# Query Evaluation and Optimization

## An Overview

---

## Steps in Query Processing

CISC 432/832

2

---

CISC 432/832

3

## Schema for Examples

Sailors (*sid*: integer, *sname*: string, *rating*: integer, *age*: real)
Reserves (*sid*: integer, *bid*: integer, *day*: dates, *rname*: string)

- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.

---

# Query Evaluation

---

## Overview

- *Evaluation Plan*: *Tree of R.A. ops, with alg for each op.*
  - Each operator typically implemented using a `pull' interface: when an operator is `pulled' for the next output tuples, it `pulls' on its inputs and computes them.
- Two main issues in query optimization:
  - For a given query, what plans are considered?
    - Need to search plan space for cheapest (estimated) plan.
  - How is the cost of a plan estimated?
- Ideally: Want to find best plan. Practically: Avoid worst plans!
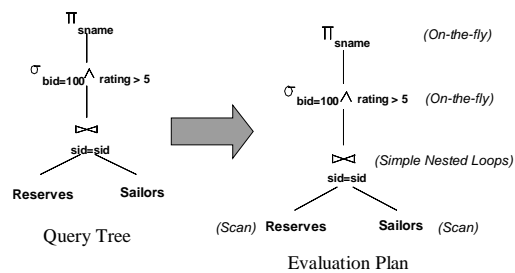- System R approach discussed in text.

## Example

SELECT S.sname
  FROM Reserves R, Sailors S
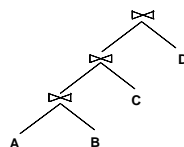  WHERE R.sid = S.sid
  AND R.bid = 100 AND S.rating > 5

*What is an equivalent relational algebra query?*

---

## Example (cont.)



Query Tree                    Evaluation Plan

$\pi_{sname}$

$\sigma_{bid=100 \wedge rating > 5}$

$\bowtie_{sid=sid}$

Reserves      Sailors

$\Pi_{sname}$   *(On-the-fly)*

$\sigma_{bid=100 \wedge rating > 5}$   *(On-the-fly)*

$\bowtie_{sid=sid}$   *(Simple Nested Loops)*

*(Scan)* Reserves      Sailors *(Scan)*

---

## Pipelined Evaluation

- Result of one operator *pipelined* to another without creating temporary table
- Lower overhead than *materialization*
- Unary operator is *on-the-fly* if input pipelined.



A      B
$\bowtie$
C
$\bowtie$
D
$\bowtie$

## Some Common Techniques

- Algorithms for evaluating relational operators use some simple ideas extensively:
  - Indexing: Can use WHERE conditions to retrieve small set of tuples (selections, joins)
  - Iteration: Sometimes, faster to scan all tuples even if there is an index. (And sometimes, we can scan the data entries in an index instead of the table itself.)
  - Partitioning: By using sorting or hashing, we can partition the input tuples and replace an expensive operation by similar operations on smaller inputs.

*Watch for these techniques as we discuss query evaluation!*

## Statistics and Catalogs

- Need information about the relations and indexes involved. **Catalogs** typically contain at least:
  - # tuples (NTuples) and # pages (NPages) for each relation.
  - # distinct key values (NKeys) and NPages for each index.
  - Index height, low/high key values (Low/High) for each tree index.
- Catalogs updated periodically.
  - Updating whenever data changes is too expensive; lots of approximation anyway, so slight inconsistency ok.
- More detailed information (e.g., histograms of the values in some field) are sometimes stored.

## Access Paths

- An <u>access path</u> is a method of retrieving tuples:
  - File scan, or index that matches a selection (in the query)

$(day<8/9/94$ AND $rname=$ 'Paul') OR $bid=5$ OR $sid=3$

- Selection conditions are first converted to *conjunctive normal form* (CNF):

$(day<8/9/94$ OR $bid=5$ OR $sid=3$ ) AND
$(rname=$ 'Paul' OR $bid=5$ OR $sid=3$)

## Access Paths (Cont.)

- A tree index _matches_ (a conjunction of) terms that involve only attributes in a _prefix_ of the search key.
  - E.g., Tree index on _<a, b, c>_ matches the selection _a=5 AND b=3_, and _a=5 AND b>6_, but not _b=3_.
- A hash index _matches_ (a conjunction of) terms that has a term _attribute = value_ for every attribute in the search key of the index.
  - E.g., Hash index on _<a, b, c>_ matches _a=5 AND b=3 AND c=5_; but it does not match _b=3_, or _a=5 AND b=3_, or _a>5 AND b=3 AND c=5_.

---

## One Approach to Selections

- Find the _most selective access path_, retrieve tuples using it, and apply any remaining terms that don't match the index:
  - _Most selective access path:_ An index or file scan that we estimate will require the fewest page I/Os.
  - Terms that match an index reduce the number of tuples _retrieved_; other terms are used to discard some retrieved tuples, but do not affect number of tuples/pages fetched.
  - Consider _day<8/9/94 AND bid=5 AND sid=3_. A B+ tree index on _day_ can be used; then, _bid=5_ and _sid=3_ must be checked for each retrieved tuple. Similarly, a hash index on _<bid, sid>_ could be used; _day<8/9/94_ must then be checked.

---

## Using an Index for Selections

- Cost depends on #qualifying tuples, and clustering.
  - Cost of finding qualifying data entries (typically small) plus cost of retrieving records (could be large w/o clustering).
  - In example, assuming uniform distribution of names, about 10% of tuples qualify (100 pages, 10000 tuples). With a clustered index, cost is little more than 100 I/Os; if unclustered, upto 10000 I/Os!

```
SELECT  *
FROM    Reserves R
WHERE   R.rname < 'C%'
```

## Projection

- The expensive part is removing duplicates.
  - SQL systems don't remove duplicates unless the keyword DISTINCT is specified in a query.
- Sorting Approach: Sort on <sid, bid> and remove duplicates. (Can optimize this by dropping unwanted information while sorting.)
- Hashing Approach: Hash on <sid, bid> to create partitions. Load partitions into memory one at a time, build in-memory hash structure, and eliminate duplicates.
- If there is an index with both R.sid and R.bid in the search key, may be cheaper to sort data entries!

| SELECT | DISTINCT |
|---|---|
| | R.sid, R.bid |
| FROM | Reserves R |

## Join: Nested Loops

foreach tuple r in R do
     foreach tuple s in S where $r_i == s_j$ do
       add <r, s> to result

- *Blocked Nested Loops*: Read a block of outer relation and match inner relation against its tuples
  - Cost: $B_R + (B_R * B_S)$

## Join: Nested Loops (cont.)

- *Indexed Nested Loops*: If there is an index on the join column of one relation (say S), can make it the inner and exploit the index.
  - Cost: $B_R + ((B_R * bfr_R) *$ cost of finding matching S tuples)
- For each R tuple, cost of probing S index is about 1.2 for hash index, 2-4 for B+ tree. Cost of then finding S tuples (assuming Alt. (2) or (3) for data entries) depends on clustering.
  - Clustered index: 1 I/O (typical), unclustered: upto 1 I/O per matching S tuple.

# Query Optimization

---

# Highlights of System R Optimizer

- Impact:
  - Most widely used currently; works well for < 10 joins.
- Cost estimation:  Approximate art at best.
  - Statistics, maintained in system catalogs, used to estimate cost of operations and result sizes.
  - Considers combination of CPU and I/O costs.
- Plan Space:  Too large, must be pruned.
  - Only the space of *left-deep plans* is considered.
    - Left-deep plans allow output of each operator to be *pipelined* into the next operator without storing it in a temporary relation.
  - Cartesian products avoided.

---

# Cost Estimation

- For each plan considered, must estimate cost:
  - Must estimate *cost* of each operation in plan tree.
    - Depends on input cardinalities.
    - We've already discussed how to estimate the cost of operations (sequential scan, index scan, joins, etc.)
  - Must also estimate *size of result* for each operation in tree!
    - Use information about the input relations.
    - For selections and joins, assume independence of predicates.

## Size Estimation and Reduction Factors

- Consider a query block:

  SELECT  attribute list
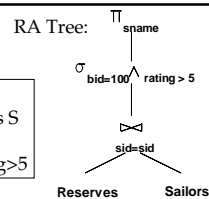  FROM  relation list
  WHERE  term1 AND ... AND termk

- Maximum # tuples in result is the product of the cardinalities of relations in the FROM clause.
- *Reduction factor (RF)* associated with each *term* reflects the impact of the *term* in reducing result size. *Result cardinality* = Max # tuples * product of all RF's.
  - Implicit assumption that *terms* are independent!
  - Term *col=value* has RF *1/NKeys(I)*, given index I on *col*
  - Term *col1=col2* has RF *1/MAX(NKeys(I1), NKeys(I2))*
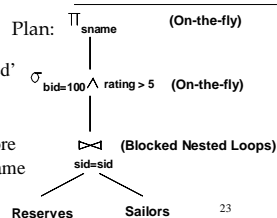  - Term *col>value* has RF *(High(I)-value)/(High(I)-Low(I))*

---

## Motivating Example

RA Tree:



SELECT  S.sname
FROM  Reserves R, Sailors S
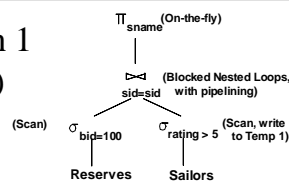WHERE  R.sid=S.sid AND
R.bid=100 AND S.rating>5

- Cost:  500+500*1000 I/Os
- By no means the worst plan!
- Misses several opportunities: selections could have been `pushed' earlier, no use is made of any available indexes, etc.
- *Goal of optimization:* To find more efficient plans that compute the same answer.

Plan:

---

## Alternative Plan 1 (Push Selects)



- Cost of plan:
  - Scan Reserves (1000) - produces 10 pages, if we have 100 boats, uniform distribution.
  - Scan Sailors (500) + write temp T1 (250 pages, if we have 10 ratings).
  - BNL: 10 * 250 = 2500
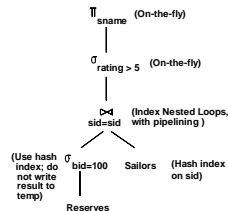  - Total:  1000 + 500 + 250 + 2500 =  4250 page I/Os.

# Alternative Plan 2
## (With Indexes)

$\Pi_{sname}$ (On-the-fly)

$\sigma_{rating > 5}$ (On-the-fly)

⋈ sid=sid (Index Nested Loops, with pipelining )

(Use hash index; do not write result to temp) $\sigma_{bid=100}$  Sailors (Hash index on sid)

Reserves

- With clustered index on *bid* of Reserves, we get 100,000/100 = 1000 tuples on 1000/100 = 10 pages.
- INL with ***pipelining*** (outer is not materialized).
  - –Projecting out unnecessary fields from outer doesn't help.
- ∨ Join column *sid* is a key for Sailors.
  - –At most one matching tuple, unclustered index on *sid* OK.
- ∨ Decision not to push *rating>5* before the join is based on availability of *sid* index on Sailors.
- ∨ Cost: Selection of Reserves tuples (10 I/Os); for each, must get matching Sailors tuple (1000*1.2); total 1210 I/Os.

---

# Summary

- There are several alternative evaluation algorithms for each relational operator.
- A query is evaluated by converting it to a tree of operators and evaluating the operators in the tree.
- Must understand query optimization in order to fully understand the performance impact of a given database design (relations, indexes) on a workload (set of queries).
- Two parts to optimizing a query:
  - Consider a set of alternative plans.
    - Must prune search space; typically, left-deep plans only.
  - Must estimate cost of each plan that is considered.
    - Must estimate size of result and cost for each plan node.
    - *Key issues*: Statistics, indexes, operator implementations.