

## Relational Query Optimization

CS186  
R & G Chapters 12/15



## Review

- **Choice of single-table operations**
  - Depends on indexes, memory, stats,...
- **Joins**
  - Blocked nested loops:
    - simple, exploits extra memory
  - Indexed nested loops:
    - best if 1 rel small and one indexed
  - Sort/Merge Join
    - good with small amount of memory, bad with duplicates
  - Hash Join
    - fast (enough memory), bad with skewed data
- **These are “rules of thumb”**
  - On their way to a more principled approach...

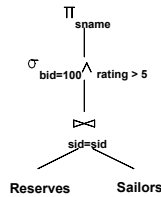


## Query Optimization Overview

- Query can be converted to relational algebra
- Rel. Algebra converts to tree, joins form branches
- Each operator has implementation choices
- Operators can also be applied in different order!

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

$\pi_{(sname)} \sigma_{(bid=100 \wedge rating > 5)}$   
(Reserves  $\bowtie$  Sailors)

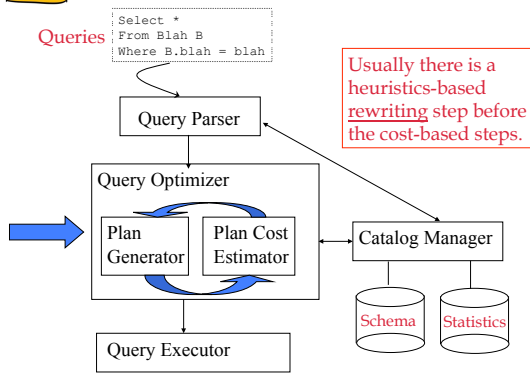


## Query Optimization Overview (cont.)

- **Plan:** Tree of R.A. ops (and some others) with choice of algorithm for each op.
  - Recall: Iterator interface (next()!)
- **Three main issues:**
  - For a given query, what plans are considered?
  - How is the cost of a plan estimated?
  - How do we “search” in the “plan space”?
- **Ideally:** Want to find best plan.
- **Reality:** Avoid worst plans!



## Cost-based Query Sub-System



## Let's go through some examples

- Just to get a flavor...



## Schema for Examples

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

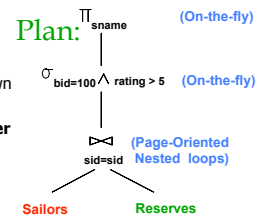
- As seen in previous lectures...
- **Reserves:**
  - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
  - Assume there are 100 boats
- **Sailors:**
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
  - Assume there are 10 different ratings
- **Assume we have 5 pages in our buffer pool!**



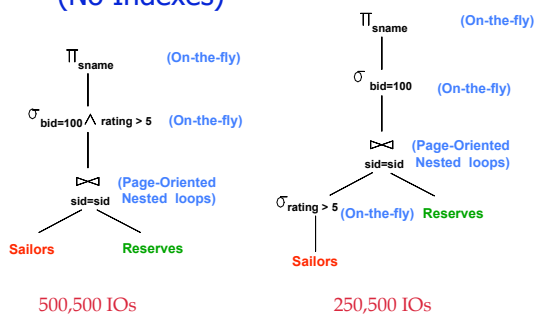
## Motivating Example

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

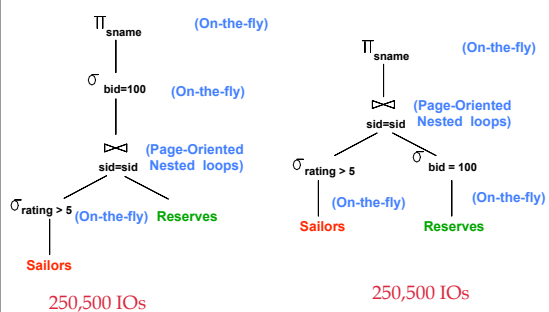
- **Cost: 500+500\*1000 I/Os**
- **By no means the worst plan!**
- **Misses several opportunities:**
  - selections could be 'pushed' down
  - no use made of indexes
- **Goal of optimization: Find faster plans that compute the same answer.**



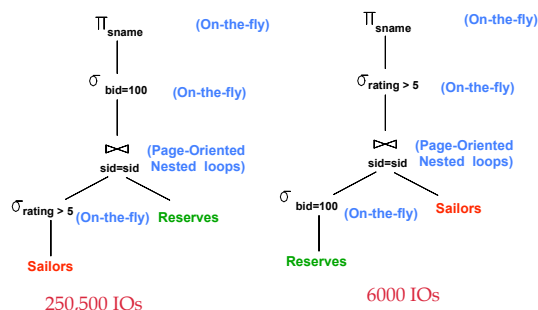
## Alternative Plans – Push Selects (No Indexes)



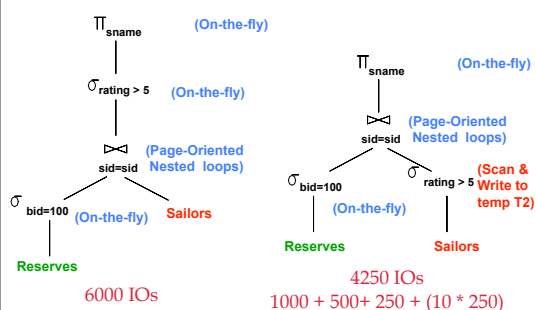
## Alternative Plans – Push Selects (No Indexes)

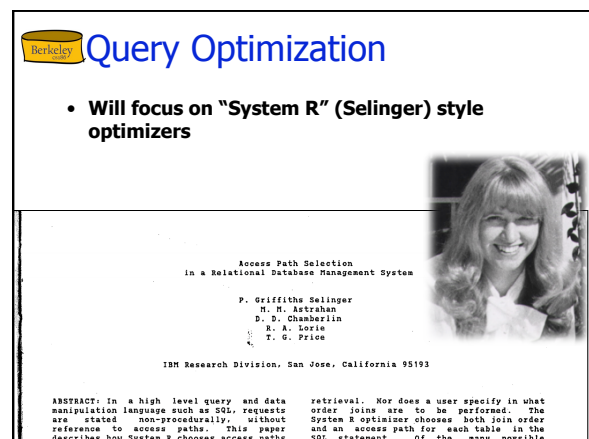
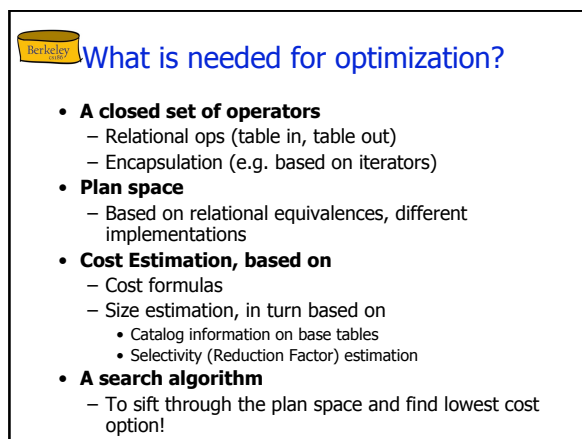
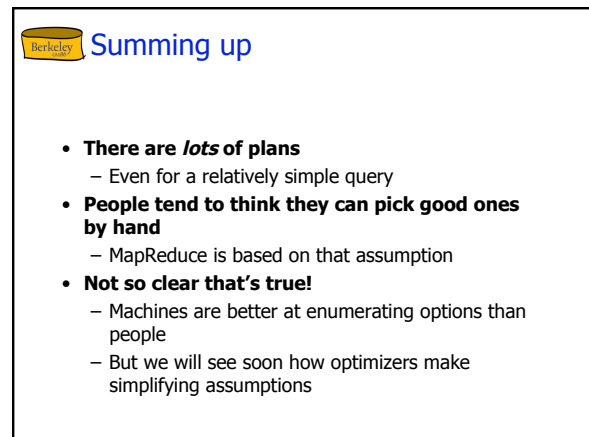
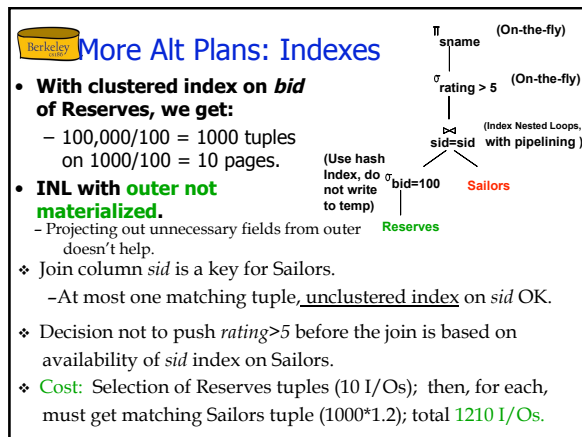
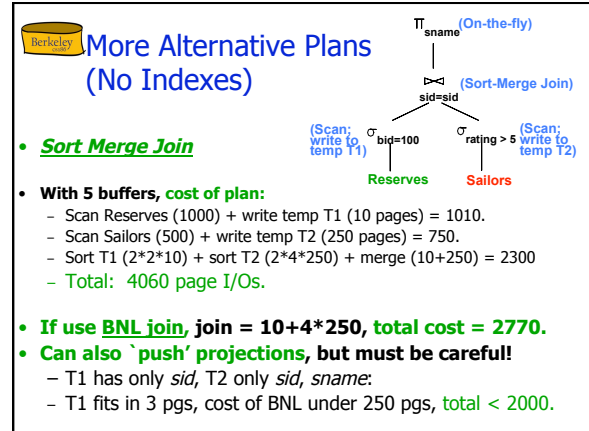
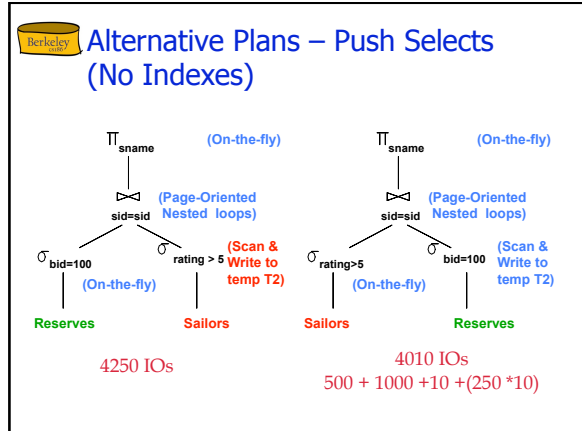


## Alternative Plans – Push Selects (No Indexes)



## Alternative Plans – Push Selects (No Indexes)







## Highlights of System R Optimizer

- **Impact:**
  - Most widely used currently; works well for 10-15 joins.
- **Cost estimation:**
  - Very inexact, but works ok in practice.
  - Statistics in system catalogs used to estimate cost of operations and result sizes.
  - Considers combination of CPU and I/O costs.
  - System R's scheme has been improved since that time.
- **Plan Space: Too large, must be pruned.**
  - Many plans share common, "overpriced" subtrees
    - ignore them all!
  - In some implementations, only the space of *left-deep plans* is considered.
  - Cartesian products avoided in some implementations.



## Query Blocks: Units of Optimization

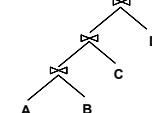
- Break query into *query blocks*
- Optimized one block at a time
- Uncorrelated nested blocks computed once
- Correlated nested blocks like function calls
  - But sometimes can be "decorrelated"
  - Beyond the scope of CS186!

```
SELECT S.sname
FROM Sailors S
WHERE S.age IN
  (SELECT MAX (S2.age)
   FROM Sailors S2
   GROUP BY S2.rating)
```

Outer block    Nested block

✧ For each block, the plans considered are:

- All available *access methods*, for each relation in FROM clause.
- All *left-deep join trees*
  - right branch always a base table
  - consider all join *orders* and join *methods*



## 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. 100 distinct bids.
- **Sailors:**
  - Each tuple is 50 bytes long, 80 tuples per page, 500 pages. 10 ratings, 40,000 sids.



## Translating SQL to Relational Algebra

```
SELECT S.sid, MIN (R.day)
FROM Sailors S, Reserves R, Boats B
WHERE S.sid = R.sid AND R.bid = B.bid AND B.color = "red"
GROUP BY S.sid
HAVING COUNT (*) >= 2
```

For each sailor with at least two reservations for red boats, find the sailor id and the earliest date on which the sailor has a reservation for a red boat.



## Translating SQL to Relational Algebra

```
SELECT S.sid, MIN (R.day)
FROM Sailors S, Reserves R, Boats B
WHERE S.sid = R.sid AND R.bid = B.bid AND B.color = "red"
GROUP BY S.sid
HAVING COUNT (*) >= 2
```

$$\pi_{S.sid, MIN(R.day)} \left( \left( HAVING COUNT(*) > 2 \left( GROUP BY_{S.sid} \left( \sigma_{B.color = "red"} (Sailors \bowtie Reserves \bowtie Boats) \right) \right) \right) \right)$$


## Relational Algebra Equivalences

- **Allow us to choose different join orders and to 'push' selections and projections ahead of joins.**

- **Selections:**

- $\sigma_{c1 \wedge \dots \wedge cn}(R) \equiv \sigma_{c1}(\dots(\sigma_{cn}(R))\dots)$  (*cascade*)
- $\sigma_{c1}(\sigma_{c2}(R)) \equiv \sigma_{c2}(\sigma_{c1}(R))$  (*commute*)

- **Projections:**

- $\pi_{a1}(R) \equiv \pi_{a1}(\dots(\pi_{a1, \dots, an}(R))\dots)$  (*cascade*)

- **Cartesian Product**

- $R \times (S \times T) \equiv (R \times S) \times T$  (*associative*)
- $R \times S \equiv S \times R$  (*commutative*)
- *This means we can do joins in any order.*
  - But...beware of cartesian product!





## Backend/optimizer/path/clausesel.c

```

/*
 *
 * THIS IS A HACK TO GET V4 OUT THE DOOR.
 * -- JMH 7/9/92
 */
s1 = (Selectivity) 0.3333333;

```

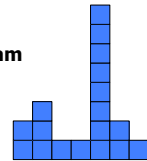


## Reduction Factors & Histograms

- For better estimation, use a histogram

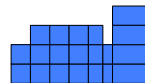
*equiwidth*

| No. of Values | 2     | 3      | 3      | 1      | 8      | 2      | 1      |
|---------------|-------|--------|--------|--------|--------|--------|--------|
| Value         | 0-.99 | 1-1.99 | 2-2.99 | 3-3.99 | 4-4.99 | 5-5.99 | 6-6.99 |



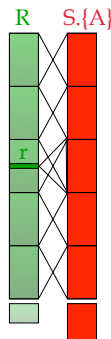
*equidepth*

| No. of Values | 2     | 3      | 3      | 3      | 3         | 2         | 4      |
|---------------|-------|--------|--------|--------|-----------|-----------|--------|
| Value         | 0-.99 | 1-1.99 | 2-2.99 | 3-4.05 | 4.06-4.67 | 4.68-4.99 | 5-6.99 |



## Think through estimation for joins

- Term  $col1=col2$ 
  - RF =  $1/MAX(NKeys(I1), NKeys(I2))$
- Q: Given R join S, what is result cardinality?
  - Two cases:
    - Key for R, Foreign Key for S
      - A common case, treat it specially! RF =  $1/|R|$
    - join on non-key {A}
      - For each  $r \in R$ , NTuples(S)/NKeys(A,S) result tuples so... (NTuples(R) \* NTuples(S)) / NKeys(A,S)
      - For each  $s \in S$ , NTuples(R)/NKeys(A,R) result tuples so... (NTuples(S) \* NTuples(R)) / NKeys(A,R)
  - If these two estimates differ, take the lower one!
    - Q: Why?



## Enumeration of Alternative Plans

- There are two main cases:
  - Single-relation plans (base case)
  - Multiple-relation plans (induction)
- Single-table queries include selects, projects, and grouping/aggregate ops:
  - Consider each available access path (file scan / index)
    - Choose the one with the least estimated cost
  - Selection/Projection done on the fly
  - Result pipelined into grouping/aggregation



## Cost Estimates for Single-Relation Plans

- Index I on primary key matches selection:
  - Cost is  $Height(I)+1$  for a B+ tree.
- Clustered index I matching one or more selects:
  - $(NPages(I)+NPAGES(R))$  \* product of RF's of matching selects.
- Non-clustered index I matching one or more selects:
  - $(NPages(I)+NTuples(R))$  \* product of RF's of matching selects.
- Sequential scan of file:
  - $NPAGES(R)$ .

☞ Recall: Must also charge for duplicate elimination if required



## Example

```

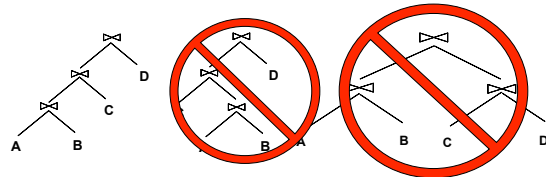
SELECT S.sid
FROM Sailors S
WHERE
  S.rating=8

```

- If we have an index on rating:
  - Cardinality =  $(1/NKeys(I)) * NTuples(R) = (1/10) * 40000$  tuples
  - Clustered index:  $(1/NKeys(I)) * (NPAGES(I)+NPAGES(R))$ 
    - $= (1/10) * (50+500) = 55$  pages are retrieved. (This is the **cost**.)
  - Unclustered index:  $(1/NKeys(I)) * (NPAGES(I)+NTuples(R))$ 
    - $= (1/10) * (50+40000) = 401$  pages are retrieved.
- If we have an index on sid:
  - Would have to retrieve all tuples/pages. With a clustered index, the cost is  $50+500$ , with unclustered index,  $50+40000$ .
- Doing a file scan:
  - We retrieve all file pages (500).

## Berkeley Queries Over Multiple Relations

- **A System R heuristic:**  
**only left-deep join trees considered.**
  - Restricts the search space
  - Left-deep trees allow us to generate all *fully pipelined plans*.
    - Intermediate results not written to temporary files.
    - Not all left-deep trees are fully pipelined (e.g., SM join).



## Berkeley Enumeration of Left-Deep Plans

- **Left-deep plans differ in**
  - the order of relations
  - the access method for each relation
  - the join method for each join.
- **Enumerated using N passes (if N relations joined):**
  - **Pass 1:** Find best 1-relation plan for each relation.
  - **Pass i:** Find best way to join result of an (i-1)-relation plan (as outer) to the i-th relation. (i between 2 and N.)
- **For each subset of relations, retain only:**
  - Cheapest plan overall, plus
  - Cheapest plan for each *interesting order* of the tuples.

## Berkeley The Dynamic Programming Table

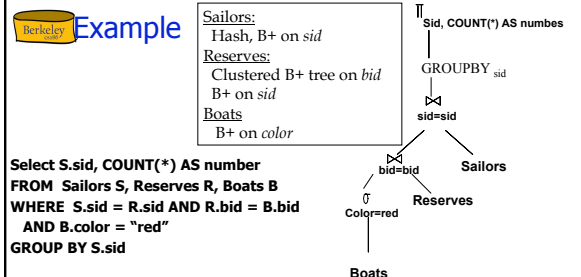
| Subset of tables in FROM clause | Interesting-order columns | Best plan      | Cost |
|---------------------------------|---------------------------|----------------|------|
| {R, S}                          | <none>                    | hashjoin(R,S)  | 1000 |
| {R, S}                          | <R.a, S.b>                | sortmerge(R,S) | 1500 |

## Berkeley A Note on "Interesting Orders"

- **An intermediate result has an "interesting order" if it is sorted by any of:**
  - ORDER BY attributes
  - GROUP BY attributes
  - Join attributes of *yet-to-be-added* (downstream) joins

## Berkeley Enumeration of Plans (Contd.)

- **Match an i-1 way plan with another table only if**
  - a) there is a join condition between them, or
  - b) all predicates in WHERE have been used up.
    - i.e., avoid Cartesian products if possible.
- **ORDER BY, GROUP BY, aggregates etc. handled as a final step**
  - via 'interestingly ordered' plan if chosen (free!)
  - or via an additional sort/hash operator
- **Despite pruning, this is exponential in #tables.**
- **Recall: in practice, COST considered is**  
**#IOs + factor \* CPU Inst**



- **Pass1: Best plan(s) for accessing each relation**
  - Reserves, Sailors: File Scan
  - Q: What about Clustered B+ on Reserves.bid???
  - Boats: B+ tree on color



## Pass 1

- **Find best plan for each relation in isolation:**
  - Reserves, Sailors: File Scan
  - Boats: B+ tree on color



## Pass 2

- **For each plan in pass 1, generate plans joining another relation as the inner, using all join methods (and matching inner access methods)**
  - File Scan Reserves (outer) with Boats (inner)
  - File Scan Reserves (outer) with Sailors (inner)
  - File Scan Sailors (outer) with Boats (inner)
  - File Scan Sailors (outer) with Reserves (inner)
  - Boats Btree on color with Sailors (inner)
  - Boats Btree on color with Reserves (inner)
- **Retain cheapest plan for each (pair of relations, order)**



## Pass 3 and beyond

- **Using Pass 2 plans as outer relations, generate plans for the next join**
  - E.g. Boats B+-tree on color with Reserves (bid) (sortmerge)  
inner Sailors (B-tree sid) sort-merge
- **Then, add cost for groupby/aggregate:**
  - This is the cost to sort the result by sid, *unless it has already been sorted by a previous operator.*
- **Then, choose the cheapest plan**



## Physical DB Design

- **Query optimizer does what it can to use indices, clustering etc.**
- **DataBase Administrator (DBA) is expected to set up physical design well**
- **Good DBAs understand query optimizers very well**



## Decisions to Make

- **What indexes should we create?**
  - Which tables, which field(s) should be the search key?  
Multiple indexes?
- **Index Clustering?**
- **Horizontal partitioning, replication, views ...**



## Index Selection

- **One approach:**
  - Consider most important queries in turn.
  - Consider best plan using the current indexes, and see if better plan is possible with an additional index.
  - If so, create it.
- **Before creating an index, must also consider the impact on updates in the workload!**
  - Trade-off: indexes can make queries go faster, updates slower. Require disk space, too.





## Issues to Consider in Index Selection

- **Attributes mentioned in a WHERE clause are candidates for index search keys.**
  - Range conditions are sensitive to clustering
  - Exact match conditions don't require clustering
    - Or do they???? :-)
- **Choose indexes that benefit many queries**
- **NOTE: only one index can be clustered per relation!**
  - So choose it wisely!



## Example 1

```
SELECT E.ename, D.mgr
FROM Emp E, Dept D
WHERE E.dno=D.dno AND D.dname='Toy'
```

- **B+ tree index on *D.dname* supports 'Toy' selection.**
  - Given this, index on *D.dno* is not needed.
- **B+ tree index on *E.dno* allows us to get matching (inner) Emp tuples for each selected (outer) Dept tuple.**
- **What if WHERE included: ``... AND E.age=25'' ?**
  - Could retrieve Emp tuples using index on *E.age*, then join with Dept tuples satisfying *dname* selection. Comparable to strategy that used *E.dno* index.
  - So, if *E.age* index is already created, this query provides much less motivation for adding an *E.dno* index.



## Example 2

```
SELECT E.ename, D.mgr
FROM Emp E, Dept D
WHERE E.sal BETWEEN 10000 AND 20000
AND E.hobby='Stamps' AND E.dno=D.dno
```

- **All selections are on Emp so it should be the outer relation in any Index NL join.**
  - Suggests that we build a B+ tree index on *D.dno*.
- **What index should we build on Emp?**
  - B+ tree on *E.sal* could be used, OR an index on *E.hobby* could be used. Only one of these is needed, and which is better depends upon the selectivity of the conditions.
    - As a rule of thumb, equality selections more selective than range selections.
- **As both examples indicate, our choice of indexes is guided by the plan(s) that we expect an optimizer to consider for a query. Have to understand optimizers!**



## Examples of Clustering

```
SELECT E.dno
FROM Emp E
WHERE E.age>40
```

- **B+ tree index on *E.age* can be used to get qualifying tuples.**
  - How selective is the condition?
  - Is the index clustered?
- **Consider the GROUP BY query.**
  - If many tuples have *E.age* > 10, using *E.age* index and sorting the retrieved tuples may be costly.
  - Clustered *E.dno* index may be better!
- **Equality queries and duplicates:**
  - Clustering on *E.hobby* helps!

```
SELECT E.dno, COUNT(*)
FROM Emp E
WHERE E.age>10
GROUP BY E.dno
```

```
SELECT E.dno
FROM Emp E
WHERE E.hobby='Stamps'
```



## Index-Only Plans

- **A number of queries can be answered without retrieving any tuples from one or more of the relations involved if a suitable index is available.**

<*E.dno*>

<*E.dno, E.eid*>

<*E.dno*>

<*E.dno, E.sal*>

B-tree trick!

<*E.age, E.sal*>

or

<*E.sal, E.age*>

```
SELECT D.mgr
FROM Dept D, Emp E
WHERE D.dno=E.dno
```

```
SELECT D.mgr, E.eid
FROM Dept D, Emp E
WHERE D.dno=E.dno
```

```
SELECT E.dno, COUNT(*)
FROM Emp E
GROUP BY E.dno
```

```
SELECT E.dno, MIN(E.sal)
FROM Emp E
GROUP BY E.dno
```

```
SELECT AVG(E.sal)
FROM Emp E
WHERE E.age=25 AND
E.sal BETWEEN 3000 AND 5000
```



## Horizontal Decompositions

- **Typical decomposition: Relation is replaced by collection of relations that are *projections*. Most important case.**
  - We will talk about this at length as part of Conceptual DB Design
- **Sometimes, might want to replace relation by a collection of relations that are *selections*.**
  - Each new relation has same schema as original, but subset of rows.
  - Collectively, new relations contain all rows of the original.
  - Typically, the new relations are disjoint.



## Horizontal Decompositions (Contd.)

- **Contracts (Cid, Sid, Jid, Did, Pid, Qty, Val)**
- **Suppose that contracts with value > 10000 are subject to different rules.**
  - So queries on Contracts will often say *WHERE val>10000*.
- **One approach: clustered B+ tree index on the val field.**
- **Second approach: replace contracts by two new relations, LargeContracts and SmallContracts, with the same attributes (CSJDPQV).**
  - Performs like index on such queries, but no index overhead.
  - Can build clustered indexes on other attributes, in addition!



## Masking Conceptual Schema Changes

```
CREATE VIEW Contracts(cid, sid, jid, did, pid, qty, val)
AS SELECT *
FROM LargeContracts
UNION
SELECT *
FROM SmallContracts
```

- **Horizontal Decomposition from above**
- **Masked by a view.**
  - NOTE: queries with condition *val>10000* must be asked wrt LargeContracts for efficiency: so some users may have to be aware of change.
    - I.e. the users who were having performance problems
    - Arguably that's OK -- they wanted a solution!



## Index Tuning "Wizards"

- **Both IBM's DB2 and MS SQL Server have automated index advisors**
  - Some info in Section 20.6 of the book
- **Basic idea:**
  - They take a workload of queries
    - Possibly based on logging what's been going on
  - They use the optimizer cost metrics to estimate the cost of the workload over different choices of sets of indexes
  - Enormous # of different choices of sets of indexes:
    - Heuristics to help this go faster



## Tuning Queries and Views

- **If a query runs slower than expected, check if an index needs to be re-clustered, or if statistics are too old.**
- **Sometimes, the DBMS may not be executing the plan you had in mind. Common areas of weakness:**
  - Selections involving **null values** (bad selectivity estimates)
  - Selections involving **arithmetic or string expressions** (ditto)
  - Selections involving **OR conditions** (ditto)
  - Complex subqueries (more on this later)
  - **Lack of evaluation features** like index-only strategies or certain join methods or poor size estimation.
- **Check the plan that is being used! Then adjust the choice of indexes or rewrite the query/view.**
  - E.g. check via POSTGRES "Explain" command
  - Some systems rewrite for you under the covers (e.g. DB2)
    - Can be confusing and/or helpful!



## More Guidelines for Query Tuning

- **Minimize the use of DISTINCT: don't need it if duplicates are acceptable, or if answer contains a key.**
- **Minimize the use of GROUP BY and HAVING:**

```
SELECT MIN (E.age)
FROM Employee E
GROUP BY E.dno
HAVING E.dno=102
```

```
SELECT MIN (E.age)
FROM Employee E
WHERE E.dno=102
```

- ❖ Consider DBMS use of index when writing arithmetic expressions: *E.age=2\*D.age* will benefit from index on *E.age*, but might not benefit from index on *D.age*!



## Guidelines for Query Tuning (Contd.)

- **Avoid using intermediate relations:**

```
SELECT E.dno, AVG(E.sal)
FROM Emp E, Dept D
WHERE E.dno=D.dno
AND D.mgrname='Joe'
GROUP BY E.dno
```

vs.

```
SELECT * INTO Temp
FROM Emp E, Dept D
WHERE E.dno=D.dno
AND D.mgrname='Joe'
```

and

```
SELECT T.dno, AVG(T.sal)
FROM Temp T
GROUP BY T.dno
```

- ❖ Does not materialize the intermediate reln Temp.
- ❖ If there is a dense B+ tree index on *<dno, sal>*, an index-only plan can be used to avoid retrieving Emp tuples in the left query!



## Points to Remember

- **Want to understand DB design (tables, indexes)?**
  - *Must* understand query optimization
- **Two parts to optimizing a query:**
  - Consider a set of alternative plans, pruning search
    - E.g., left-deep plans only
    - avoid Cartesian products.
    - Prune plans with *interesting orders* separate from unordered plans
  - Must estimate cost of each plan that is considered.
    - Output cardinality and cost for each plan node.
    - *Key issues:* Statistics, indexes, operator implementations.



## Points to Remember

- **Single-relation queries:**
  - All access paths considered, cheapest is chosen.
  - *Issues:*
    - Selections that *match* index
    - whether index key has all needed fields
    - whether index provides tuples in an interesting order.



## More Points to Remember

- **Multiple-relation queries:**
  - All single-relation plans are first enumerated.
    - Selections/projections considered as early as possible.
  - Use best 1-way plans to form 2-way plans. Prune losers.
  - Use best  $(i-1)$ -way plans and best 1-way plans to form  $i$ -way plans
  - At each level, for each subset of relations, retain:
    - best plan for each interesting order (including no order)



## Summary

- Optimization is the reason for the lasting power of the relational system
- But it is primitive in some ways
- New areas: many!
  - Smarter summary statistics (fancy histograms and "sketches")
  - Auto-tuning statistics,
  - Adaptive runtime re-optimization (e.g. *eddies*),
  - Multi-query optimization,
  - And parallel scheduling issues, etc.