

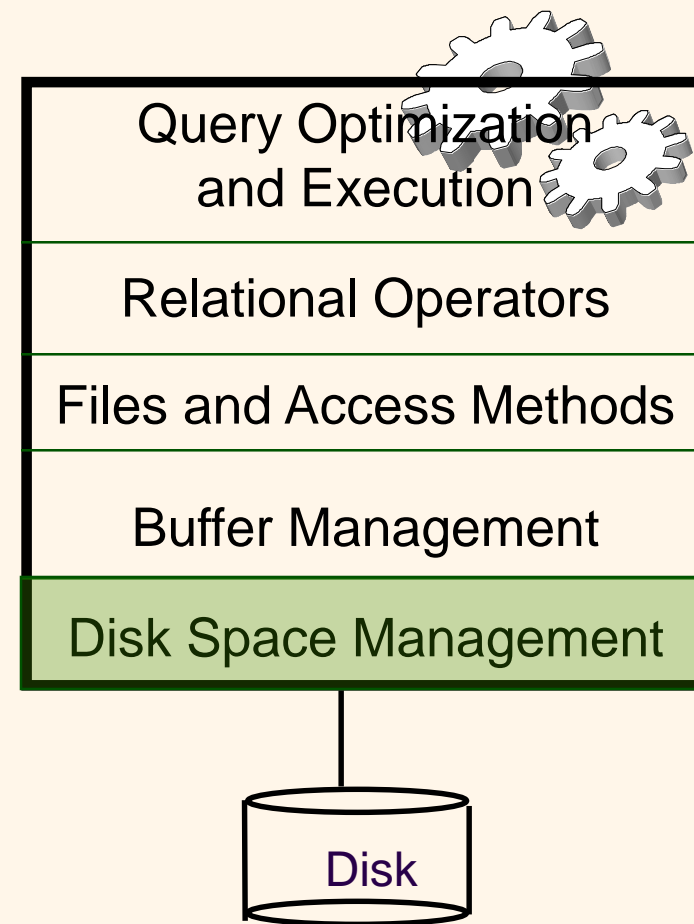
Overview of Storage and Indexing

Storing Data: Disks and Files

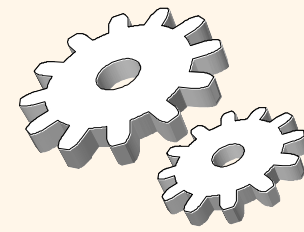
Chapters 8-9

Disk Space Management

- ❖ Lowest layer of DBMS software manages space on disk.
- ❖ Higher levels call upon this layer to:
 - allocate/de-allocate a page
 - read/write a page



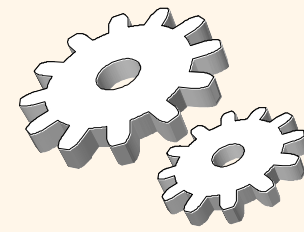
- ❖ Request for a *sequence* of pages must be satisfied by allocating the pages sequentially on disk! Higher levels don't need to know how this is done, or how free space is managed.



Alternative File Organizations

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

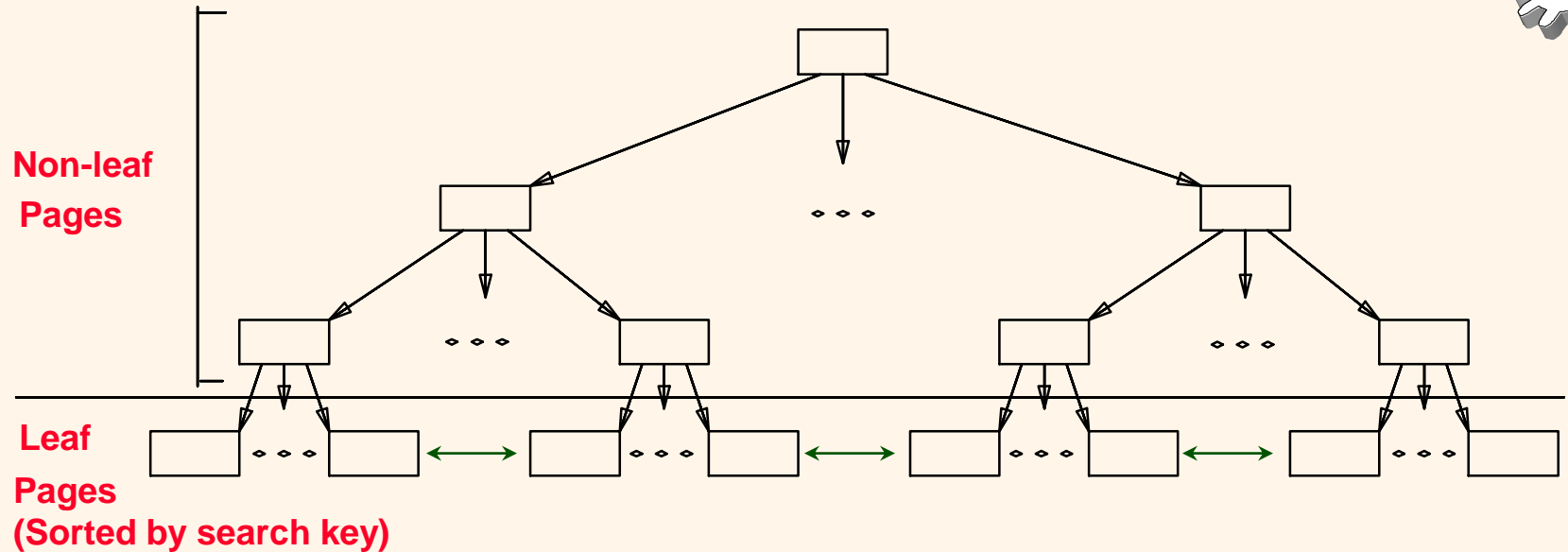
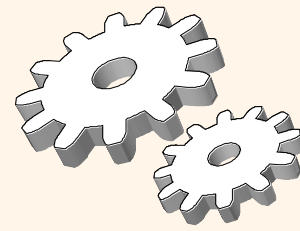
- Heap (random order) files: Suitable when typical access is a file scan retrieving all records.
- Sorted Files: Best if records must be retrieved in some order, or only a `range` of records is needed.
- Indexes: Data structures to organize records via trees or hashing.
 - Like sorted files, they speed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



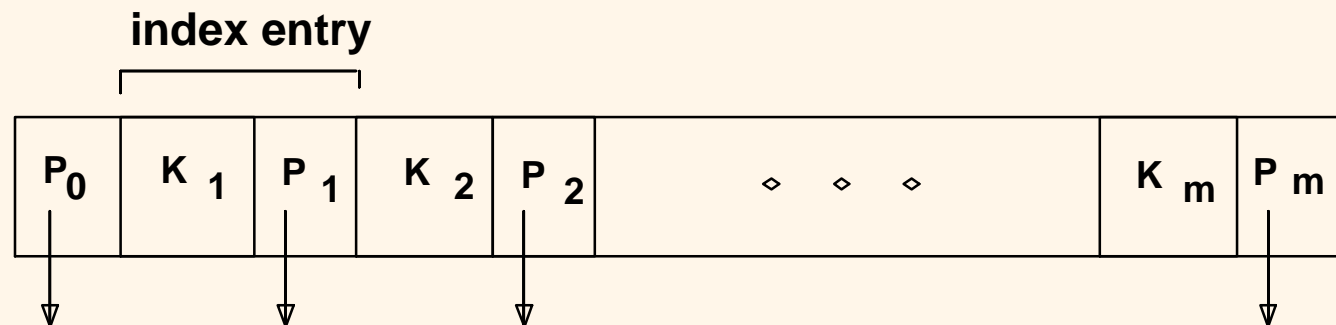
Indexes

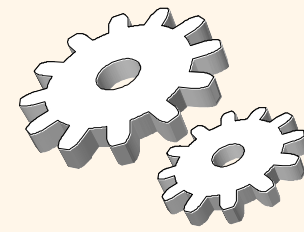
- ❖ An index on a file speeds up selections on the *search key fields* for the index.
 - Any subset of the fields of a relation can be the search key for an index on the relation.
 - *Search key* is **not** the same as *key* (minimal set of fields that uniquely identify a record in a relation).
- ❖ An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

B+ Tree Indexes

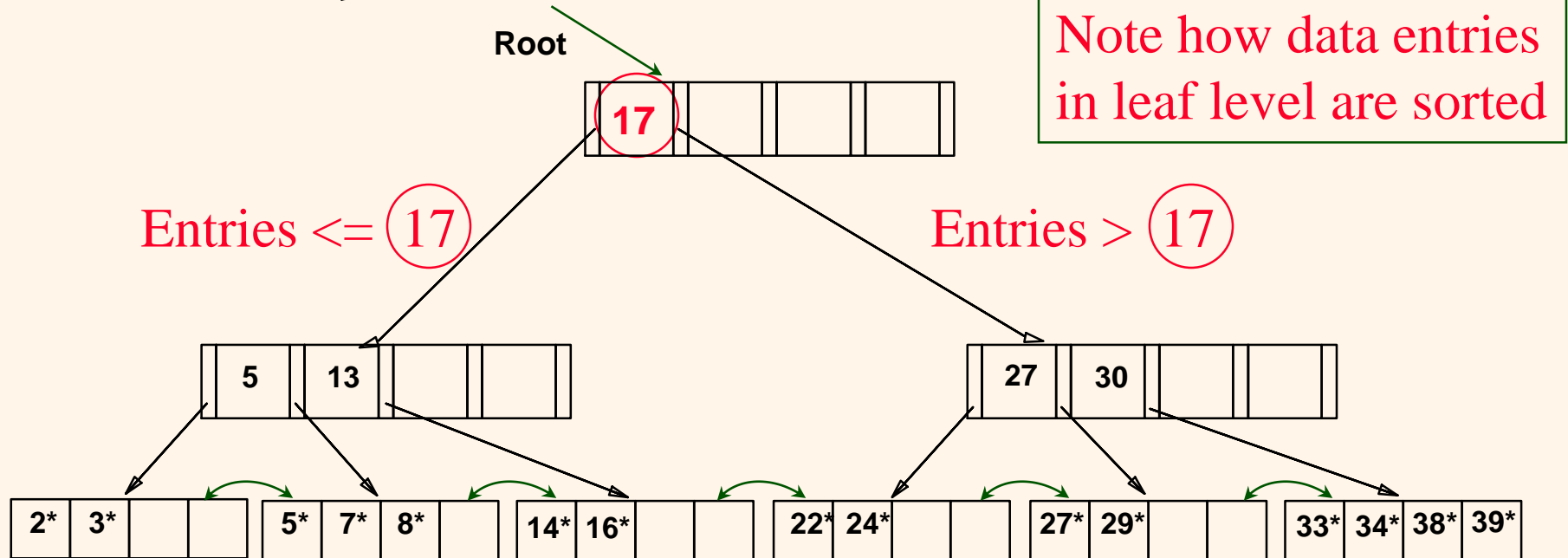


- ❖ Leaf pages contain *data entries*, and are chained (prev & next)
- ❖ Non-leaf pages have *index entries*; only used to direct searches:

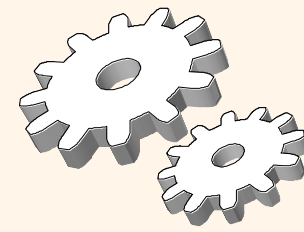




Example B+ Tree

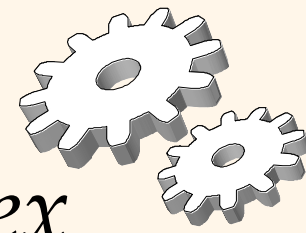


- ❖ Find 28*? 29*? All $> 15^*$ and $< 30^*$
- ❖ Insert/delete: Find data entry in leaf, then change it. Need to adjust parent sometimes.
 - And change sometimes bubbles up the tree



Hash-Based Indexes

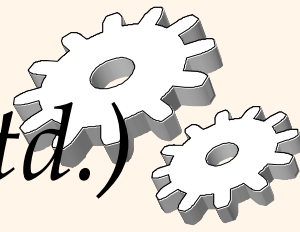
- ❖ Good for equality selections.
- ❖ Index is a collection of buckets.
 - Bucket = *primary page* plus zero or more *overflow pages*.
 - Buckets contain data entries.
- ❖ *Hashing function h* : $h(r)$ = bucket in which (data entry for) record r belongs. h looks at the *search key* fields of r .
 - *No need for “index entries” in this scheme.*



Alternatives for Data Entry k^ in Index*

- ❖ In a data entry k^* we can store:
 - Data record with key value k , or
 - $\langle k, \text{rid of data record with search key value } k \rangle$, or
 - $\langle k, \text{list of rids of data records with search key } k \rangle$
- ❖ 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

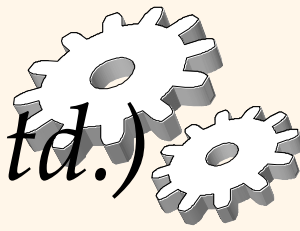
Alternatives for Data Entries (Contd.)



❖ **Alternative 1:**

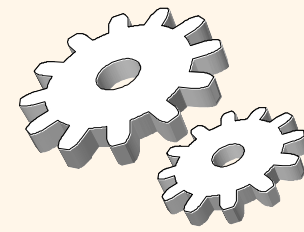
- If this is used, index structure is a file organization for data records (instead of a Heap file or sorted file).
- At most one index on a given collection of data records can use Alternative 1. (Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency.)
- If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

Alternatives for Data Entries (Contd.)



❖ 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.

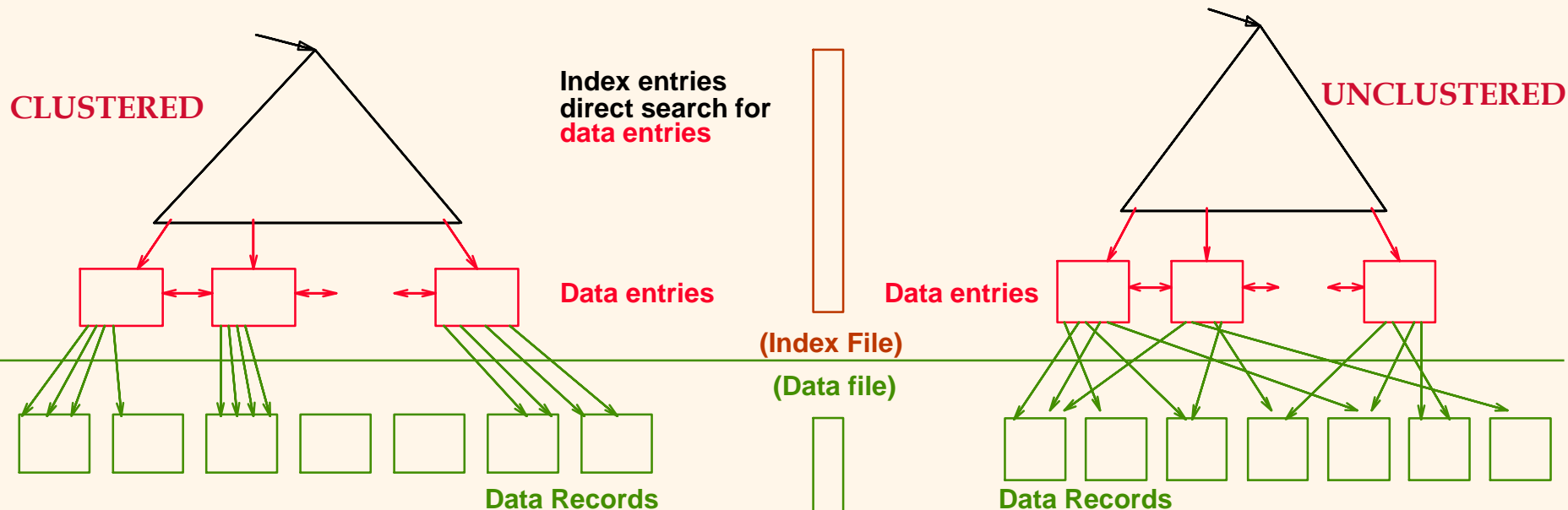


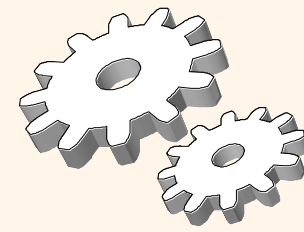
Index Classification

- ❖ *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!

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.)



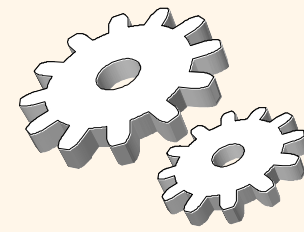


Cost Model for Our Analysis

We ignore CPU costs, for simplicity:

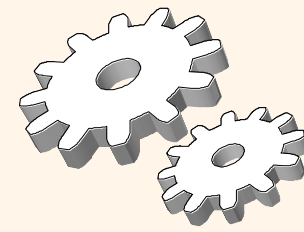
- **B:** The number of data pages
- **R:** Number of records per page
- **D:** (Average) time to read or write disk page
- Measuring number of page I/O's ignores gains of pre-fetching a sequence of pages; thus, even I/O cost is only approximated.
- Average-case analysis; based on several simplistic assumptions.
- The size of the data entry is 10% of the corresponding record

** Good enough to show the overall trends!*



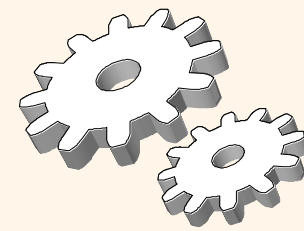
Comparing File Organizations

- ❖ **Heap files** (random order; insert at eof)
- ❖ **Sorted files**, sorted on $\langle \text{age}, \text{sal} \rangle$
- ❖ **Clustered B+ tree file**, Alternative (1), search key $\langle \text{age}, \text{sal} \rangle$
- ❖ **Heap file with unclustered B + tree index** on search key $\langle \text{age}, \text{sal} \rangle$
- ❖ **Heap file with unclustered hash index** on search key $\langle \text{age}, \text{sal} \rangle$



Operations to Compare

- ❖ Scan: Fetch all records from disk
- ❖ Equality search
- ❖ Range selection
- ❖ Insert a record
- ❖ Delete a record



Assumptions in Our Analysis

❖ Heap Files:

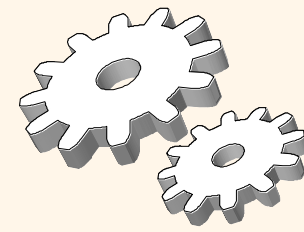
- Equality selection on key; exactly one match.

❖ Sorted Files:

- Files compacted after deletions.

❖ Indexes:

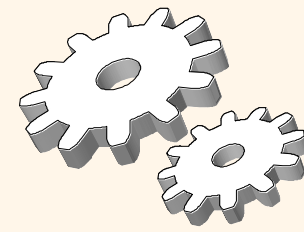
- Alt (2), (3): data entry size = 10% size of record
- Hash: No overflow buckets.
 - 80% page occupancy => File size = 1.25 data size
- Tree: 67% occupancy (this is typical).
 - Implies file size = 1.5 data size



Cost of Operations

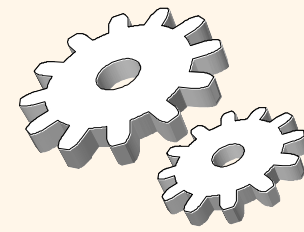
	(a) Scan	(b) Equality	(c) Range	(d) Insert	(e) Delete
(1) Heap	BD	0.5BD	BD	2D	Search +D
(2) Sorted	BD	$D \log_2 B$	$D(\log_2 B + \# \text{ pgs with match recs})$	Search + BD	Search +BD
(3) Clustered	1.5BD	$D \log_F 1.5B$	$D(\log_F 1.5B + \# \text{ pgs w. match recs})$	Search + D	Search +D
(4) Unclust. Tree index	$BD(R+0.15)$	$D(1 + \log_F 0.15B)$	$D(\log_F 0.15B + \# \text{ match recs})$	Search + 2D	Search + 2D
(5) Unclust. Hash index	$BD(R+0.125)$	2D	BD	Search + 2D	Search + 2D

** Several assumptions underlie these (rough) estimates!*



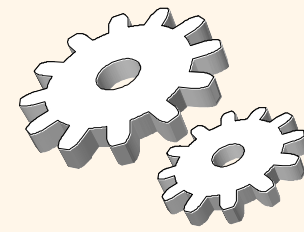
Understanding the Workload

- ❖ 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.



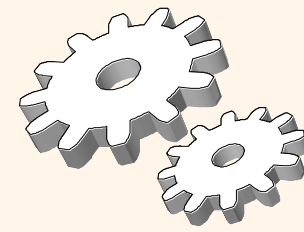
Choice of Indexes

- ❖ 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?



Choice of Indexes (Contd.)

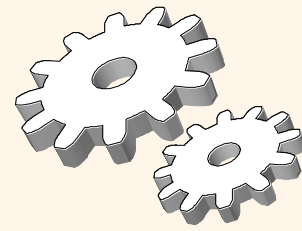
- ❖ **One approach:** Consider the most important queries in turn. Consider the best plan using the current indexes, and see if a better plan is possible with an additional index. If so, create it.
 - Obviously, this implies that we must understand how a DBMS evaluates queries and creates **query evaluation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



Index Selection Guidelines

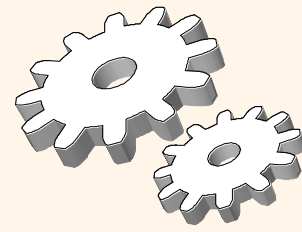
- ❖ Attributes in WHERE clause are candidates for index keys.
 - Exact match condition suggests hash index.
 - Range query suggests tree index.
 - Clustering is especially useful for range queries; can also help on equality queries if there are many duplicates.

Index Selection Guidelines



- ❖ 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.

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?

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

❖ 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!

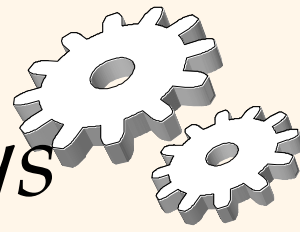
```
SELECT E.dno, COUNT (*)  
FROM Emp E  
WHERE E.age>10  
GROUP BY E.dno
```

❖ Equality queries and duplicates:

- Clustering on *E.hobby* helps!

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

Indexes with Composite Search Keys

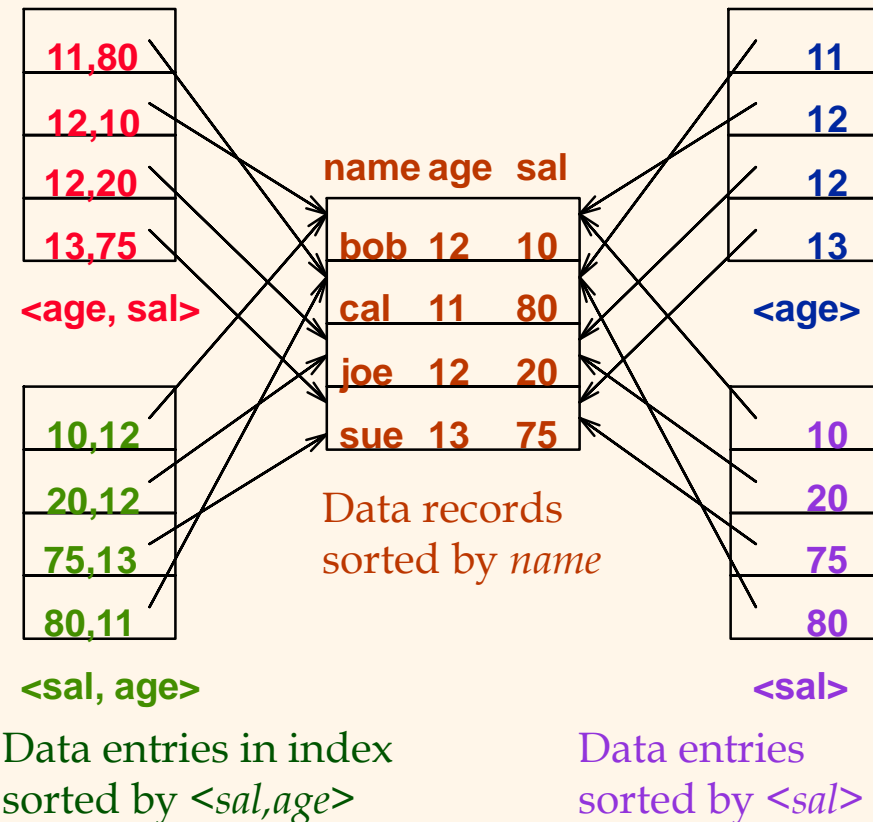


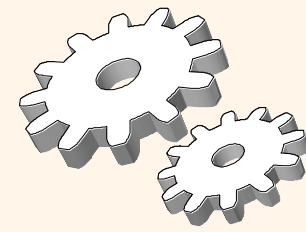
❖ **Composite Search Keys:** Search on a combination of fields.

- **Equality query:** Every field value is equal to a constant value. E.g. wrt $\langle \text{sal}, \text{age} \rangle$ 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.

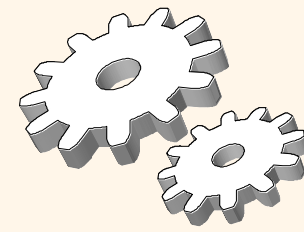
Examples of composite key indexes using lexicographic order.





Composite Search Keys

- ❖ To retrieve Emp records with $age=30$ AND $sal=4000$, an index on $\langle age, sal \rangle$ 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 $\langle age, sal \rangle$ or $\langle sal, age \rangle$ is best.
- ❖ If condition is: $age=30$ AND $3000 < sal < 5000$:
 - Clustered $\langle age, sal \rangle$ tree index much better than $\langle sal, age \rangle$ index!
- ❖ Composite indexes are larger, updated more often.



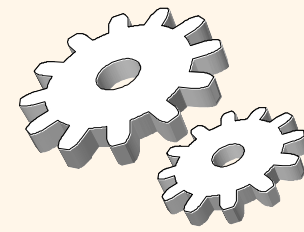
Index-Only Plans

- ❖ 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.sal>*
Tree index!
- <E. age, E.sal>*
or
<E.sal, E.age>
Tree index!

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

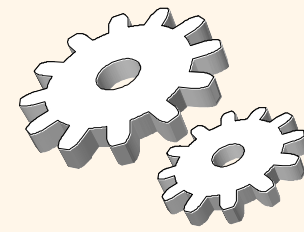
```
SELECT E.dno, MIN(E.sal)  
FROM Emp E  
GROUP BY E.dno
```

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



Summary

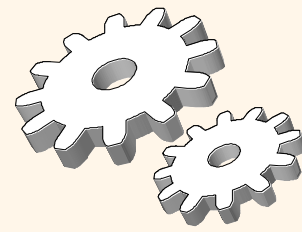
- ❖ Many alternative file organizations exist, each appropriate in some situation.
- ❖ If selection queries are frequent, sorting the file or building an *index* is important.
 - Hash-based indexes only good for equality search.
 - Sorted files and tree-based indexes best for range search; also good 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.



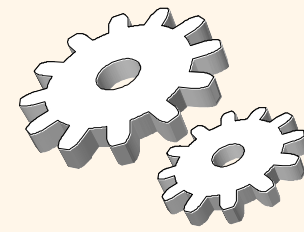
Summary (Contd.)

- ❖ Data entries can be actual data records, $\langle \text{key}, \text{rid} \rangle$ pairs, or $\langle \text{key}, \text{rid-list} \rangle$ 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, and primary vs. secondary. Differences have important consequences for utility/performance.

Summary (Contd.)

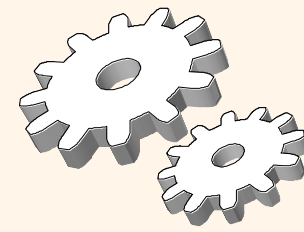


- ❖ Understanding the nature of the *workload* for the application, and the performance goals, is essential to developing a good design.
 - What are the important queries and updates? What attributes/relations are involved?
- ❖ Indexes must be chosen to speed up important queries (and perhaps some updates!).
 - 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.



Summary

- ❖ Disks provide cheap, non-volatile storage.
 - Random access, but cost depends on location of page on disk; important to arrange data sequentially to minimize *seek* and *rotation* delays.
- ❖ Buffer manager brings pages into RAM.
 - Page stays in RAM until released by requestor.
 - Written to disk when frame chosen for replacement (which is sometime after requestor releases the page).
 - Choice of frame to replace based on *replacement policy*.
 - Tries to *pre-fetch* several pages at a time.



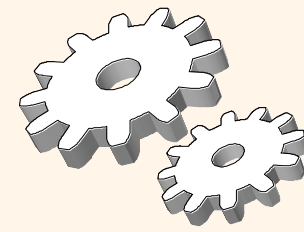
Summary (Contd.)

❖ DBMS vs. OS File Support

- DBMS needs features not found in many OS's, e.g., forcing a page to disk, controlling the order of page writes to disk, files spanning disks, ability to control pre-fetching and page replacement policy based on predictable access patterns, etc.

❖ Variable length record format with field offset directory offers support for direct access to i'th field and null values.

❖ Slotted page format supports variable length records and allows records to move on page.



Summary (Contd.)

- ❖ File layer keeps track of pages in a file, and supports abstraction of a collection of records.
 - Pages with free space identified using linked list or directory structure (similar to how pages in file are kept track of).
- ❖ Indexes support efficient retrieval of records based on the values in some fields.
- ❖ Catalog relations store information about relations, indexes and views. (*Information that is common to all records in a given collection.*)