



CS330, March 9, 2004

Indexing, Query Processing, and
Transactions

1



Some Logistics

- ❖ Next two homework assignments out today
- ❖ Extra lab session:
 - This Thursday, after class, in this room
 - Bring your laptop fully charged
- ❖ Extra homework:
 - Everybody who received ≤ 80 points on Assignment 2b, go and talk to the TAs.
- ❖ Recall: Prelim next week Thursday, 3/18!

2



Three Topics

1. Storage and Indexing
2. Query Processing
3. Transaction Management

3

Data on External Storage



- ❖ **Disks:** Can retrieve random page at fixed cost
 - But reading several consecutive pages is much cheaper than reading them in random order
- ❖ **Tapes:** Can only read pages in sequence
 - Cheaper than disks; used for archival storage
- ❖ **File organization:** Method of arranging a file of records on external storage.
 - **Record id (rid)** is sufficient to physically locate record
 - **Indexes** are data structures that allow us to find the record ids of records with given values in **index search key** fields
- ❖ **Architecture: Buffer manager** stages pages from external storage to main memory buffer pool. File and index layers make calls to the buffer manager.

4

Disks and Files



- ❖ DBMS stores information on (“hard”) disks.
- ❖ This has major implications for DBMS design!
 - **READ:** transfer data from disk to main memory (RAM).
 - **WRITE:** transfer data from RAM to disk.
 - Both are high-cost operations, relative to in-memory operations, so must be planned carefully!

5

Why Not Store Everything in Main Memory



- ❖ *Costs too much?*
- ❖ **Main memory is volatile.** We want data to be saved between runs. (Obviously!)
- ❖ Typical storage hierarchy:
 - Main memory (RAM) for currently used data.
 - Disk for the main database (secondary storage).
 - Tapes for archiving older versions of the data (tertiary storage).

6

Disks



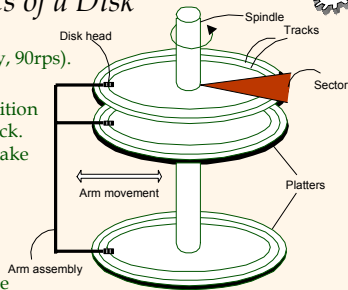
- ❖ Secondary storage device of choice.
- ❖ Main advantage over tapes: *random access* vs. *sequential*.
- ❖ Data is stored and retrieved in units called *disk blocks* or *pages*.
- ❖ Unlike RAM, time to retrieve a disk page varies depending upon location on disk.
 - Therefore, relative placement of pages on disk has major impact on DBMS performance!

7

Components of a Disk



- ❖ The platters spin (say, 90rps).
- ❖ The arm assembly is moved in or out to position a head on a desired track. Tracks under heads make a *cylinder* (imaginary!).
- ❖ Only one head reads/writes at any one time.
- ❖ *Block size* is a multiple of *sector size* (which is fixed).



8

Accessing a Disk Page



- ❖ Time to access (read/write) a disk block:
 - *seek time* (moving arms to position disk head on track)
 - *rotational delay* (waiting for block to rotate under head)
 - *transfer time* (actually moving data to/from disk surface)
- ❖ Seek time and rotational delay dominate.
 - Seek time varies from about 1 to 20msec
 - Rotational delay varies from 0 to 10msec
 - Transfer rate is about 1msec per 4KB page
- ❖ Key to lower I/O cost: *reduce seek/rotation delays!*
Hardware vs. software solutions?

9

Arranging Pages on Disk



- ❖ *'Next'* block concept:
 - blocks on same track, followed by
 - blocks on same cylinder, followed by
 - blocks on adjacent cylinder
- ❖ Blocks in a file should be arranged sequentially on disk (by *'next'*), to minimize seek and rotational delay.
- ❖ For a **sequential scan**, *pre-fetching* several pages at a time is a big win!

10

RAID



- ❖ **Disk Array:** Arrangement of several disks that gives abstraction of a single, large disk.
- ❖ **Goals:** Increase performance and reliability.
- ❖ **Two main techniques:**
 - **Data striping:** Data is partitioned; size of a partition is called the striping unit. Partitions are distributed over several disks.
 - **Redundancy:** More disks => more failures. Redundant information allows reconstruction of data if a disk fails.

11

RAID Levels



- ❖ **Level 0:** No redundancy
- ❖ **Level 1:** Mirrored (two identical copies)
 - Each disk has a mirror image (check disk)
 - Parallel reads, a write involves two disks.
 - Maximum transfer rate = transfer rate of one disk
- ❖ **Level 0+1:** Striping and Mirroring
 - Parallel reads, a write involves two disks.
 - Maximum transfer rate = aggregate bandwidth

12

RAID Levels (Contd.)



- ❖ Level 3: Bit-Interleaved Parity
 - Striping Unit: One bit. One check disk.
 - Each read and write request involves all disks; disk array can process one request at a time.
- ❖ Level 4: Block-Interleaved Parity
 - Striping Unit: One disk block. One check disk.
 - Parallel reads possible for small requests, large requests can utilize full bandwidth
 - Writes involve modified block and check disk
- ❖ Level 5: Block-Interleaved Distributed Parity
 - Similar to RAID Level 4, but parity blocks are distributed over all disks

13

Alternative File Organizations



Many alternatives exist, *each ideal for some situations, and not so good in others:*

- **Heap (random order) files:** Suitable when typical access is a file scan retrieving all records.
- **Sorted Files:** Best if records must be retrieved in some order, or only a 'range' of records is needed.
- **Indexes:** Data structures to organize records via trees or hashing.
 - Like sorted files, they speed up searches for a subset of records, based on values in certain ("search key") fields
 - Updates are much faster than in sorted files.

14

Indexes



- ❖ An **index** on a file speeds up selections on the **search key fields** for the index.
 - Any subset of the fields of a relation can be the search key for an index on the relation.
 - **Search key** is **not** the same as **key** (minimal set of fields that uniquely identify a record in a relation).
- ❖ An index contains a collection of **data entries**, and supports efficient retrieval of all data entries **k*** with a given key value **k**.

15

Alternatives for Data Entry k^* in Index

- ❖ Three alternatives:
 - Data record with key value k
 - $\langle k, \text{rid of data record with search key value } k \rangle$
 - $\langle k, \text{list of rids of data records with search key } k \rangle$
- ❖ Choice of alternative for data entries is orthogonal to the indexing technique used to locate data entries with a given key value k .
 - Examples of indexing techniques: B+ trees, hash-based structures
 - Typically, index contains auxiliary information that directs searches to the desired data entries

16

Alternatives for Data Entries (Cont'd.)

- ❖ **Alternative 1:**
 - If this is used, index structure is a file organization for data records (instead of a Heap file or sorted file).
 - At most one index on a given collection of data records can use Alternative 1. (Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency.)
 - If data records are very large, # of pages containing data entries is high. Implies size of auxiliary information in the index is also large, typically.

17

Alternatives for Data Entries (Cont'd.)

- ❖ **Alternatives 2 and 3:**
 - Data entries typically much smaller than data records. So, better than Alternative 1 with large data records, especially if search keys are small. (Portion of index structure used to direct search, which depends on size of data entries, is much smaller than with Alternative 1.)
 - Alternative 3 more compact than Alternative 2, but leads to variable sized data entries even if search keys are of fixed length.

18

Index Classification



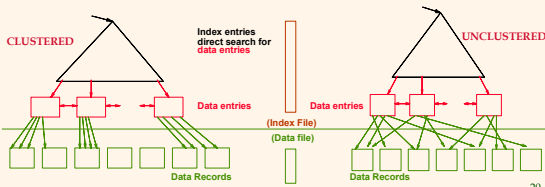
- ❖ **Primary vs. secondary:** If search key contains primary key, then called primary index.
 - **Unique** index: Search key contains a candidate key.
- ❖ **Clustered vs. unclustered:** If order of data records is the same as, or 'close to', order of data entries, then called clustered index.
 - Alternative 1 implies clustered; in practice, clustered also implies Alternative 1 (since sorted files are rare).
 - A file can be clustered on at most one search key.
 - Cost of retrieving data records through index varies *greatly* based on whether index is clustered or not!

19

Clustered vs. Unclustered Index



- ❖ Suppose that Alternative (2) is used for data entries, and that the data records are stored in a Heap file.
 - To build clustered index, first sort the Heap file (with some free space on each page for future inserts).
 - Overflow pages may be needed for inserts. (Thus, order of data recs is 'close to', but not identical to, the sort order.)



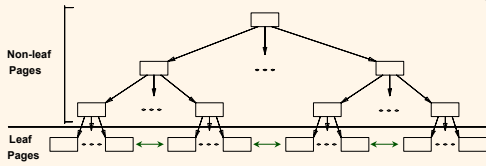
Hash-Based Indexes



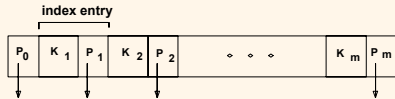
- ❖ Good for equality selections.
 - Index is a collection of **buckets**. Bucket = **primary page** plus zero or more **overflow pages**.
 - **Hashing function h**: $h(r)$ = bucket in which record r belongs. h looks at the **search key** fields of r .
- ❖ If Alternative (1) is used, the buckets contain the data records; otherwise, they contain <key, rid> or <key, rid-list> pairs.

21

B+ Tree Indexes

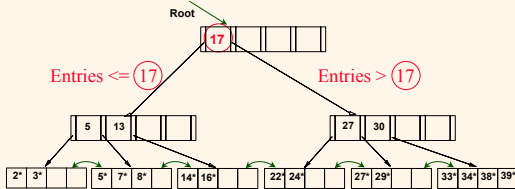


- ❖ Leaf pages contain *data entries*, and are chained (prev & next)
- ❖ Non-leaf pages contain *index entries* and direct searches:



22

Example B+ Tree



- ❖ Find 28*? 29*? All > 15* and < 30*
- ❖ Insert/delete: Find data entry in leaf, then change it. Need to adjust parent sometimes.
 - And change sometimes bubbles up the tree

23

Understanding the Workload



- ❖ For each query in the workload:
 - Which relations does it access?
 - Which attributes are retrieved?
 - Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?
- ❖ For each update in the workload:
 - Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?
 - The type of update (INSERT/DELETE/UPDATE), and the attributes that are affected.

24

Choice of Indexes



- ❖ What indexes should we create?
 - Which relations should have indexes? What field(s) should be the search key? Should we build several indexes?
- ❖ For each index, what kind of an index should it be?
 - Clustered? Hash/tree?

25

Choice of Indexes (Contd.)



- ❖ **One approach:** Consider the most important queries in turn. Consider the best plan using the current indexes, and see if a better plan is possible with an additional index. If so, create it.
 - Obviously, this implies that we must understand how a DBMS evaluates queries and creates **query evaluation plans!**
 - For now, we discuss simple 1-table queries.
- ❖ 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.

26

Index Selection Guidelines



- ❖ Attributes in WHERE clause are candidates for index keys.
 - Exact match condition suggests hash index.
 - Range query suggests tree index.
 - Clustering is especially useful for range queries; can also help on equality queries if there are many duplicates.
- ❖ Multi-attribute search keys should be considered when a WHERE clause contains several conditions.
 - Order of attributes is important for range queries.
 - Such indexes can sometimes enable **index-only** strategies for important queries.
 - For index-only strategies, clustering is not important!
- ❖ Try to choose indexes that benefit as many queries as possible. Since only one index can be clustered per relation, choose it based on important queries that would benefit the most from clustering.

27

Examples of Clustered Indexes



❖ 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
FROM Emp E
WHERE E.age>40
```

```
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
```

28

Indexes with Composite Search Keys



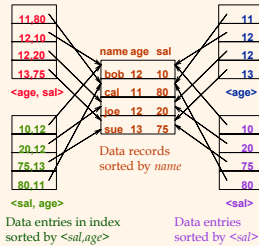
❖ **Composite Search Keys:** Search on a combination of fields.

- **Equality query:** Every field value is equal to a constant value. E.g. wrt <sal,age> index:
 - age=20 and sal =75
- **Range query:** Some field value is not a constant. E.g.:
 - age =20; or age=20 and sal > 10

❖ Data entries in index sorted by search key to support range queries.

- **Lexicographic order,** or
- **Spatial order.**

Examples of composite key indexes using lexicographic order.



29

Composite Search Keys



❖ To retrieve Emp records with age=30 AND sal=4000, an index on <age,sal> would be better than an index on age or an index on sal.

- Choice of index key orthogonal to clustering etc.

❖ If condition is: 20<age<30 AND 3000<sal<5000:

- Clustered tree index on <age,sal> or <sal,age> is best.

❖ If condition is: age=30 AND 3000<sal<5000:

- Clustered <age,sal> index much better than <sal,age> index!

❖ Composite indexes are larger, updated more often.

30

Summary



- ❖ Many alternative file organizations exist, each appropriate in some situation.
- ❖ If selection queries are frequent, sorting the file or building an *index* is important.
 - Hash-based indexes only good for equality search.
 - Sorted files and tree-based indexes best for range search; also good for equality search. (Files rarely kept sorted in practice; B+ tree index is better.)
- ❖ Index is a collection of data entries plus a way to quickly find entries with given key values.

31

Summary (Contd.)



- ❖ Data entries can be actual data records, <key, rid> pairs, or <key, rid-list> pairs.
 - Choice orthogonal to *indexing technique* used to locate data entries with a given key value.
- ❖ Can have several indexes on a given file of data records, each with a different search key.
- ❖ Indexes can be classified as clustered vs. unclustered, primary vs. secondary, and dense vs. sparse. Differences have important consequences for utility/performance.

32

Summary (Contd.)



- ❖ Understanding the nature of the *workload* for the application, and the performance goals, is essential to developing a good design.
 - What are the important queries and updates? What attributes/relations are involved?
- ❖ Indexes must be chosen to speed up important queries (and perhaps some updates!).
 - Index maintenance overhead on updates to key fields.
 - Choose indexes that can help many queries, if possible.
 - Build indexes to support index-only strategies.
 - Clustering is an important decision; only one index on a given relation can be clustered!
 - Order of fields in composite index key can be important.

33

Three Topics



1. Storage and Indexing
2. Query Processing
3. Transaction Management

34

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.

35

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!

36

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.

37

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 $c=5$.

38

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.

39

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.

40

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%'
```

41

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!

42

Join: Index Nested Loops



```
foreach tuple r in R do
  foreach tuple s in S where ri == sj do
    add <r, s> 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: $M + (M * p_s) * \text{cost of finding matching S tuples}$
 - $M = \# \text{pages of R}$, $p_r = \# \text{R tuples per page}$
- ❖ 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.

43

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.

44

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



- ❖ Sort R and S on the join column, then scan them to do a "merge" (on join col.), and output result tuples.
 - 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.)

45

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:** $M \log M + N \log N + (M+N)$
 - 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.

46

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.

47

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.

48

Size Estimation and Reduction Factors



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

- ❖ Consider a query block:
- ❖ 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 $(High(I)-value)/(High(I)-Low(I))$

49

Schema for Examples



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

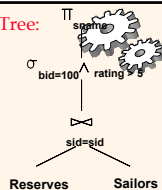
- ❖ Similar to old schema; *rname* added for variations.
- ❖ Reserves:
 - Each tuple is 40 bytes long, 100 tuples per page, 1000 pages.
- ❖ Sailors:
 - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.

50

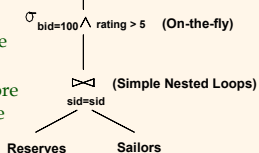
Motivating Example

RA Tree:

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

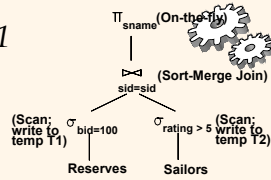


- ❖ Cost: $500+500*1000 I/Os$
- ❖ By no means the worst plan! **Plan:** π_{sname} (On-the-fly)
- ❖ Misses several opportunities: selections could have been 'pushed' earlier, no use is made of any available indexes, etc.
- ❖ Goal of optimization: To find more efficient plans that compute the same answer.



51

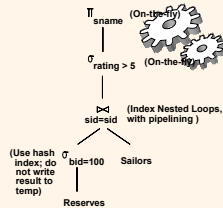
Alternative Plans 1 (No Indexes)



- ❖ **Main difference: push selects.**
- ❖ **With 5 buffers, cost of plan:**
 - Scan Reserves (1000) + write temp T1 (10 pages, if we have 100 boats, uniform distribution).
 - Scan Sailors (500) + write temp T2 (250 pages, if we have 10 ratings).
 - Sort T1 ($2 \times 2 \times 10$), sort T2 ($2 \times 3 \times 250$), merge (10+250)
 - **Total: 3560 page I/Os.**
- ❖ **If we used BNL join, join cost = $10 + 4 \times 250$, total cost = 2770.**
- ❖ **If we 'push' projections, T1 has only *sid*, T2 only *sid* and *sname*:**
 - T1 fits in 3 pages, cost of BNL drops to under 250 pages, **total < 2000.**

52

Alternative Plans 2 With Indexes



- ❖ 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.**

53

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.

54

Three Topics



1. Storage and Indexing
2. Query Processing
3. Transaction Management

55

Transactions



- ❖ Concurrent execution of user programs is essential for good DBMS performance.
 - Because disk accesses are frequent, and relatively slow, it is important to keep the cpu humming by working on several user programs concurrently.
- ❖ A user's program may carry out many operations on the data retrieved from the database, but the DBMS is only concerned about what data is read/written from/to the database.
- ❖ A transaction is the DBMS's abstract view of a user program: a sequence of reads and writes.

56

Concurrency in a DBMS



- ❖ Users submit transactions, and can think of each transaction as executing by itself.
 - Concurrency is achieved by the DBMS, which interleaves actions (reads/writes of DB objects) of various transactions.
 - Each transaction must leave the database in a consistent state if the DB is consistent when the transaction begins.
 - DBMS will enforce some ICs, depending on the ICs declared in CREATE TABLE statements.
 - Beyond this, the DBMS does not really understand the semantics of the data. (e.g., it does not understand how the interest on a bank account is computed).
- ❖ Issues: Effect of *interleaving* transactions, and *crashes*.

57

Atomicity of Transactions



- ❖ A transaction might *commit* after completing all its actions, or it could *abort* (or be aborted by the DBMS) after executing some actions.
- ❖ A very important property guaranteed by the DBMS for all transactions is that they are *atomic*. That is, a user can think of a Xact as always executing all its actions in one step, or not executing any actions at all.
 - DBMS *logs* all actions so that it can *undo* the actions of aborted transactions.

58

Example



- ❖ Consider two transactions (*Xacts*):

```
T1: BEGIN A=A+100, B=B-100 END
T2: BEGIN A=1.06*A, B=1.06*B END
```

- ❖ Intuitively, the first transaction is transferring \$100 from B's account to A's account. The second is crediting both accounts with a 6% interest payment.
- ❖ There is no guarantee that T1 will execute before T2 or vice-versa, if both are submitted together. However, the net effect *must* be equivalent to these two transactions running serially in some order.

59

Example (Contd.)



- ❖ Consider a possible interleaving (*schedule*):

```
T1: A=A+100, B=B-100
T2: A=1.06*A, B=1.06*B
```

- ❖ This is OK. But what about:

```
T1: A=A+100, B=B-100
T2: A=1.06*A, B=1.06*B
```

- ❖ The DBMS's view of the second schedule:

```
T1: R(A), W(A), R(B), W(B)
T2: R(A), W(A), R(B), W(B)
```

60

Scheduling Transactions



- ❖ **Serial schedule:** Schedule that does not interleave the actions of different transactions.
 - ❖ **Equivalent schedules:** For any database state, the effect (on the set of objects in the database) of executing the first schedule is identical to the effect of executing the second schedule.
 - ❖ **Serializable schedule:** A schedule that is equivalent to some serial execution of the transactions.
- (Note: If each transaction preserves consistency, every serializable schedule preserves consistency.)

61

Anomalies with Interleaved Execution



- ❖ Reading Uncommitted Data (WR Conflicts, “dirty reads”):

T1:	R(A), W(A),	R(B), W(B), Abort
T2:	R(A), W(A), C	

- ❖ Unrepeatable Reads (RW Conflicts):

T1:	R(A),	R(A), W(A), C
T2:	R(A), W(A), C	

62

Anomalies (Continued)



- ❖ Overwriting Uncommitted Data (WW Conflicts):

T1:	W(A),	W(B), C
T2:	W(A), W(B), C	

63

Lock-Based Concurrency Control



- ❖ **Strict Two-phase Locking (Strict 2PL) Protocol:**
 - Each Xact must obtain a **S (shared)** lock on object before reading, and an **X (exclusive)** lock on object before writing.
 - All locks held by a transaction are released when the transaction completes
 - **(Non-strict) 2PL Variant:** Release locks anytime, but cannot acquire locks after releasing any lock.
 - If an Xact holds an X lock on an object, no other Xact can get a lock (S or X) on that object.
- ❖ Strict 2PL allows only serializable schedules.
 - Additionally, it simplifies transaction aborts
 - **(Non-strict) 2PL** also allows only serializable schedules, but involves more complex abort processing

64

Aborting a Transaction



- ❖ If a transaction T_i is aborted, all its actions have to be undone. Not only that, if T_j reads an object last written by T_i , T_j must be aborted as well!
- ❖ Most systems avoid such **cascading aborts** by releasing a transaction's locks only at commit time.
 - If T_i writes an object, T_j can read this only after T_i commits.
- ❖ In order to *undo* the actions of an aborted transaction, the DBMS maintains a *log* in which every write is recorded. This mechanism is also used to recover from system crashes: all active Xacts at the time of the crash are aborted when the system comes back up.

65

The Log



- ❖ The following actions are recorded in the log:
 - **T_i writes an object:** the old value and the new value.
 - Log record must go to disk **before** the changed page!
 - **T_i commits/aborts:** a log record indicating this action.
- ❖ Log records are chained together by Xact id, so it's easy to undo a specific Xact.
- ❖ Log is often *duplexed* and *archived* on stable storage.
- ❖ All log related activities (and in fact, all CC related activities such as lock/unlock, dealing with deadlocks etc.) are handled transparently by the DBMS.

66

Recovering From a Crash



- ❖ There are 3 phases in the *Aries* recovery algorithm:
 - **Analysis:** Scan the log forward (from the most recent *checkpoint*) to identify all Xacts that were active, and all dirty pages in the buffer pool at the time of the crash.
 - **Redo:** Redoes all updates to dirty pages in the buffer pool, as needed, to ensure that all logged updates are in fact carried out and written to disk.
 - **Undo:** The writes of all Xacts that were active at the crash are undone (by restoring the *before value* of the update, which is in the log record for the update), working backwards in the log. (Some care must be taken to handle the case of a crash occurring during the recovery process!)

67

Summary



- ❖ Concurrency control and recovery are among the most important functions provided by a DBMS.
- ❖ Users need not worry about concurrency.
 - System automatically inserts lock/unlock requests and schedules actions of different Xacts in such a way as to ensure that the resulting execution is equivalent to executing the Xacts one after the other in some order.
- ❖ Write-ahead logging (WAL) is used to undo the actions of aborted transactions and to restore the system to a consistent state after a crash.
 - *Consistent state:* Only the effects of committed Xacts seen.

68
