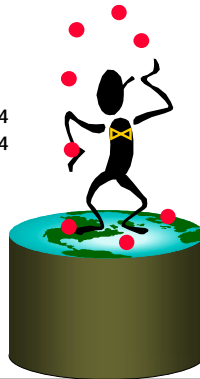


## Implementation of Relational Operations

15-415, Spring 2003, lecture 14  
R&G - Chapters 12 and 14

First comes thought; then organization of that thought, into ideas and plans; then transformation of those plans into reality. The beginning, as you will observe, is in your imagination.

Napolean Hill

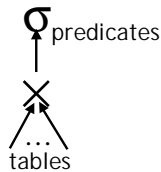


## Introduction

- We've covered the basic underlying storage, buffering, and indexing technology.
  - Now we can move on to query processing.
- Some database operations are **EXPENSIVE**
- Can greatly improve performance by being "smart"
  - e.g., can speed up 1,000,000x over naive approach
- Main weapons are:
  - clever implementation techniques for operators
  - exploiting "equivalencies" of relational operators
  - using statistics and cost models to choose among these.

## A Really Bad Query Optimizer

- For each Select-From-Where query block
  - Create a plan that:
    - Forms the cartesian product of the FROM clause
    - Applies the WHERE clause
    - Incredibly inefficient
      - Huge intermediate result!
- Then, as needed:
  - Apply the GROUP BY clause
  - Apply the HAVING clause
  - Apply any projections and output expressions
  - Apply duplicate elimination and/or ORDER BY

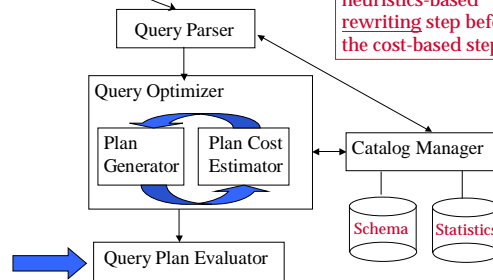


## Cost-based Query Sub-System

Queries  

```
Select *
From Blah B
Where B.blah = blah
```

Usually there is a heuristics-based rewriting step before the cost-based steps.



## The Query Optimization Game

- "Optimizer" is a bit of a misnomer...
- Goal is to pick a "good" (i.e., low expected cost) plan.
  - Involves choosing access methods, physical operators, operator orders, ...
  - Notion of cost is based on an abstract "cost model"
- Roadmap for this topic:
  - First: basic operators
  - Then: joins
  - After that: optimizing multiple operators

## Relational Operations

- We will consider how to implement:
  - Selection** ( $\sigma$ ) Selects a subset of rows from relation.
  - Projection** ( $\pi$ ) Deletes unwanted columns from relation.
  - Join** ( $\times$ ) Allows us to combine two relations.
  - Set-difference** ( $-$ ) Tuples in reln. 1, but not in reln. 2.
  - Union** ( $\cup$ ) Tuples in reln. 1 and in reln. 2.
  - Aggregation** (SUM, MIN, etc.) and GROUP BY
- Since each op returns a relation, ops can be **composed!** After we cover the operations, we will discuss how to **optimize** queries formed by composing them.



## Schema for Examples

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

- Similar to old schema; *rname* added for variations.
- Sailors:
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
- Reserves:
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.



## Simple Selections

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

- Of the form  $\sigma_{R.attr op value}(R)$
- Question: how best to perform? Depends on:
  - what indexes/access paths are available
  - what is the expected size of the result (in terms of number of tuples and/or number of pages)
- Size of result approximated as  
*size of R \* reduction factor*
  - “reduction factor” is usually called *selectivity*.
  - estimate of reduction factors is based on statistics – we will discuss shortly.



## Alternatives for Simple Selections

- With no index, unsorted:
  - Must essentially scan the whole relation
  - cost is  $M$  (#pages in R). For “reserves” = 1000 I/Os.
- With no index, sorted:
  - cost of binary search + number of pages containing results.
  - For reserves =  $10 \text{ I/Os} + \lceil \text{selectivity} * \# \text{pages} \rceil$
- With an index on selection attribute:
  - Use index to find qualifying data entries,
  - then retrieve corresponding data records.
  - (Hash index useful only for equality selections.)



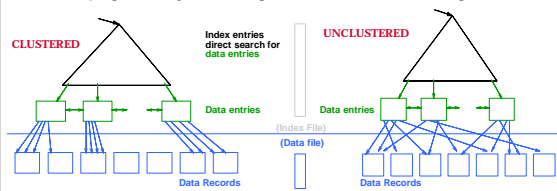
## Using an Index for Selections

- Cost depends on #qualifying tuples, and clustering.
  - Cost:
    - finding qualifying data entries (typically small)
    - plus cost of retrieving records (could be large w/o clustering).
  - In example “reserves” relation, if 10% of tuples qualify (100 pages, 10000 tuples).
    - With a *clustered* index, cost is little more than 100 I/Os;
    - if *unclustered*, could be up to 10000 I/Os! unless...



## Selections using Index (cont)

- Important refinement for unclustered indexes:
  1. Find qualifying data entries.
  2. Sort the rid's of the data records to be retrieved.
  3. Fetch rids in order. This ensures that each data page is looked at just once (though # of such pages likely to be higher than with clustering).



## General Selection Conditions

- \*  $(day < 8/9/94 \text{ AND } rname = \text{'Paul'}) \text{ OR } bid = 5 \text{ OR } sid = 3$
- Such selection conditions are first converted to *conjunctive normal form (CNF)*:
  - $(day < 8/9/94 \text{ OR } bid = 5 \text{ OR } sid = 3) \text{ AND } (rname = \text{'Paul'} \text{ OR } bid = 5 \text{ OR } sid = 3)$
- We only discuss the case with no ORs (a conjunction of terms of the form *attr op value*).
- A *B-tree* index *matches* (a conjunction of) terms that involve only attributes in a *prefix* of the search key.
  - Index on  $\langle a, b, c \rangle$  matches  $a=5 \text{ AND } b=3$ , but not  $b=3$ .
- For *Hash* index, must have all attributes in search key



## Two Approaches to General Selections

- **First approach:** 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** this 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.



## Most Selective Index - Example

- 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  $\langle bid, sid \rangle$  could be used;
  - Then,  $day < 8/9/94$  must be checked.
- **How about a B+tree on  $\langle rname, day \rangle$ ?**
- **How about a B+tree on  $\langle day, rname \rangle$ ?**
- **How about a Hash index on  $\langle day, rname \rangle$ ?**



## Intersection of Rids

- **Second approach:** if we have 2 or more matching indexes (w/Alternatives (2) or (3) for data entries):
  - Get **sets of rids** of data records using **each** matching index.
  - Then **intersect** these **sets of rids**.
  - Retrieve the records and apply any remaining terms.
  - Consider  $day < 8/9/94$  AND  $bid=5$  AND  $sid=3$ . With a **B+ tree index on *day*** and an **index on *sid***, we can retrieve rids of records satisfying  $day < 8/9/94$  using the first, rids of recs satisfying  $sid=3$  using the second, **intersect**, retrieve records and check  $bid=5$ .
  - Note: commercial systems use various tricks to do this:
    - bit maps, bloom filters, index joins



## The Halloween Problem – An Aside.

- Story from the early days of System R.
- While testing the optimizer on 10/31/75(?), the following update was run:

```
UPDATE payroll
SET salary = salary*1.1
WHERE salary > 20K;
```



- AND IT NEVER STOPPED!
- Can you guess why???(hint: it was an optimizer bug...)



## Schema for Examples

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

- Similar to old schema; *rname* added for variations.
- **Sailors:**
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - $N=500$ ,  $p_s=80$ .
- **Reserves:**
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - $M=1000$ ,  $p_r=100$ .



## The Projection Operation

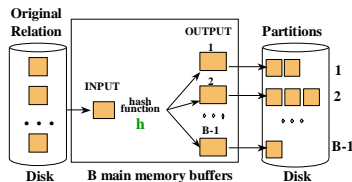
```
SELECT DISTINCT
  R.sid, R.bid
FROM Reserves R
```

- Issue is removing **duplicates**.
- **Basic approach is to use sorting**
  - 1. Scan R, extract only the needed attrs (why do this first?)
  - 2. Sort the resulting set
  - 3. Remove adjacent duplicates
  - **Cost:** Reserves with size ratio 0.25 = 250 pages. With 20 buffer pages can sort in 2 passes, so  $1000 + 250 + 2 * 2 * 250 + 250 = 2500$  I/Os
- **Can improve by modifying external sort algorithm (see chapter 13):**
  - **Modify Pass 0 of external sort to eliminate unwanted fields.**
  - **Modify merging passes to eliminate duplicates.**
  - **Cost:** for above case: read 1000 pages, write out 250 in runs of 40 pages, merge runs =  $1000 + 250 + 250 = 1500$ .



## Projection Based on Hashing

- **Partitioning phase:** Read R using one input buffer. For each tuple, discard unwanted fields, apply hash function  $h1$  to choose one of B-1 output buffers.
  - Result is B-1 partitions (of tuples with no unwanted fields).
  - 2 tuples from different partitions guaranteed to be distinct.



## DupElim with Hashing (continued)

- **Duplicate elimination phase:** For each partition, read it and build an in-memory hash table, using hash fn  $h2 (< > h1)$  on all fields, while discarding duplicates.
  - If partition does not fit in memory, can apply hash-based projection algorithm recursively to this partition.
- **Cost: For previous case**
  - assuming partitions fit in memory (i.e. #bufs  $\geq$  square root of the #of pages of projected tuples)
  - read 1000 pages and write out partitions of projected tuples (250 pages)
  - Do dup elim on each partition (total 250 page reads)
  - Total : 1500 I/Os.



## Discussion of Projection

- Sort-based approach is the standard; better handling of **skew** and result is **sorted**.
- If enough buffers, both have same I/O cost:  $M + 2T$  where M is #pgs in R, T is #pgs of R with unneeded attributes removed.
  - Although many systems don't use the specialized sort.
- If an index on the relation contains all wanted attributes in its search key, can do **index-only** scan.
  - Apply projection techniques to data entries (much smaller!)
- If an ordered (i.e., tree) index contains all wanted attributes as **prefix** of search key, can do even better:
  - Retrieve data entries in order (index-only scan), discard unwanted fields, compare adjacent tuples to check for duplicates.



## Joins

- Joins are very common.
- Joins can be very expensive (cross product in worst case).
- Many approaches to reduce join cost.



## Equality Joins With One Join Column

```
SELECT *
FROM Reserves R1, Sailors S1
WHERE R1.sid=S1.sid
```

- In algebra:  $R \bowtie S$ . Common! Must be carefully optimized.  $R \times S$  is large; so,  $R \times S$  followed by a selection is inefficient.
- Remember, join is associative and commutative.
- Assume:
  - M pages in R,  $p_R$  tuples per page.
  - N pages in S,  $p_S$  tuples per page.
  - In our examples, R is Reserves and S is Sailors.
- We will consider more complex join conditions later.
- **Cost metric:** # of I/Os. We will ignore output costs.



## Simple Nested Loops Join

```
foreach tuple r in R do
  foreach tuple s in S do
    if r1 == s1 then add <r, s> to result
```

- For each tuple in the **outer** relation R, we scan the entire **inner** relation S.
- How much does this Cost?
- $(p_R * M) * N + M = 100 * 1000 * 500 + 1000$  I/Os.
  - At 10ms/I/O, Total: ???
- What if smaller relation (S) was outer?
- What assumptions are being made here?

**Q: What is cost if one relation can fit entirely in memory?**



## Page-Oriented Nested Loops Join

```

foreach page  $b_R$  in R do
  foreach page  $b_S$  in S do
    foreach tuple  $r$  in  $b_R$  do
      foreach tuple  $s$  in  $b_S$  do
        if  $r_i == s_j$  then add  $\langle r, s \rangle$  to result
  
```

- For each *page* of R, get each *page* of S, and write out matching pairs of tuples  $\langle r, s \rangle$ , where  $r$  is in R-page and  $S$  is in S-page.
- What is the cost of this approach?
- $M * N + M = 1000 * 500 + 1000$ 
  - If smaller relation (S) is outer, cost =  $500 * 1000 + 500$



## Index Nested Loops Join

```

foreach tuple  $r$  in R do
  foreach tuple  $s$  in S where  $r_i == s_j$  do
    add  $\langle r, s \rangle$  to result
  
```

- If there is an index on the join column of one relation (say S), can make it the inner and exploit the index.
  - Cost:  $M + (M * p_R) * \text{cost of finding matching S tuples}$
- 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 per page of matching S tuples.
- **Unclustered:** up to 1 I/O per matching S tuple.



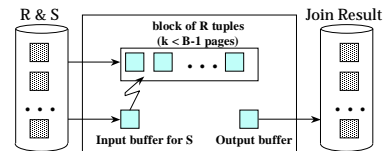
## Examples of Index Nested Loops

- **Hash-index (Alt. 2) on *sid* of Sailors (as inner):**
  - Scan Reserves: 1000 page I/Os,  $100 * 1000$  tuples.
  - For each Reserves tuple: 1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching Sailors tuple. **Total:**
- **Hash-index (Alt. 2) on *sid* of Reserves (as inner):**
  - Scan Sailors: 500 page I/Os,  $80 * 500$  tuples.
  - For each Sailors tuple: 1.2 I/Os to find index page with data entries, plus cost of retrieving matching Reserves tuples. **Assuming uniform distribution**, 2.5 reservations per sailor ( $100,000 / 40,000$ ). Cost of retrieving them is 1 or 2.5 I/Os depending on whether the index is clustered.
  - **Totals:**



## Block Nested Loops Join

- Page-oriented NL doesn't exploit extra buffers.
- **Alternative approach:** Use one page as an input buffer for scanning the inner S, one page as the output buffer, and use all remaining pages to hold "block" of outer R.
- For each matching tuple  $r$  in R-block,  $s$  in S-page, add  $\langle r, s \rangle$  to result. Then read next R-block, scan S, etc.



## Examples of Block Nested Loops

- **Cost:** Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \# \text{ of pages of outer} / \text{blocksize} \rceil$
- **With Reserves (R) as outer, and 100 pages of R:**
  - Cost of scanning R is 1000 I/Os; a total of 10 blocks.
  - Per block of R, we scan Sailors (S);  $10 * 500$  I/Os.
  - If space for just 90 pages of R, we would scan S 12 times.
- **With 100-page block of Sailors as outer:**
  - Cost of scanning S is 500 I/Os; a total of 5 blocks.
  - Per block of S, we scan Reserves;  $5 * 1000$  I/Os.
- With **sequential reads** considered, analysis changes: may be best to divide buffers evenly between R and S.



## Sort-Merge Join ( $R \bowtie_{i=j} S$ )

- Sort R and S on the join column, then scan them to do a "merge" (on join col.), and output result tuples.
- **Useful if**
  - one or both inputs are already sorted on join attribute(s)
  - output is required to be sorted on join attributes(s)
- "Merge" phase can require some back tracking if duplicate values appear in join column
- R is scanned once; each S group is scanned once per matching R tuple. (Multiple scans of an S group will probably find needed pages in buffer.)



## Example of Sort-Merge Join

sid	sname	rating	age	sid	bid	day	rname
22	dustin	7	45.0	28	103	12/4/96	guppy
28	yuppy	9	35.0	28	103	11/3/96	yuppy
31	lubber	8	55.5	31	101	10/10/96	dustin
44	guppy	5	35.0	31	102	10/12/96	lubber
58	rusty	10	35.0	31	101	10/11/96	lubber
				58	103	11/12/96	dustin

- **Cost: Sort R + Sort S + (M+N)**
  - The cost of scanning, M+N, could be M\*N (very unlikely!)
- **With 35, 100 or 300 buffer pages, both Reserves and Sailors can be sorted in 2 passes; total join cost: 7500.**

(BNL cost: 2500 to 15000 I/Os)



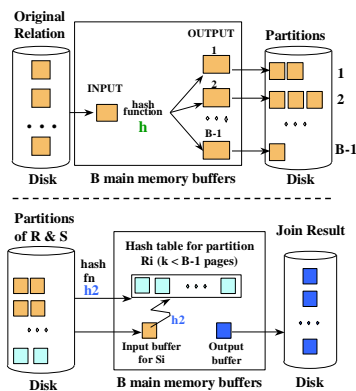
## Refinement of Sort-Merge Join

- We can combine the merging phases in the *sorting* of R and S with the merging required for the join.
  - Allocate 1 page per run of each relation, and 'merge' while checking the join condition
  - With  $B > \sqrt{L}$ , where L is the size of the larger relation, using the sorting refinement that produces runs of length 2B in Pass 0, #runs of each relation is  $< B/2$ .
  - **Cost:** read+write each relation in Pass 0 + read each relation in (only) merging pass (+ writing of result tuples).
  - In example, cost goes down from 7500 to 4500 I/Os.
- In practice, cost of sort-merge join, like the cost of external sorting, is *linear*.



## Hash-Join

- Partition both relations using hash fn h: R tuples in partition i will **only** match S tuples in partition i.



- Read in a partition of R, hash it using  $h_2$  ( $< h$ ). Scan matching partition of S, probe hash table for matches.



## Observations on Hash-Join

- #partitions  $k < B$ , and  $B-1 >$  size of largest partition to be held in memory. Assuming uniformly sized partitions, and maximizing k, we get:
  - $k = B-1$ , and  $M/(B-1) < B-2$ , i.e., B must be  $> \sqrt{M}$
- Since we build an in-memory hash table to speed up the matching of tuples in the second phase, a little more memory is needed.
- If the hash function does not partition uniformly, one or more R partitions may not fit in memory. Can apply hash-join technique recursively to do the join of this R-partition with corresponding S-partition.



## Cost of Hash-Join

- In partitioning phase, read+write both relns;  $2(M+N)$ . In matching phase, read both relns;  $M+N$  I/Os.
- In our running example, this is a total of 4500 I/Os.
- Sort-Merge Join vs. Hash Join:
  - Given a minimum amount of memory (*what is this, for each?*) both have a cost of  $3(M+N)$  I/Os. Hash Join superior on this count if relation sizes differ greatly. Also, Hash Join shown to be highly parallelizable.
  - Sort-Merge less sensitive to data skew; result is sorted.



## General Join Conditions

- Equalities over several attributes (e.g.,  $R.sid=S.sid$  AND  $R.rname=S.rname$ ):
  - For Index NL, build index on  $\langle sid, rname \rangle$  (if S is inner); or use existing indexes on *sid* or *rname*.
  - For Sort-Merge and Hash Join, sort/partition on combination of the two join columns.
- Inequality conditions (e.g.,  $R.rname < S.rname$ ):
  - For Index NL, need (clustered!) B+ tree index.
    - Range probes on inner; # matches likely to be much higher than for equality joins.
  - Hash Join, Sort Merge Join not applicable!
  - Block NL quite likely to be the best join method here.



## Set Operations

- **Intersection and cross-product special cases of join.**
- **Union (Distinct) and Except similar; we'll do union.**
- **Sorting based approach to union:**
  - Sort both relations (on combination of all attributes).
  - Scan sorted relations and merge them.
  - *Alternative:* Merge runs from Pass 0 for *both* relations.
- **Hash based approach to union:**
  - Partition R and S using hash function  $h$ .
  - For each S-partition, build in-memory hash table (using  $h2$ ), scan corr. R-partition and add tuples to table while discarding duplicates.



## Aggregate Operations (AVG, MIN, etc.)

- **Without grouping:**
  - In general, requires scanning the relation.
  - Given index whose search key includes all attributes in the SELECT or WHERE clauses, can do index-only scan.
- **With grouping:**
  - Sort on group-by attributes, then scan relation and compute aggregate for each group. (Can improve upon this by combining sorting and aggregate computation.)
  - Similar approach based on hashing on group-by attributes.
  - Given tree index whose search key includes all attributes in SELECT, WHERE and GROUP BY clauses, can do index-only scan; if group-by attributes form prefix of search key, can retrieve data entries/tuples in group-by order.



## Impact of Buffering

- If several operations are executing concurrently, estimating the number of available buffer pages is guesswork.
- Repeated access patterns interact with buffer replacement policy.
  - e.g., Inner relation is scanned repeatedly in Simple Nested Loop Join. With enough buffer pages to hold inner, replacement policy does not matter. Otherwise, MRU is best, LRU is worst (*sequential flooding*).
  - Does replacement policy matter for Block Nested Loops?
  - What about Index Nested Loops? Sort-Merge Join?



## Summary

- A virtue of relational DBMSs: *queries are composed of a few basic operators*; the implementation of these operators can be carefully tuned (and it is important to do this!).
- Many alternative implementation techniques for each operator; no universally superior technique for most operators.
- Must consider available alternatives for each operation in a query and choose best one based on system statistics, etc. This is part of the broader task of optimizing a query composed of several ops.