tree index stores 20 pointers per (inner) node, then the B^+ -tree index must have between **3 and 5** leaf nodes and the entire tree has a depth of 2
 \implies a total of 22 blocks must be read.

Selections Involving Comparisons

- Selections of the form $\sigma_{A \leq v}(R)$ or $\sigma_{A \geq v}(R)$ are implemented using a file scan or binary search, or by using either a
 - **S6** – A **primary index** on A , or
 - **S7** – A **secondary index** on A (in this case, typically a linear file scan may be cheaper; but this depends on the selectivity of A)

Complex Selections

- General pattern:
 - Conjunction – $\sigma_{\theta_1 \wedge \dots \wedge \theta_n}(R)$
 - Disjunction – $\sigma_{\theta_1 \vee \dots \vee \theta_n}(R)$
 - Negation – $\sigma_{\neg \theta}(R)$
- 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 tuples in R that satisfy Θ_i , then Θ_i 's selectivity is estimated as s_i/N_R .

Join Operations

- There are several different algorithms that can be used to implement joins (natural-, equi-, condition-join)
 - Nested-Loop Join
 - Block Nested-Loop Join
 - Index Nested-Loop Join
 - Sort-Merge Join
 - Hash-Join
- Choice of a particular algorithm is based on cost estimate
- For this, **join size** estimates are required and in particular cost estimates for outer-level operations in a relational algebra expression.

- **Example:**

Assume the query $CUSTOMERS \bowtie ORDERS$ (with join attribute only being CName)

- $N_{CUSTOMERS} = 5,000$ tuples
- $F_{CUSTOMERS} = 20$, i.e., $B_{CUSTOMERS} = 5,000/20 = 250$ blocks
- $N_{ORDERS} = 10,000$ tuples
- $F_{ORDERS} = 25$, i.e., $B_{ORDERS} = 400$ blocks
- $V(CName, ORDERS) = 2,500$, meaning that in this relation, on average, each customer has **four** orders
- Also assume that CName in ORDERS is a foreign key on CUSTOMERS

Estimating the Size of Joins

- The Cartesian product $R \times S$ results in $N_R * N_S$ tuples; each tuple requires $S_R + S_S$ bytes.
- If $\text{schema}(R) \cap \text{schema}(S) = \text{primary key for } R$, then a tuple of S will match with at most one tuple from R .
Therefore, the number of tuples in $R \bowtie S$ is not greater than N_S

If $\text{schema}(R) \cap \text{schema}(S) = \text{foreign key in } S \text{ referencing } R$, then the number of tuples in $R \bowtie S$ is exactly N_S .

Other cases are symmetric.

- In the example query $\text{CUSTOMERS} \bowtie \text{ORDERS}$, $CName$ in ORDERS is a foreign key of CUSTOMERS ; the result thus has exactly $N_{\text{ORDERS}} = 10,000$ tuples
- If $\text{schema}(R) \cap \text{schema}(S) = \{A\}$ is not a key for R or S ; assume that every tuple in R produces tuples in $R \bowtie S$. Then 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 is $\frac{N_R * N_S}{V(A, R)}$

and the lower of the two estimates is probably the more accurate one.

- Size estimates for CUSTOMERS \bowtie ORDERS without using information about foreign keys:
 - $V(\text{CName}, \text{CUSTOMERS}) = 5,000$, and
 $V(\text{CName}, \text{ORDERS}) = 2,500$
 - The two estimates are $5,000 * 10,000 / 2,500 = 20,000$ and $5,000 * 10,000 / 5,000 = 10,000$.
- We choose the lower estimate, which, in this case, is the same as our earlier computation using foreign key information.

Nested-Loop Join

- Evaluate the condition join $R \bowtie_C S$
 - **for each** tuple t_R **in** R **do begin**
 - for each** tuple t_S **in** S **do begin**
 - check whether pair (t_R, t_S) satisfies join condition
 - if they do, add $t_R \circ t_S$ to the result
 - end**
 - end**
- R is called the **outer** and S the **inner** relation of the join.
- Requires no indexes and can be used with any kind of join condition.
- **Worst case:** db buffer can only hold one block of each relation
 $\implies B_R + N_R * B_S$ disk accesses
- **Best case:** both relations fit into db buffer
 $\implies B_R + B_S$ disk accesses.

An Improvement: Block Nested-Loop Join

- Evaluate the condition join $R \bowtie_C S$
- **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**
 for each tuple t_S **in** B_S **do**
 check whether pair (t_R, t_S)
 satisfies join condition
 if they do, add $t_R \circ t_S$ to the result
 end end end end
- Also requires no indexes and can be used with any kind of join condition.
- **Worst case:** db buffer can only hold one block of each relation
 $\implies B_R + B_R * B_S$ disk accesses.
- **Best case:** both relations fit into db buffer
 $\implies B_R + B_S$ disk accesses.
- If **smaller relation** completely fits into db buffer, use that as inner relation. Reduces the cost estimate to $B_R + B_S$ disk accesses.

Block Nested-Loop Join (cont.)

- Some improvements of block nested-loop algorithm
 - If equi-join attribute is the key on inner relation, stop inner loop with first match
 - Use $M - 2$ disk blocks as blocking unit for outer relation, where M = db buffer size in blocks; use remaining two blocks to buffer inner relation and output.
Reduces number of scans of inner relation greatly.
 - Scan inner loop forward and backward alternately, to make use of blocks remaining in buffer (with LRU replacement strategy)
 - Use index on inner relation, if available . . .

Index Nested-Loop Join

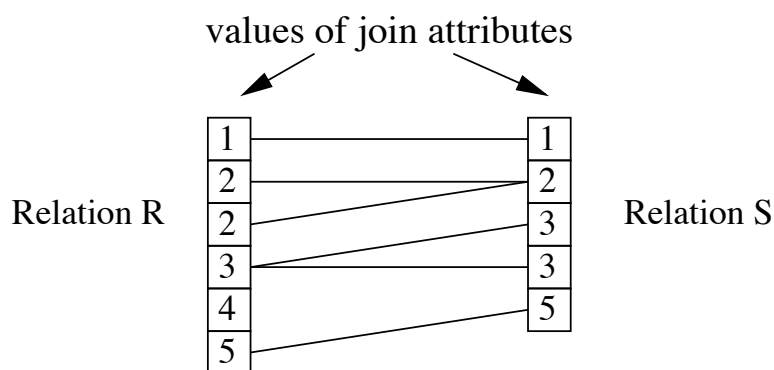
- If an index is available on the inner loop's join attribute and join is an equi-join or natural join, more efficient index lookups can replace file scans.
- It is even possible (reasonable) to construct index just to compute a join.
- For each tuple t_R in the outer relation R, use the index to lookup tuples in S that satisfy join condition with t_R
- **Worst case:** db buffer has space for only one page of R and one page of the index associated with S:
 - B_R disk accesses to read R, and for each tuple in R, perform index lookup on S.
 - Cost of the join: $B_R + N_R * c$, where c is the cost of a single selection on S using the join condition.
- If indexes are available on both R and S, use the one with the fewer tuples as the outer relation.

- **Example:**

- Compute $CUSTOMERS \bowtie ORDERS$, with $CUSTOMERS$ as the outer relation.
- Let $ORDERS$ have a primary B^+ -tree index on the join-attribute $CName$, which contains 20 entries per index node
- Since $ORDERS$ has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data records (based on tuple identifier).
- Since $N_{CUSTOMERS}$ is 5,000, the total cost is $250 + 5000 * 5 = 25,250$ disk accesses.
- This cost is lower than the 100,250 accesses needed for a block nested-loop join.

Sort-Merge Join

- Basic idea: first sort both relations on join attribute (if not already sorted this way)
- Join steps are similar to the merge stage in the external sort-merge algorithm (discussed later)
- Every pair with same value on join attribute must be matched.



- If no repeated join attribute values, each tuple needs to be read only once. As a result, each block is read only once. Thus, the number of block accesses is $B_R + B_S$ (plus the cost of sorting, if relations are unsorted).
- Worst case: all join attribute values are the same. Then the number of block accesses is $B_R + B_R * B_S$.
- If one relation is sorted and the other has a secondary B^+ -tree index on the join attribute, a *hybrid merge-join* is possible. The sorted relation is merged with the leaf node entries of the B^+ -tree. The result is sorted on the addresses (rids) of the unsorted relation's tuples, and then the addresses can be replaced by the actual tuples efficiently.

Hash-Join

- only applicable in case of equi-join or natural join
- a hash function is used to partition tuples of both relations into sets that have the same hash value on the join attribute

Partitioning Phase: $2 * (B_R + B_S)$ block accesses

Matching Phase: $B_R + B_S$ block accesses

(under the assumption that one partition of each relation fits into the database buffer)

Cost Estimates for other Operations

Sorting:

- If whole relation fits into db buffer \rightsquigarrow quick-sort
- Or, build index on the relation, and use index to read relation in sorted order.
- Relation that does not fit into db buffer \rightsquigarrow external sort-merge
 1. Phase: Create *runs* by sorting portions of the relation in db buffer
 2. Phase: Read runs from disk and merge runs in sort order

Duplicate Elimination:

- Sorting: remove all but one copy of tuples having identical value(s) on projection attribute(s)
- Hashing: partition relation using hash function on projection attribute(s); then read partitions into buffer and create in-memory hash index; tuple is only inserted into index if not already present

Set Operations:

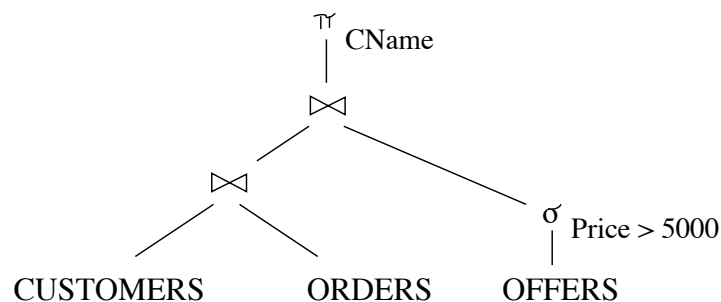
- Sorting or hashing
- Hashing: Partition both relations using the same hash function; use in-memory index for partitions R_i
 $R \cup S$: if tuple in R_i or in S_i , add tuple to result
 \cap : if tuple in R_i and in S_i , . . .
 $-$: if tuple in R_i and not in S_i , . . .

Grouping and aggregation:

- Compute groups via sorting or hashing.
- Hashing: while groups (partitions) are built, compute partial aggregate values (for group attribute A , $V(A,R)$ tuples to store values)

Evaluation of Expressions

- *Strategy 1: materialization.* Evaluate one operation at a time, starting at the lowest level. Use intermediate results materialized in temporary relations to evaluate next level operation(s).



- First compute and store $\sigma_{\text{Price} > 5000}(\text{OFFERS})$; then compute and store join of CUSTOMERS and ORDERS; finally, join the two materialized relations and project on to CName.
- *Strategy 2: pipelining.* evaluate several operations simultaneously, and pass the result (tuple- or block-wise) on to the next operation.

In the example above, once a tuple from OFFERS satisfying selection condition has been found, pass it on to the join. Similarly, don't store result of (final) join, but pass tuples directly to projection.

- Much cheaper than materialization, because temporary relations are not generated and stored on disk.

Evaluation of Expressions (cont.)

- Pipelining is not always possible, e.g., for all operations that include sorting (*blocking operation*).
- Pipelining can be executed in either *demand driven* or *producer driven* fashion.

Transformation of Relational Expressions

- Generating a query-evaluation plan for an expression of the relational algebra involves two steps:
 1. generate logically equivalent expressions
 2. annotate these evaluation plans by specific algorithms and access structures to get alternative query plans
- Use *equivalence rules* to transform a relational algebra expression into an equivalent one.
- Based on estimated cost, the most cost-effective annotated plan is selected for evaluation. The process is called **cost-based query optimization**.

Equivalence of Expressions

Result relations generated by two equivalent relational algebra expressions have the same set of attributes and contain the same set of tuples, although their attributes may be ordered differently.

Equivalence Rules (for expressions E, E_1, E_2 , conditions F_i)

Applying distribution and commutativity of relational algebra operations

1. $\sigma_{F_1}(\sigma_{F_2}(E)) \equiv \sigma_{F_1 \wedge F_2}(E)$
2. $\sigma_F(E_1 [\cup, \cap, -] E_2) \equiv \sigma_F(E_1) [\cup, \cap, -] \sigma_F(E_2)$
3. $\sigma_F(E_1 \times E_2) \equiv \sigma_{F_0}(\sigma_{F_1}(E_1) \times \sigma_{F_2}(E_2));$
 $F \equiv F_0 \wedge F_1 \wedge F_2, F_i$ contains only attributes of $E_i, i = 1, 2.$
4. $\sigma_{A=B}(E_1 \times E_2) \equiv E_1 \bowtie_{A=B} E_2$
5. $\pi_A(E_1 [\cup, \cap, -] E_2) \not\equiv \pi_A(E_1) [\cup, \cap, -] \pi_A(E_2)$
6. $\pi_A(E_1 \times E_2) \equiv \pi_{A_1}(E_1) \times \pi_{A_2}(E_2),$
with $A_i = A \cap \{\text{attributes in } E_i\}, i = 1, 2.$
7. $E_1 [\cup, \cap] E_2 \equiv E_2 [\cup, \cap] E_1$
 $(E_1 \cup E_2) \cup E_3 \equiv E_1 \cup (E_2 \cup E_3)$ (the analogous holds for \cap)
8. $E_1 \times E_2 \equiv \pi_{A_1, A_2}(E_2 \times E_1)$
 $(E_1 \times E_2) \times E_3 \equiv E_1 \times (E_2 \times E_3)$
 $(E_1 \times E_2) \times E_3 \equiv \pi((E_1 \times E_3) \times E_2)$
9. $E_1 \bowtie E_2 \equiv E_2 \bowtie E_1$ $(E_1 \bowtie E_2) \bowtie E_3 \equiv E_1 \bowtie (E_2 \bowtie E_3)$

The application of equivalence rules to a relational algebra expression is also sometimes called **algebraic optimization**.

Examples:

- Selection:

- Find the name of all customers who have ordered a product for more than \$5,000 from a supplier located in Davis.

$$\pi_{CName}(\sigma_{SAddress \text{ like } '%Davis\%'} \wedge Price > 5000 \\ (CUSTOMERS \bowtie (ORDERS \bowtie (OFFERS \bowtie SUPPLIERS))))$$

Perform selection as early as possible (but take existing indexes on relations into account)

$$\pi_{CName}(CUSTOMERS \bowtie (ORDERS \bowtie \\ (\sigma_{Price > 5000}(OFFERS) \bowtie (\sigma_{SAddress \text{ like } '%Davis\%'}(SUPPLIERS))))))$$

- Projection:

- $\pi_{CName,account}(CUSTOMERS \bowtie \sigma_{Prodname='CD-ROM'}(ORDERS))$

Reduce the size of argument relation in join

$$\pi_{CName,account}(CUSTOMERS \bowtie \pi_{CName}(\sigma_{Prodname='CD-ROM'}(ORDERS)))$$

Projection should not be shifted before selections, because minimizing the number of tuples in general leads to more efficient plans than reducing the size of tuples.

Join Ordering

- For relations R_1, R_2, R_3 ,

$$(R_1 \bowtie R_2) \bowtie R_3 \equiv R_1 \bowtie (R_2 \bowtie R_3)$$

- 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 a smaller temporary relation is computed and materialized

- Example: List the name of all customers who have ordered a product from a supplier located in Davis.

$$\pi_{CName}(\sigma_{SAddress \text{ like } '%Davis\%'}(\text{SUPPLIERS} \bowtie \text{ORDERS} \bowtie \text{CUSTOMERS}))$$

$\text{ORDERS} \bowtie \text{CUSTOMERS}$ is likely to be a large relation. Because it is likely that only a small fraction of suppliers are from Davis, we compute the join

$$\sigma_{SAddress \text{ like } '%Davis\%'}(\text{SUPPLIERS} \bowtie \text{ORDERS})$$

first.

Summary of Algebraic Optimization Rules

1. Perform selection as early as possible
2. Replace Cartesian Product by join whenever possible
3. Project out useless attributes early.
4. If there are several joins, perform *most restrictive* join first

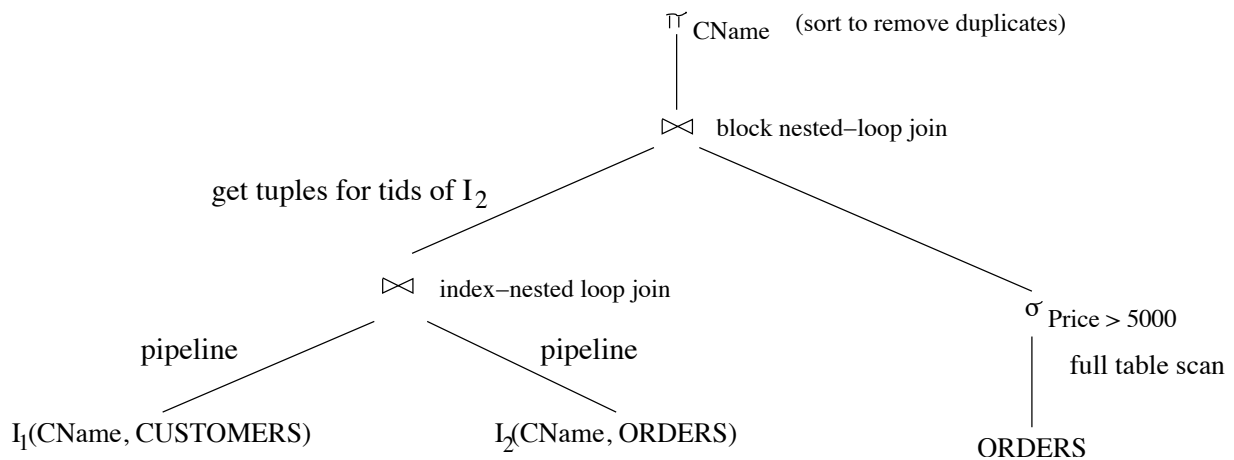
Evaluation Plan

An evaluation plan for a query exactly defines what algorithm is used for each operation, which access structures are used (tables, indexes, clusters), and how the execution of the operations is coordinated.

Example of Annotated Evaluation Plan

- Query: List the name of all customers who have ordered a product that costs more than \$5,000.

Assume that for both CUSTOMERS and ORDERS an index on CName exists: $I_1(\text{CName}, \text{CUSTOMERS})$, $I_2(\text{CName}, \text{ORDERS})$.



Choice of an Evaluation Plan

- Must consider interaction of evaluation techniques when choosing evaluation plan: choosing the algorithm with the least cost for each operation independently may not yield the best overall algorithm.
- Practical query optimizers incorporate elements of the following two optimization approaches:
 - Cost-based: enumerate all the plans and choose the best plan in a cost-based fashion.
 - Rule-based: Use rules (heuristics) to choose plan.
- Remarks on cost-based optimization:
 - Finding a join order for $R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$:
 - $n!$ different *left-deep* join orders
 - * For example, for $n = 9$, the number is 362880.
 - \rightsquigarrow use of dynamic programming techniques
- Heuristic (or rule-based) optimization transforms a given query tree by using a set of rules that typically (but not in all cases) improve execution performance:
 - Perform selection early (reduces number of tuples)
 - Perform projection early (reduces number of attributes)
 - Perform most restrictive selection and join operations before other similar operations.