

ECS 165B: Database System Implementation

Lecture 7

UC Davis
April 12, 2010

Acknowledgements: portions based on slides by Raghu Ramakrishnan and Johannes Gehrke.

Class Agenda

- Last time:
 - Dynamic aspects of B+ Trees
- Today:
 - Summary: tree-structured indices
 - Overview of query evaluation
- Reading
 - Chapter 12 in Ramakrishnan and Gehrke
 - (or Chapter 13 in Silberschatz, Korth, and Sudarshan)

Announcements

Expanded set of tests posted:

`/home/cs165b/DavisDB/TestRM.cpp`

Page file manager bugfixes:

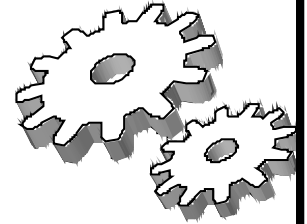
`/home/cs165b/DavisDB/PageFileManager.cpp, h`

`/home/cs165b/DavisDB/FileHandle.cpp`

Have **you** done an svn commit lately?

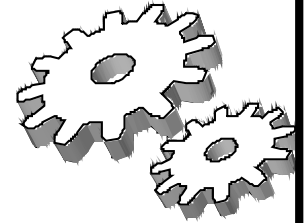
Summary: Tree-Structured Indices

Summary



- ❖ Tree-structured indexes are ideal for range-searches, also good for equality searches.
- ❖ ISAM is a static structure.
 - Only leaf pages modified; overflow pages needed.
 - Overflow chains can degrade performance unless size of data set and data distribution stay constant.
- ❖ B+ tree is a dynamic structure.
 - Inserts/deletes leave tree height-balanced; $\log_F N$ cost.
 - High fanout (F) means depth rarely more than 3 or 4.
 - Almost always better than maintaining a sorted file.

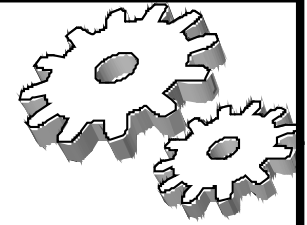
Summary (Contd.)



- Typically, 67% occupancy on average.
- Usually preferable to ISAM, modulo *locking* considerations; adjusts to growth gracefully.
- If data entries are data records, splits can change rids!
- ❖ Key compression increases fanout, reduces height.
- ❖ Bulk loading can be much faster than repeated inserts for creating a B+ tree on a large data set.
- ❖ Most widely used index in database management systems because of its versatility. One of the most optimized components of a DBMS.

Overview of Query Evaluation

Reading: Chapter 12 of Ramakrishnan and Gehrke
(Chapter 13 of Silberschatz et al)



Overview of Query Evaluation

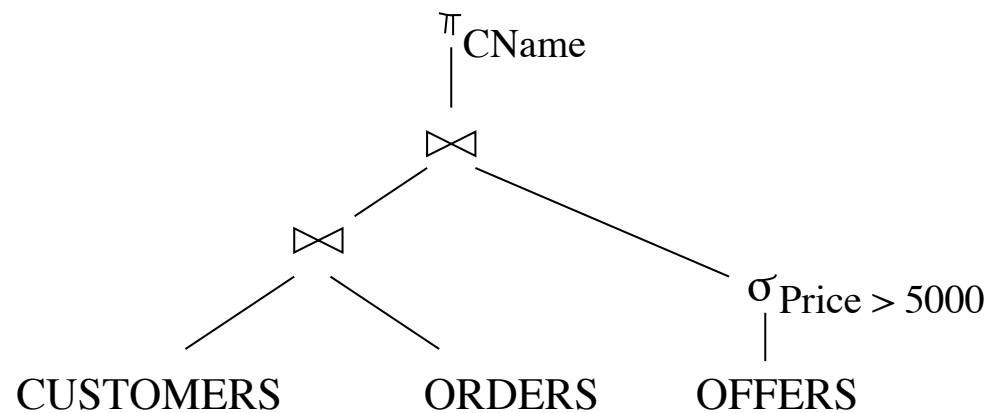
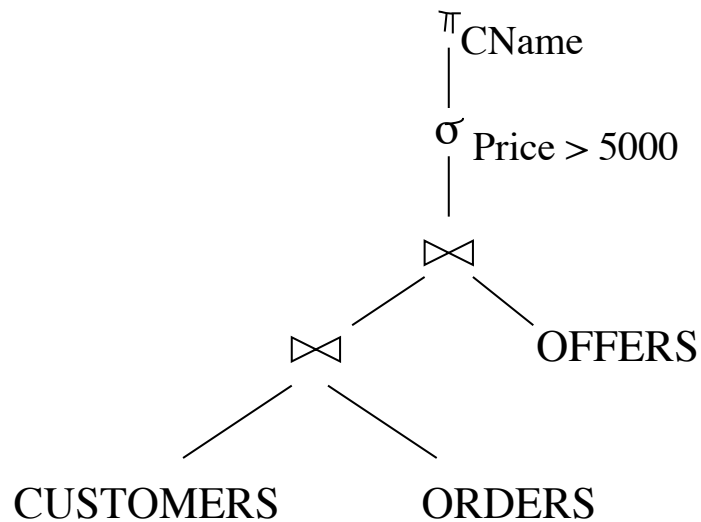
- ❖ Plan: Tree of R.A. ops, with choice of alg for each op.
 - Each operator typically implemented using a 'pull' interface: when an operator is 'pulled' for the next output tuples, it 'pulls' on its inputs and computes them.
- ❖ Two main issues in query optimization:
 - For a given query, what plans are considered?
 - Algorithm 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!
- ❖ We will study the System R approach.

Recall: "Evaluation Plan" (ECS 165A)

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

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

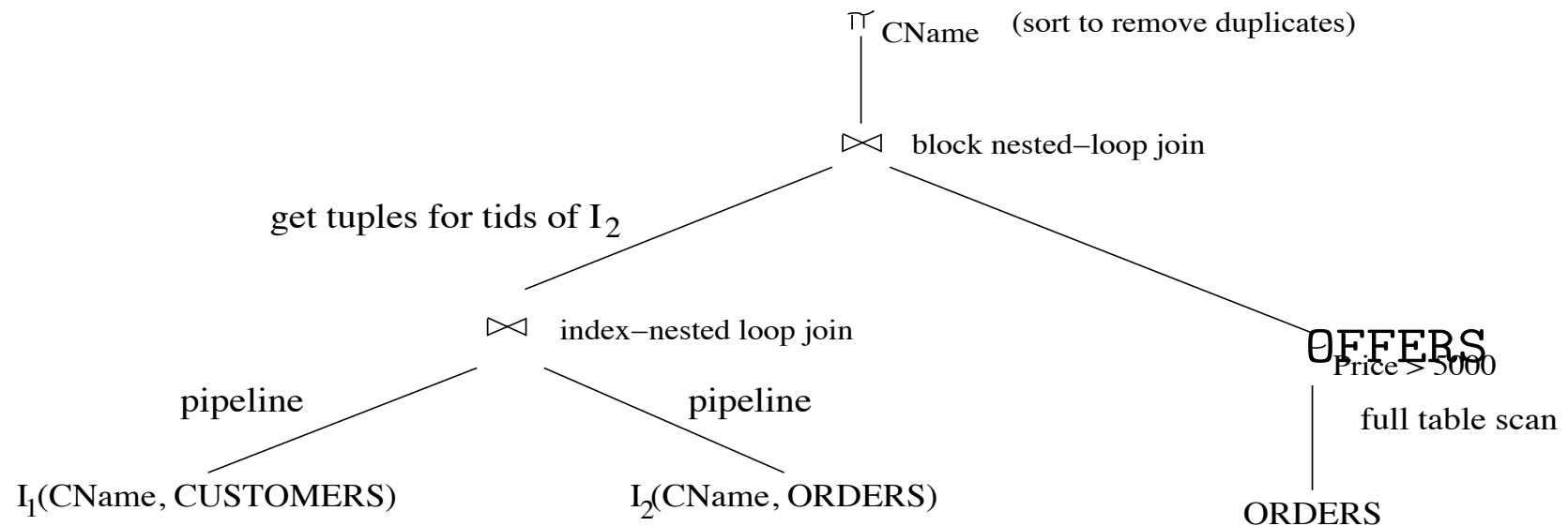
Representation as *evaluation plan* (query tree):

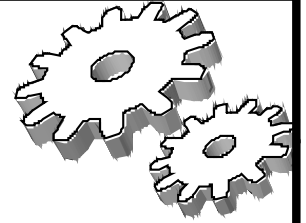


Recall: "Annotated Evaluation Plan" (ECS 165A)

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

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

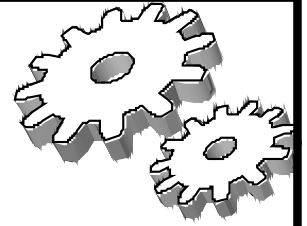




Some Common Techniques

- ❖ 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.

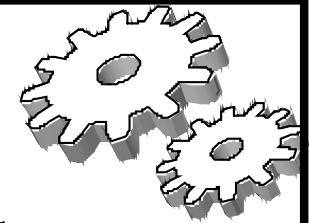
Catalogs in DavisDB

- Two catalog tables: `relations` and `attributes`

`relations` : relation name, tuple length, number of attributes, cardinality, ...

`attributes`: relation name, attribute name, offset in tuple, attribute type, attribute length, index name, ...

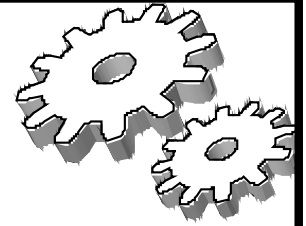
- More details when Part 3 is assigned



Access Paths

- ❖ An access path is a method of retrieving tuples:
 - File scan, or index that matches a selection (in the query)
- ❖ A tree index matches (a conjunction of) terms that involve only attributes in a *prefix* of the search key.
 - E.g., Tree index on $\langle a, b, c \rangle$ 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 $\langle a, b, c \rangle$ 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 $c=5$.

A Note on Complex Selections

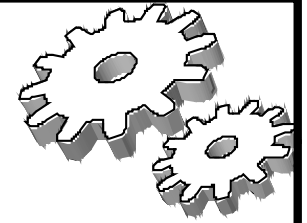


$(day < 8/9/94 \text{ AND } rname = 'Paul') \text{ OR } bid = 5 \text{ OR } sid = 3$

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

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

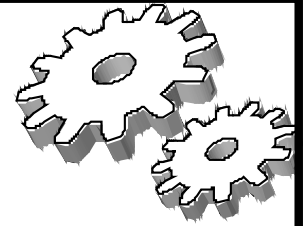
- ❖ We only discuss case with no ORs; see text if you are curious about the general case.



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 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.
 - 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

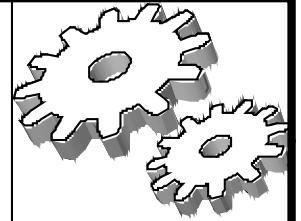


- ❖ Cost depends on #qualifying tuples, and clustering.
 - Cost of finding qualifying data entries (typically small) plus cost of retrieving records (could be large w/o clustering).
 - In example, assuming uniform distribution of names, about 10% of tuples qualify (100 pages, 10000 tuples). With a clustered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

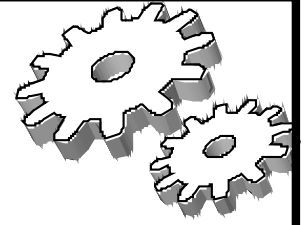
```
SELECT *  
FROM   Reserves R  
WHERE  R.name < 'C%'
```

Projection

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



- ❖ 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!

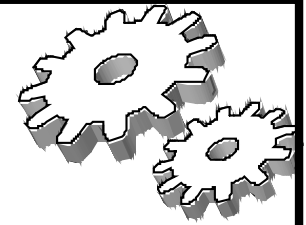


Join: Index Nested Loops

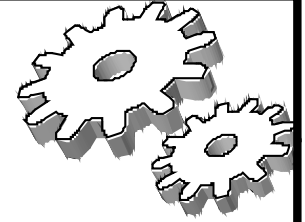
```
foreach tuple r in R do
    foreach tuple s in S where  $r_i == s_j$  do
        add  $\langle r, s \rangle$  to result
```

- ❖ 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 (typical), unclustered: upto 1 I/O per matching S tuple.

Examples of Index Nested Loops



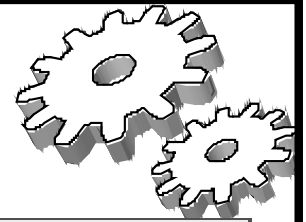
- ❖ Hash-index (Alt. 2) on *sid* of Sailors (as inner):
 - Scan Reserves: 1000 page I/Os, 100*1000 tuples.
 - For each Reserves tuple: 1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching Sailors tuple. Total: 220,000 I/Os.
- ❖ Hash-index (Alt. 2) on *sid* of Reserves (as 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.



Join: Sort-Merge ($R \sqcup_{i=j} S$)

- ❖ Sort R and S on the join column, then scan them to do a “merge” (on join col.), and output result tuples.
 - Advance scan of R until current R-tuple \geq current S tuple, then advance scan of S until current S-tuple \geq current R tuple; do this until current R tuple = current S tuple.
 - At this point, all R tuples with same value in R_i (*current R group*) and all S tuples with same value in S_j (*current S group*) match; output $\langle r, s \rangle$ for all pairs of such tuples.
 - Then resume scanning R and S.
- ❖ R is scanned once; each S group is scanned once per matching R tuple. (Multiple scans of an S group are likely to find needed pages in buffer.)

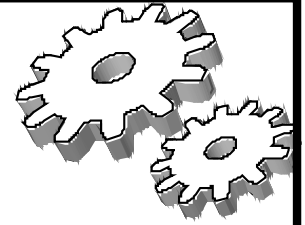
Example of Sort-Merge Join



<u>sid</u>	sname	rating	age	<u>sid</u>	<u>bid</u>	<u>day</u>	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: $M \log M + N \log N + (M+N)$
 - 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.

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.