

# Chapter - 5

Date.

## Query Processing and Optimization.

### 5.1 Query Cost Estimation

### 5.2 Query Operations

### 5.3 Evaluation of Expressions

### 5.4 Query Optimization

### 5.5 Query Decomposition

### 5.6 Performance Tuning.

## Query Processing.

- ↳ Query processing is a translation of high-level queries into low level expression.
- ↳ It refers to the range of activities that are involved in extracting data from the database.
- ↳ It includes translation of queries in high-level database languages into expressions that can be implemented at the physical level of the file system.
- ↳ It is a step wise process that can be used at the physical level of the

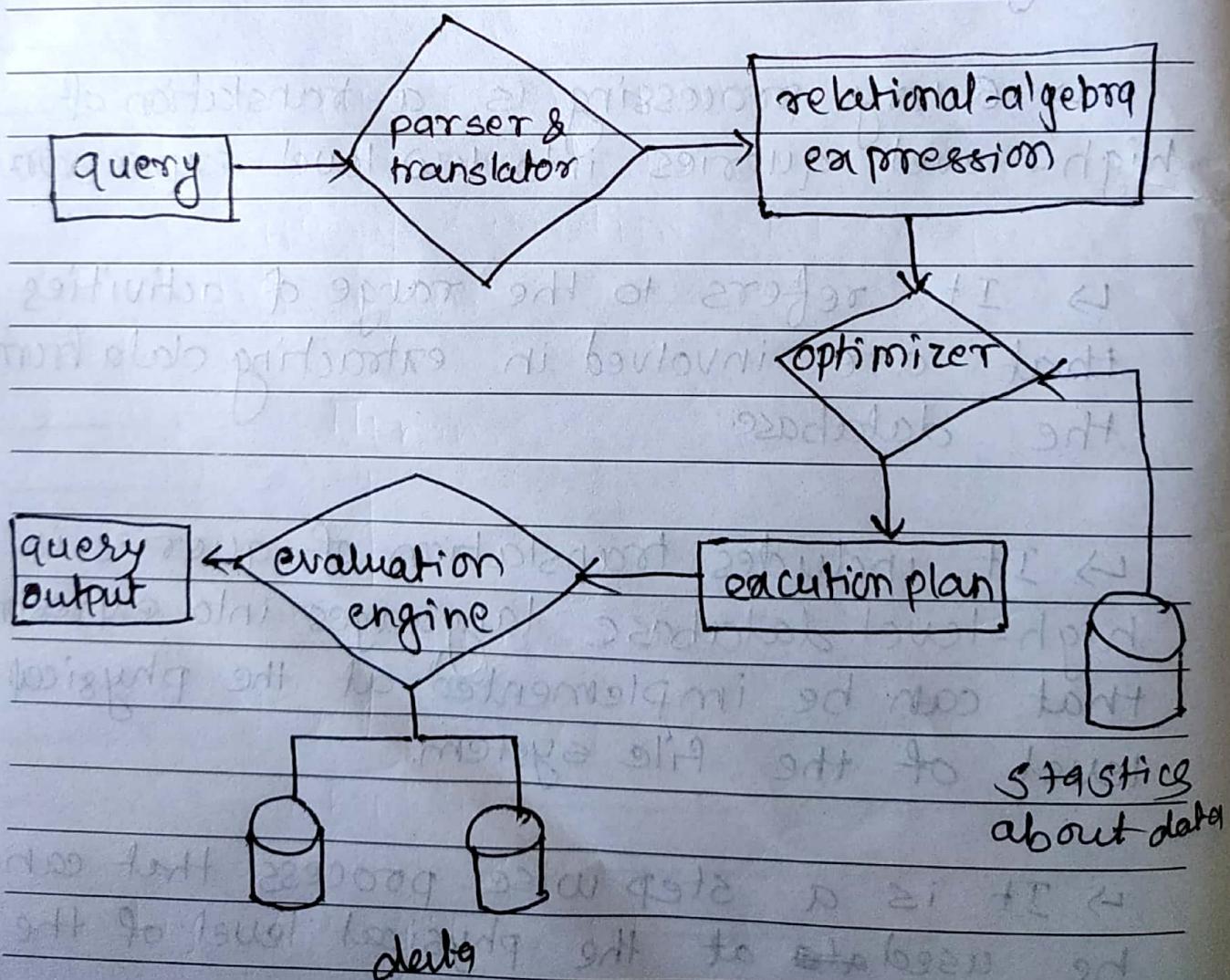
Date. \_\_\_\_\_

file system, query optimization and actual execution of query to get the result.

↳ It requires the basic concepts of relational algebra and file structure.

## Basic steps in Query Processing.

- i. Parsing and translation.
- ii. Optimization.
- iii. Evaluation.



## Query parsing and translation (query compiler)

- ↳ A query is translated into SQL and into a relational algebraic expression.
- ↳ parser checks the syntax and verifies the relations and the attributes which are used in the query.
- ↳ Different possible relational algebra expression for a single query.

**Example:**

SELECT Ename FROM Employee  
WHERE salary > 5000;

Translated into Relational Algebra Expression

$\sigma_{\text{salary} > 5000} (\pi_{\text{Ename}} (\text{Employee}))$

OR

$\pi_{\text{Ename}} (\sigma_{\text{salary} > 5000} (\text{Employee}))$

## Query optimization (query optimizer)

- ↳ 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.
- ↳ A relational-algebra operation annotated with instructions on how much to evaluate it is called an evaluation primitive.
- ↳ A sequence of primitive operations that can be used to evaluate a query is a query execution plan or query evaluation plan.

$\Pi_{\text{Ename}}$

$\sigma_{\text{Salary} > 500}$

$\text{Employee}$

Fig: Query Execution plan.

## ③ Query evaluation ( command processor )

- ↳ The query-execution engine takes a query-evaluation plan, executes that plan and returns the answers to the query.
- ↳ specify which access path to follow.
- ↳ Specify which algorithm to use to evaluate operators.
- ↳ Specify how operators interleave.

## Query cost Estimation.

- ↳ The cost of query evaluation can be measured in terms of a number of different resources, including disk accesses, CPU time to execute a query, and, in distributed or parallel database system, the cost of communication. The response time for a query-evaluation plan, assuming no other activity is going on on the computer, would account for all these costs, and could be used as a good measure of the cost of the plan.

- ↪ CPU cost is difficult to calculate.
- ↪ CPU speed is at faster rate than disk speed.
- ↪ Typically disk access is the predominant cost, and is also relatively easy to estimate.
- ↪ Disk access cost is measured by taking into account following:
  - i. Number of seeks [no. of attempts to find]
  - ii. Number of blocks read.
  - iii. Number of blocks write.

cost to write is greater than cost to read a block.

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

$t_T$  - time to transfer one block  
 $+s$  - time for one seek.

Cost of  $b$  block transfer plus  $s$  seek  
 $b * t_T + s * ts$

We ignore CPU costs for simplicity.  
 Real systems do take CPU cost

into account.

Date.

We do not include cost to writing output in our cost formulae.

↳ 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.

## Algorithms for RA operator

### Selection operator:

File scan

Read one file at a time sequentially

Index Scan

Read file randomly

**File scan:**

search algorithms that locate & retrieve records that fulfill a selection condition.

### A1(Linear search)

↳ Here each record is read from the beginning of the file till search record is reached

↳ It checks each record for filtering condition one after the other.

↳ An initial seek is required to access the first block of the file.

5  
Date.  
↳ In case blocks of the file are not stored contiguously, extra seeks may be required, but we ignore this effect for simplicity

cost of  $b_r$  block transfers plus 1 seek.

$$\text{Cost estimate} = t_s + b_r * t_T$$

$b_r$  denotes the number of blocks in the file.

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

Average Case:

cost of  $(b_r/2)$  block transfers plus 1 seek.

$$t_s + (b_r/2) * t_T = \text{ind}$$

for worst case  $b_r$  blocks are required.

**Index scan:**

↳ Search algorithms that use an index are referred to as index scans.

predicate to

↳ we use the selection, guide us in the choice of index to use in processing

the query.

## A2(Primary index, equality on key.)

↳ In this method search key value will be the primary key and index will be created on it.

↳ These types of indexes are usually stored in B+ tree structure with height 'h'.

↳ Worst case of the query is given as:

$$(h_i + 1) * (t_T + t_S)$$

$h_i$  = height of the index.

↳  $(h_i + 1)$  is given because to seek to traverse the B+ tree of height  $h$  and we need one seek to traverse the data file where the record is stored.

### A3(Primary index, equality on nonkey)

Date.

↳ we can retrieve multiple records by using a primary index when the selection condition specifies an equality comparison on a nonkey attribute.

↳ cost of query becomes:

$$h_i * (TT + TS) + b * TT$$

↳ seek and traversal time for B+ tree of height  $h$  to get index of the record.

The we have to seek  $b$  times the data files and traverse  $b$  times all the fetched <sup>records</sup> from the data files

### A4(Secondary index, equality on nonkey)

↳ Selections specifying an equality condition can use a secondary index.

↳ This strategy can retrieve a single record if the equality condition is on key; multiple records may be returned if the indexing field is not a key

Date.

### A3(Primary index, equality on nonkey)

↳ we can retrieve multiple records by using a primary index when the selection condition specifies an equality comparison on a nonkey attribute.

↳ cost of query becomes:

$$h_i * (TT + TS) + b * TT$$

↳ seek and traversal time for B+ tree of height h to get index of the record.

Then we have to seek  $b$  times the data files and traverse  $b$  times all the fetched <sup>records</sup> from the data files

### A4(Secondary index, equality on nonkey)

↳ Selections specifying an equality condition can use a secondary index.

↳ This strategy can retrieve a single record if the equality condition is on key; multiple records may be retrieved if the indexing field is not a key

cost when search key is a candidate key

$$(h_i + 1) * (tT + ts)$$

cost when search key is not a candidate key.

$$(h_i + n) * (tT + ts)$$

where  $n$  is numbers of records fetched.

$$tT + dt + (st + rt) * n$$

↳ 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

- i. a linear file scan
- ii. By using indices

## A5 (Primary index, comparison.)

↳ In this method comparison operators like  $<$ ,  $\leq$ , and  $\geq$  are used to retrieve more than one record from the file.

→ This comparison operator is used on the primary key, on which index is created

cost:  ~~$h_i * (t_T + t_S)$~~

$$h_i * (t_T + t_S) + b * t_T$$

### A<sub>6</sub> (Secondary index, comparison)

(i) cost:  $(h_i + h) * (t_T + t_S)$

### Implementation of complex selections

#### Conjunction:

A conjunctive selection is a selection of the form

$$\sigma_{Q_1 \wedge Q_2 \wedge \dots \wedge Q_n}(r)$$

### A<sub>7</sub> (Conjunctive selection using one index)

→ We first determine whether an access path is available for an attribute in one of the simple conditions

↳ if one is, one of the selection algorithms A<sub>2</sub> through A<sub>6</sub> can retrieve records satisfying that condition.

$$(2f + r) * (n + id)$$

↳ The cost of algorithm A<sub>7</sub> is given by the cost of chosen algorithm.

(from previous problem you have seen)

### A<sub>8</sub> (conjunctive selection using composite index)

↳ An appropriate composite index (i.e. index on multiple attributes) may be available for some conjunctive selections

↳ The type of index determines which of algorithms A<sub>2</sub>, A<sub>3</sub> or A<sub>4</sub> will be used.

### A<sub>9</sub> (conjunctive selection by intersection of identifiers)

↳ Conjunctive selection operation involves the use of record pointers & record identifiers

↳ This algorithm requires indices with

record pointers on the fields involved in individual conditions

↳ The algorithm scans index for pointers to tuples that satisfy individual condition.

↳ The intersection of all retrieved pointers is the set of pointers to tuples that satisfy the conjunctive condition.

↳ The algorithm then uses the pointers to retrieve actual records

### Disjunction:

A disjunctive selection (e.g.) a selection of the form

$$\sigma_{Q_1 \vee Q_2 \vee \dots \vee Q_n}^{\{r\}}$$

Also (Disjunctive selection by union of identifiers)

If access paths are available on all the conditions of a disjunctive selection, each index is scanned for pointers to tuples that satisfy the individual condition.

↳ The union of all the retrieved pointers yields the set of pointers to all tuples that satisfy disjunctive condition.

↳ We then use the pointers to retrieve the actual records.

### Negation

↳ The result of a selection  $\sigma_{\neg Q}(r)$  is the set of tuples of  $r$  for which that condition  $\neg Q$  evaluates to false.

↳ In the absence of nulls, this set is simply the set of tuples of  $r$  that are not in  $\sigma_Q(r)$ .

### Sorting

↳ We can sort a relation by building an index on the sort key, and then using that index to read the relation in sorted order.

↳ Such process orders that relation only logically, through an index, rather

than physically

- ↳ For relations that fit in memory that is records are completely in main 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, which cumulatively sorts multiple runs of the data based on amount that fits in memory at one time.

### Join operation:

Different algorithms to implement joins.

- i. Nested-loop join.
- ii. Block nested-loop join.
- iii. Indealed nested-loop join
- iv. Merge-join
- v. Hash-join

choice is based on cost estimate.

## Nested loop join.

↳ Simple algorithm to compute the theta join  $r \bowtie s$ , of two relations  $r$  and  $s$ . This algorithm is called nested-loop join algorithm, since it basically consists of a pair of nested for loops.

```

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 ◊
        if they do, add  $t_r, t_s$  to the result.
    end
end.

```

↳  $r$  is called the outer relation &  $s$  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 two relations.

↳ In the worst case, if there is enough

Date.

memory only to hold one block of each relation, the estimated cost is.

$n_r * b_s + b_r$ , block transfers plus  
 $n_r + b_r$  seeks

Where  $n_r$  denotes the number of tuples in  $r$  and  $b_s$  denotes number of blocks containing tuples  $r$  and  $s$ , respectively.

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

cost:  $b_r + b_s$  block transfers and 2 seeks.

Example:

Number of records of

$$n_{\text{Student}} = 5,000 \quad n_{\text{takes}} = 10,000$$

Number of blocks of:

$$b_{\text{student}} = 100 \quad b_{\text{takes}} = 400$$

worst case: (with student as outer relation)

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

$$5000 + 100 = 5100 \text{ seeks.}$$

with takes as outer relation

$$1000 * 100 + 400 = 1,000,400 \text{ block transfers}$$

$$10,400 \text{ seeks.}$$

If student fits entirely in memory,  
the cost estimate will be 500 block transfers.

### 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  $(tr, ts)$  satisfy the join condition

if they do, add  $tr, ts$  to the result.

end

end

end

end

worst case estimate

$b_r * b_{S+Br}$  block transfers plus 2 seeks.

Best case  $b_r + b_S$  block transfers plus 2 seeks.

### Merge Join:

→ sort both relations on their join attributes (if not already sorted on the join attributes,

→ The Merge Join operator is one of four operators that join data from two input streams into a single combined output stream.

- ↳ It has two inputs, called the left and right input.
- ↳ Merge join can't be used without adding extra sort operators. These extra sort increase the total plan cost.
- ↳ The Merge join operator supports all ten logical join operations: inner join; left, right and full outer join; left and right semi and anti semi join; as well a concatenation & union.
- ↳ Algorithm requires at least one equality based join predicate.
- ↳ can be used only for equi-joins and natural joins.

cost of merge join is:

$$br + bs \text{ block transfers} + [br/b_b] + [bs/b_b]$$

+ The cost of sorting if relations are unsorted.

## hybrid merge join

- ↳ If one relation is sorted, and the other has a secondary B<sup>+</sup>-tree index on the join attribute.
  - ↳ Merge the sorted relation with the leaf entries of the B<sup>+</sup>-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 address by the actual tuples.
  - ↳ Sequential scan more efficient than random lookup.
- Hash join:**
- ↳ A hash function  $h$  is used to partition tuples of both relations. The basic idea is to partition the tuples of each of the relations into sets that have same hash value on the join attributes.

we assume that :

- ↳  $h$  is a hash function mapping joinAttrs values to  $\{0, 1, \dots, n\}$ , where joinAttrs denotes
- ↳  $r_0, r_1, \dots, r_n$  denote partitions of  $\sigma$  tuples; each initially empty. Each tuple  $t \in \sigma$  is put in partition  $r_i$  where  $i = h(t[\text{joinAttrs}])$ .
- ↳  $s_0, s_1, \dots, s_n$  denote partitions  $S_i$ , where  $i = h(t[\text{joinAttrs}])$ .
- ↳ suppose that an  $r$  tuple and an  $s$  tuple satisfy the join condition; then, they 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  $s$  tuples in  $s_i$ .
- ↳  $\sigma$  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.

↳ If the value of  $n_h$  is greater than or equal to the number of page frames of memory, the partitioning cannot be done in one pass. Since there will not be enough buffer pages. Instead partitioning has to be done in repeated passes.

↳ If recursive partitioning is not required cost of hash join is

$$3(br+bs) + 4 * n_h \text{ block transfers} + \\ 2(br/bb) + (bs/bb) \text{ seeks.}$$

↳ If recursive partitioning is required.

Total cost estimate

$$2(br+bs) [\log_{bb} M - 1] (b/M) + br+bs \\ \text{block transfers} + 2([br/bb] + [bs/bb]) \\ \log_{bb} M - 1 (bs/M) \text{ seeks.}$$

↳ If the entire build input can be kept in main memory no partitioning is required. Cost estimate goes down to  $br+bs$

**Example:**

$b_{instructor}$   $\bowtie$   $b_{teaches}$

Assume that memory size is 20 blocks

$$b_{instructor} = 100 \text{ and } b_{teaches} = 400$$

$b_{instructor}$  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  $b_{teaches}$  into five partitions  
each of size 80. This is also done in  
one pass.

Therefore total cost, ignoring cost of  
(writing partially filled blocks)

$$3(100 + 400) = 1500 \text{ block transfers} +$$

$$2([100/3] + [400/3]) = 336 \text{ seeks}$$

## Evaluation of Expressions:

Date.

→ In reality expressions cannot exist as a single operation. instead a typical query combines them.

For example: select customers-name from account natural join customer where balance < 30 00

This example consists of a selection, a natural join and a projection.

so Alternatives for evaluating an entire expression tree:

- i. Materialization
- ii. Pipelining.

### Materialization:

→ In materialization an expression is evaluated one operation at a time in an appropriate order.

→ The result of each evaluation is

materialized in a temporary relation for subsequent use.

↳ Evaluation starts with one operation at a time, at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.

Consider the expressions

i.  $\Pi_{\text{customer-name}} (\sigma_{\text{balance} < 2500} (\text{account}) \bowtie \text{customer})$

$\Pi_{\text{customer-name}}$

$\sigma_{\text{balance} < 2500}$

account

customer

Fig: Pictorial representation of an expression.

ii.  $\sigma_{building = "watson"} (department) \bowtie instructor$

$\Pi name$

$\sigma_{building = "watson"} \bowtie instructor$

department

↳ Materialized evaluation is always applicable.

In the example 2, we start from the lowest-level operations in the expressions (at the bottom of the tree). Here only one such operation exists; the selection operation on account. The inputs to the lowest-level operations are relations in the database. We execute these operations by using suitable algorithm, we store the results in temporary relations.

~~the~~

we can use these temporary relations to execute the operations at the next level up in the tree, where the inputs are now either temporary relations or relations stored in the database.

→ In our example, the inputs to the join are the customer relation and temporary relation created by the selection on account. The join can now be evaluated creating another temporary relation.

→ By repeating the process, we will eventually evaluate the operation at the root of the tree giving the final result of the expression. Here we get final result by executing the projection operation at the root of the tree, using as input the temporary relation created by the join.

→ 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 and reduces execution time.

### Pipelining:

↳ we can improve query-evaluation efficiency by reducing the number of temporary files that are produced.

↳ we achieve this reduction by combining several relational operations into a pipeline of operations, in which the result of one operation are passed along to the next operation in the

Pipeline Evaluation as just described is called pipelined evaluation.

- ↳ For example: consider the expression  $(\Pi_{a_1, a_2} (r \bowtie s))$

If materialization were applied, evaluation would involve creating temporary relation to hold the result of the join, and then back in the result to perform the projection. These operations can be combined: when the join operation generates a tuple of its result, it passes that tuple immediately to the project operation for processing. By combining the join and the projection, we avoid creating the intermediate result, and instead create the final result directly.

- ↳ Pipelining is 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 either of two ways:

i. Demand driven.

ii. Producer driven.

### Demand driven:

In demand driven or lazy evaluation,

↳ the system makes repeated requests for tuples from the operation at the top of the pipeline.

↳ Each time that an operation receives a request for tuples, it computes the next tuple (or tuples) to be returned, and then returns that tuple.

↳ In between calls, operation has to maintain "state" so it knows what to return next.

## Producers - driven or eager pipeling.

In producer driven or eager pipeling

↳ Operations do not wait for requests to produce tuples, but instead generate the tuples eagerly.

↳ Each operation at the bottom of a pipeline continually generates output tuples, and puts them in its output buffer until the buffer is full.

↳ An operation at any other level of a pipeline generates output tuples when it gets input tuples from lower down in the pipeline until its output buffer is full.

↳ Once the operation uses a tuple from a pipelined input, it removes the tuple from its input buffer.

↳ The operation repeats this process until all the output tuples have been generated.

Alternative name: pull and push models of pipelining.

### Transformation of Relational Expressions

- ↳ A query can be expressed in several different ways, with different costs of evaluation
- ↳ Rather than taking the relational expression as given, we consider alternative, equivalent expressions.
- ↳ Two-relational-algebra expressions are said to be equivalent if, on every legal database instance, the two expressions generate the same set of tuples.
- ↳ In SQL, inputs and outputs are multisets of tuples, two expressions in the multiset version of the relational algebra are said to be equivalent if on every legal database instance the two expressions generate the same multiset of tuples.
- ↳ An equivalence rule says that expressions of two forms are equivalent; can replace

expression of first form by second, or  
vice versa.

↳ Generation of query-evaluation plans for an expression involves several steps:

- i. Generating logically equivalent expressions
  - ii. Use equivalence rules to transform an expression into an equivalent one
  - iii. Annotating resultant expressions to get alternative query plans.
  - iv. choosing the cheapest plan based on estimated cost.
- ↳ The overall process is called cost based optimization.

↳ Equivalence Rules:

An equivalence rule says that expressions of two forms are equivalent.

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.

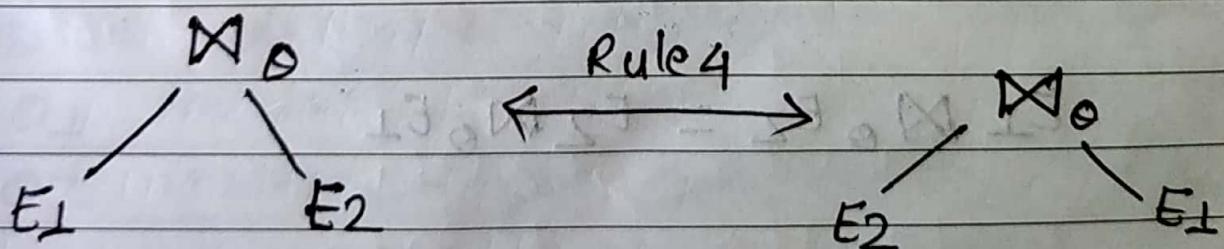
$$\pi_{t_1}(\pi_{t_2}(\dots(\pi_{t_n}(E))\dots)) = \pi_{t_1}(E)$$

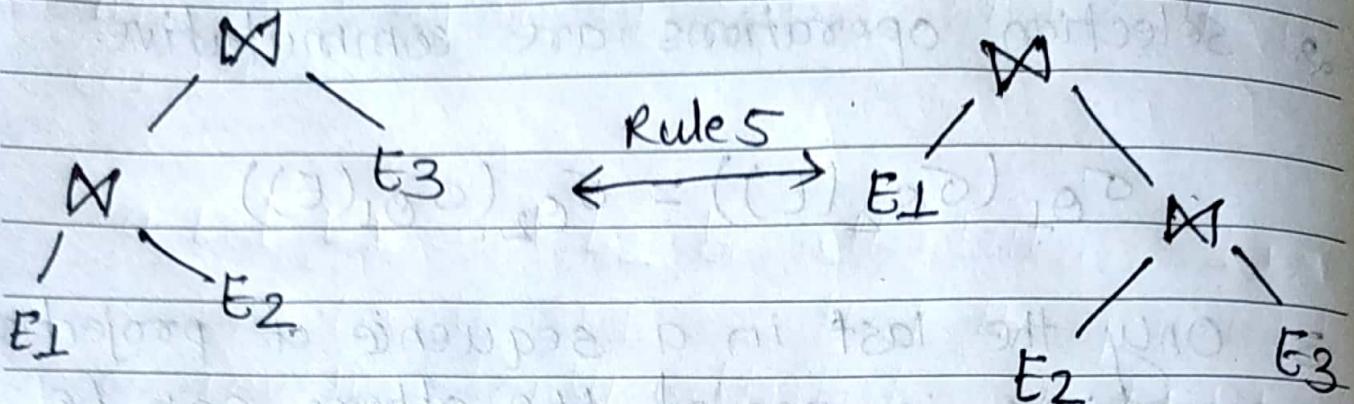
4. Selections can be combined with cartesian products and theta joins.

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

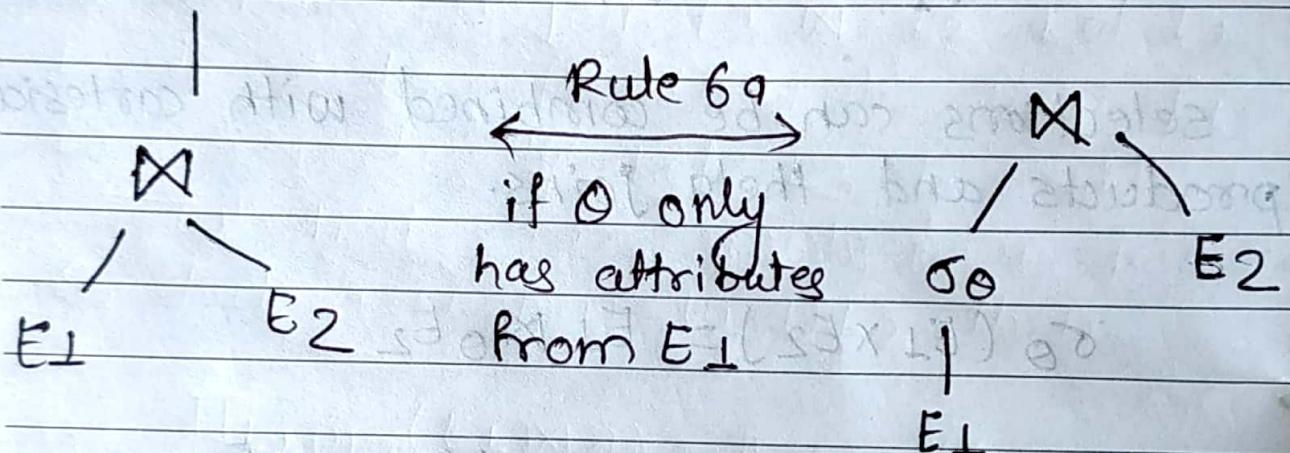
$$\sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) = E_1 \bowtie_{\theta_1 \bowtie \theta_2} E_2$$

Pictorial Depiction of Equivalence Rules:





$$(\exists) \sigma_\theta = ((\exists)_{E_1} \wedge (\exists)_{E_2} \wedge \dots \wedge (\exists)_{E_n}) \perp \perp$$



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

$$E_1 \Delta_\theta E_2 = E_2 \Delta_\theta E_1$$

6 a. Natural joins are 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:

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 = E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3)$$

corrected

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

7. The selection operation distributes over the theta-join operation under the following two conditions:

a. It distributes when all the attributes in selection condition  $\theta_0$  involve only the attributes of one of the expressions (say  $E_1$ ) being joined.

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

b. It distributes when selection condition  $\theta_1$  involves only the attributes of  $E_1$  &  $\theta_2$  involves only attributes of  $E_2$

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

8. The projection operation distributes over

the theta-join operation under the following conditions.

- a. Let  $L_1$  and  $L_2$  be attributes of  $E_1$  &  $E_2$  respectively. Suppose that the join condition  $\Theta$  involves only attributes in  $L_1 \cup L_2$ . Then

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

- 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  $\Theta$ , but are not in  $L_1 \cup L_2$ . Then

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

- c. The set operations union and intersection are commutative.

$$E_1 \cup E_2 = E_2 \cup E_1$$

$$E_1 \cap E_2 = E_2 \cap E_1$$

set difference is not commutative.

10. Set Union and intersection are associative

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

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

11. The selection operation distributes over the union, intersection, and set-difference operations.

$$\sigma_p(E_1 - E_2) = \sigma_p(E_1) - \sigma_p(E_2)$$

Similarly, the preceding equivalence, with  $-$  replaced with either  $\cup$  or  $\cap$  also holds. Further:

$$\sigma_p(E_1 - E_2) = \sigma_p(E_1) - E_2$$

The preceding equivalence with  $-$  replaced by  $\cap$  also holds, but does not hold if  $-$  is replaced by  $\cup$ .

12. The projection operation distributes over the union operation.

$$\pi_L(E_1 \cup E_2) = (\pi_L(E_1) \sqcup \pi_L(E_2))$$

## Transformation Example:

Query: Find the names of all customers who have an account at some branch located in kathmandu.

$\Pi_{\text{customer-name}} (\sigma_{\text{branch-city} = \text{"kathmandu"} } \text{branch} \bowtie \text{account} \bowtie \text{depositor})$

Transformation using rule 7a.

$\Pi_{\text{customer-name}} ((\sigma_{\text{branch-city} = \text{"kathmandu"} } (\text{branch})) \bowtie \text{account} \bowtie \text{depositor})$

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

Query: Find the names of all customers with an account at kathmandu branch whose account balance is over 10 000

$\Pi_{\text{customers-name}} (\sigma_{\text{branch-city} = \text{"kathmandu"} } \wedge \text{balance} > 10000 (\text{branch} \bowtie \text{account} \bowtie \text{depositor}))$

## Date.

### Transformation using join associatively (Rule 6a)

$\Pi \text{customer-name} ((\sigma_{\text{branch-city} = \text{"kathmandu"}} \wedge \sigma_{\text{balance} > 1000}) (\text{branch} \bowtie \text{account})) \bowtie \text{depositor}$ .

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

$\sigma_{\text{branch-city} = \text{"kathmandu"}} (\text{branch}) \bowtie \sigma_{\text{balance} > 1000} (\text{account})$

$\Pi \text{customer-name}$

$\sqcap$

$\sigma_{\text{branch-city} = \text{kathmandu}}$   
 $\wedge \sigma_{\text{balance} > 1000}$

$\bowtie$

$\text{branch}$

$\bowtie$

$\text{account}$

$\text{depositor}$

$\Pi \text{customer-name}$

$\bowtie$

$\text{depositor}$

$\bowtie$

$\sigma_{\text{branch-city} = \text{"kathmandu"}}$

$"\text{kathmandu}"$

$\text{account}$

$\text{branch}$

Initial Expression Tree

Tree after multiple Transformations

## Query Decomposition

The Query decomposition is the first phase of query processing whose aims are to transfer a high-level query into a relational algebra query & to check whether that query is syntactically and semantically correct.

### Stages of Query Decomposition.

- i. Normalization
- ii. Analysis
- iii. Elimination of Redundancy
- iv. Rewriting

#### Normalization.

Converts the query into a normalized form that can be more easily manipulated.

There are two different normal forms, conjunctive normal form and disjunctive normal form.

## Conjunctive normal form.

A sequence of conjuncts that are connected with the and operator. Each conjunct contains one or more terms connected by the or operator.

for example:

(position = 'Manager'  $\vee$  salary > 20000)  $\wedge$  branch  
NO = 'B003'.

## Disjunctive normal form.

A sequence of disjuncts that are connected by the or operator. Each disjunct contains one or more terms connected by the and operator.

for example:

(position = 'Manager'  $\wedge$  branchNO = 'B003')  $\vee$   
(salary  $\geq$  20000  $\wedge$  branchNO = 'B003')

## Analysis:

The query is lexically and syntactically analyzed using the techniques of programming language compilers.

- ↳ verifies that the relations and attributes specified in the query are defined in the system catalog.
- ↳ verifies that any operations applied to database objects are appropriate for the object type.
- ↳ On completion of the analysis, the high-level query has been transformed into some internal representation (query tree) that is more suitable for processing.

Root

Intermediate operation



Leaves

## Simplification.

To detect redundant qualifications, eliminate common subexpressions, and

transform the query to a semantically equivalent but more easily and efficiently computed form.

- ↳ Access restrictions, view definitions, & integrity constraints are considered at this stage.
- ↳ Expression replacement already used in views.

### Query restructuring / rewriting.

- ↳ The final stage of query decomposition.
- ↳ The query is restructured to provide a more efficient implementation.
- ↳ Operator Tree is used.  
 Leaf nodes are operand relations  
 Non leaf are intermediate tables produced as a result of some relational operators.

Example:

Select ename FROM EMP, ASG, PROJ  
 Where

$\text{EMP.eN} = \text{ASG.eNo}$

AND  $\text{ASG.pNo} = \text{PROJ.pNo}$

AND  $eName \neq 'Saleem'$

AND  $pName = 'CAD/CAM'$

AND  $(dur = 36 \text{ or } dur = 24)$

$\Pi eName$

} projection

$\sigma_{dur = 12 \vee dur = 24}$



$\sigma_{pName = 'CAD/CAM'}$

$\sigma_{eName = 'Saleem'}$

} select

$\bowtie pNo$

$\bowtie eNo$

$\bowtie ASG$

$\bowtie eNo$

$\bowtie EMP$

} join

Second

## Query Optimization.

Date.

- ↳ Query optimization is the process of selecting the most efficient query-evaluation plan from among many strategies usually possible for processing a given query, especially query is complex.
  - ↳ we do not expect users to write their queries so that they can be processed efficiently. Rather, we expect the system to construct a query evaluation plan that minimizes the cost of query evaluation. This is where query optimization comes into play.
  - ↳ An important aspect of query processing is query optimization.
  - ↳ The aim of query optimization is to choose the one that minimizes resource usage.
- query → equivalent query 1  
                 ↑  
                 ↓  
                 equivalent query 2      ←→ faster query  
                 ↓  
                 ↓  
                 equivalent query n

Aim: transfer query into faster, equivalent query.

- ↳ Every method of query optimization depend on database statistics.
- ↳ The statistics cover information about relations, attribute, and indexes.
- ↳ Keeping the statistics current can be problematic.
- ↳ If the DBMS updates the statistics every time a tuple is inserted, updated, or deleted, this would have a significant impact on performance during peak period.
- ↳ An alternative approach is to update the statistics on a periodic basis, for example nightly or whenever the system is idle.
- ↳ After a query has been scanned, parsed and validated, the DBMS must choose from a set of possible execution strategies for the query, and this is known as

## query optimization.

Date.

### Types:

- i. Heuristic (Logical) query optimization.
- ii. cost-based (Physical) query optimization.

### Cost Based optimization.

- ↳ A cost based optimization explores the space of all query-evaluation plans that are equivalent to the given query, and chooses the one with the least estimated cost.
- ↳ cost of physical plans include processor time and communication time. The most important factor to consider is disk I/Os because it is the most time consuming action.
- ↳ some other costs associated are: operations (joins, unions, intersections).

### Steps in cost-based (query) optimization.

- i. generate logically equivalent expressions using equivalence rules.

ii. Annotate resultant expressions to get alternative query plans.

iii. 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 formula for algorithms, computed using statistics.

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

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

- ↳ No need to generate all the join orders. Using dynamic programming, the least-cost join order for any subset of  $\{r_1, r_2, \dots, r_n\}$  is computed only once and stored for future use.
- ↳ Cost based optimization is expensive, but worthwhile for queries on large datasets (typically queries have small  $n$ , generally  $\leq 10$ )

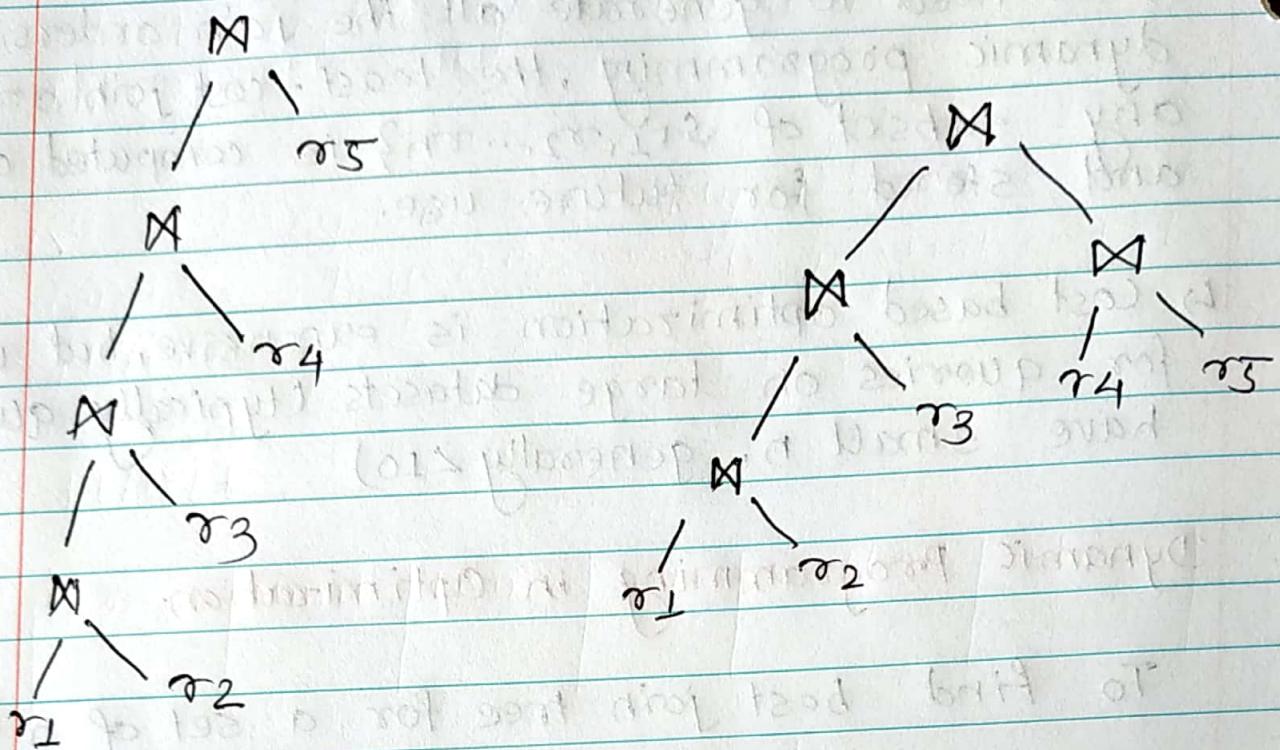
## Dynamic Programming in Optimization

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  $2^{n-1}$  alternatives.
- ↳ When plan for any subset is computed, store it and reuse it when it is required again, instead of recomputing it.

## 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.



Left-deep join tree. not a non-left-deep join tree.

### Cost of optimization.

With dynamic programming time complexity with bushy tree 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 other relations

as left hand side input

- ↳ Using (recursively computed and stored) least-cost join order for each alternative on left-hand side, choose the cheapest of the  $n$  alternatives.

If only left-deep trees are considered, time complexity of finding best join order is  $O(n^{2n})$

Space complexity remains at  $O(2^n)$

### Heuristic Optimization.

- ↳ Heuristic optimization transforms the query-tree by using a set of rules (Heuristics) that typically (but not in all cases) improve execution performance.
  - i. Perform selection early (reduces the number of tuples)
  - ii. Perform projection early (reduces the number of attributes)
  - iii. Perform most restrictive selection & join operations (i.e. with smallest result size) before the other similar operations.

- ↳ Generate initial query tree from SQL statement & a single query tree is involved.

- ↳ Transform query tree into more efficient query tree, via series of tree modifications, each of

which hopefully reduces the execution time.

↳ Some systems use only heuristics; others combine heuristics with partial cost based optimization.

### PLAYER

Fname	Minit	Lname	<u>SSN</u>	Draft date	Address	Salary

coach-ssn	Team Name
-----------	-----------

### UNIT

Uname	Unumber	coach-ssn

### PLAY

PlayName	<u>Play Num</u>	PlayType	Unum

### WORKS-ON

<u>PSSN</u>	PlayNo	Hours

Query:

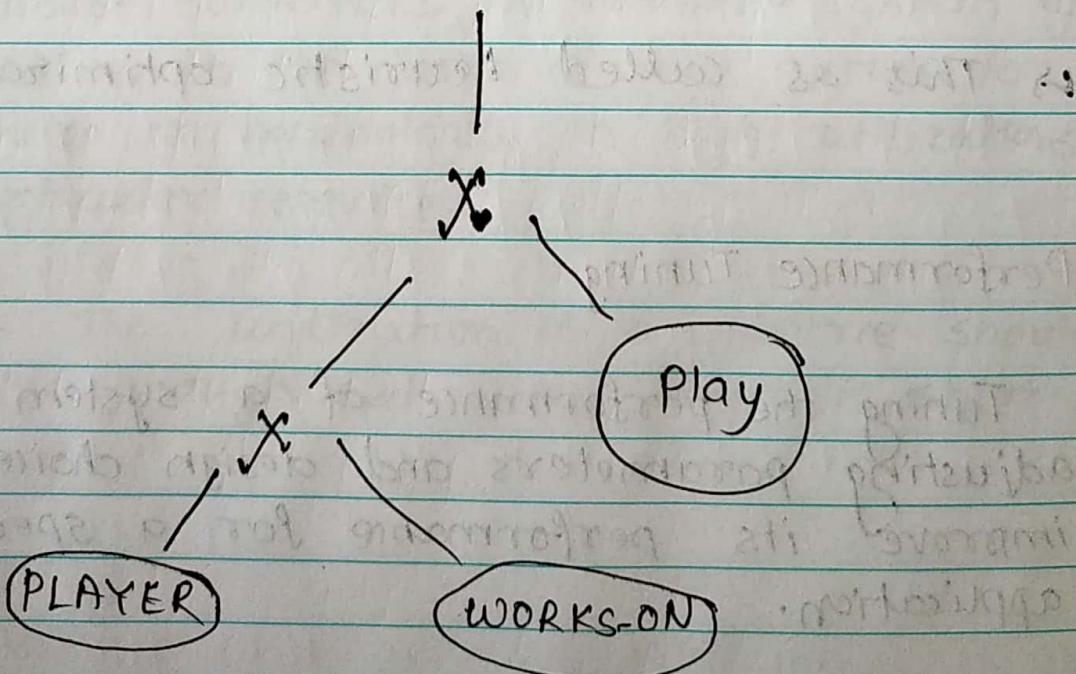
Find the last names of players who were drafted after 2000 who have worked on a play called 'halfback-pass'.

SQL:

```
SELECT Lname  
FROM PLAYER, WORKS-ON, PLAY  
WHERE playName = 'halfback-pass' AND  
PlayNo = PlayNum AND PSSn = SSN AND  
DraftDate > '2000-12-31'
```

PROJECT (Lname)

```
SELECT (PlayName = 'halfback-pass' AND playNo = playNum AND  
PSSn = SSN AND DraftDate > '2000-12-31')
```



- ↳ It should work, but it is not optimal
- ↳ We have to find the cartesian product of the PLAYER, WORKS-ON, and PLAY files. Depending on the size and number of records in each table, cartesian product could be huge and would contain lots of unneeded data.
- ↳ First, move SELECT operation down the query tree.
- ↳ Second, perform the more restrictive SELECT operations first.
- ↳ Third replace CARTESIAN PRODUCT and SELECT combination with join operations.
- ↳ Finally, move PROJECT operations down the query tree.
- ↳ This is called heuristic optimization.

## Performance Tuning:

Tuning the performance of a system involves adjusting parameters and design choices to improve its performance for a specific application.

↳ Various aspects of a database-system design - ranging from high level aspects such as the schema and transaction design, to database parameters such as buffer sizes, down to hardware issues such as number of disks - affect the performance of an application. Each of these aspects can be adjusted so that performance is improved.

i. Location of Bottleneck

ii. Tunable parameters

iii. Tuning of Hardware

iv. Tuning of the schema

v. Tuning of indices

vi. Using Materialized views

vii. Tuning of Transactions

viii. Performance simulation

### Location of Bottleneck.

↳ As a result of the numerous queues in the database, bottlenecks in database system typically show up in the form of long queues for particular service or equivalently in high utilizations for a particular service.

↳ The utilization of a resource should be kept low enough that queue length is short.

## Tunable parameters

- ↳ Database administrators can tune database system at three levels.
- ↳ The lowest level is at the hardware level
  - ↳ Options for tuning at this level includes adding disks or using a RAID system if disk I/O is bottleneck, adding more memory if the disk buffer size is a bottleneck, or moving to a faster processor if CPU use is bottleneck.
- ↳ The second level consists of the database system parameters, such as buffersize and check pointing intervals.
- ↳ The exact set of database-system parameters that can be tuned depends on the specific database system.
- ↳ Most database-system manuals provide information on what database-system parameters can be adjusted, and how you should choose value for the parameters.
- ↳ The third level is the highest level. It includes the schema and transactions.
- ↳ The administrator can tune the design

of the schema, the indices that are created, and the transactions that are executed, to improve performance. Tuning at this level is comparatively system independent.