

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.

System Catalogs



- ❖ For each index:
 - structure (e.g., B+ tree) and search key fields
- ❖ For each relation:
 - name, file name, file structure (e.g., Heap file)
 - attribute name and type, for each attribute
 - index name, for each index
 - integrity constraints
- ❖ For each view:
 - view name and definition
- ❖ Plus statistics, authorization, buffer pool size, etc.
 - * *Catalogs are themselves stored as relations!*

Attr_Cat(attr_name, rel_name, type, position)



attr_name	rel_name	type	position
attr_name	Attribute_Cat	string	1
rel_name	Attribute_Cat	string	2
type	Attribute_Cat	string	3
position	Attribute_Cat	integer	4
sid	Students	string	1
name	Students	string	2
login	Students	string	3
age	Students	integer	4
gpa	Students	real	5
fid	Faculty	string	1
fname	Faculty	string	2
sal	Faculty	real	3

Alternative File Organizations



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

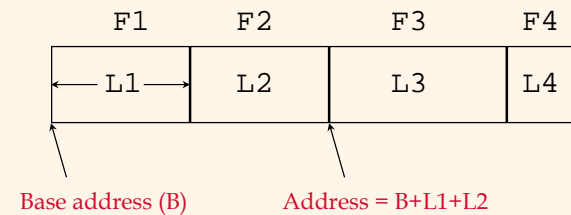
- **Heap (random order) files:** Suitable when typical access is a file scan retrieving all records.
- **Sorted Files:** Best if records must be retrieved in some order, or only a 'range' of records is needed.
- **Indexes:** Data structures to organize records via trees or hashing.
 - Like sorted files, 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.

Files of Records



- ❖ Page or block is OK when doing I/O, but higher levels of DBMS operate on *records*, and *files of records*.
- ❖ **FILE**: A collection of pages, each containing a collection of records. Must support:
 - insert/delete/modify record
 - read a particular record (specified using *record id*)
 - scan all records (possibly with some conditions on the records to be retrieved)

Record Formats: Fixed Length

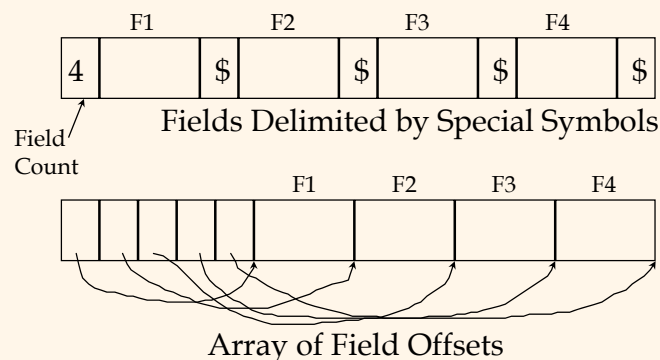


- ❖ Information about field types same for all records in a file; stored in *system catalogs*.
- ❖ Finding *i*'th field does not require scan of record.

Record Formats: Variable Length

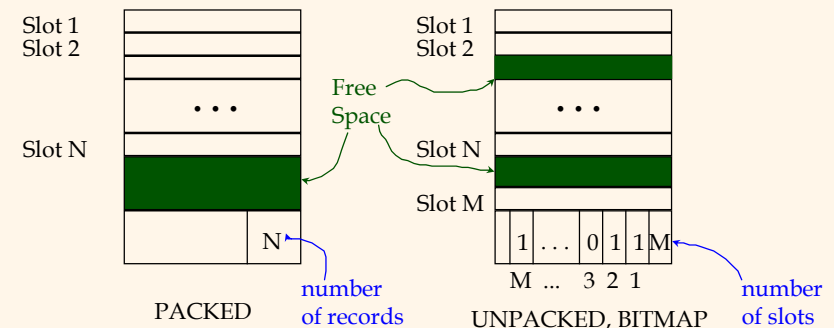


- ❖ Two alternative formats (# fields is fixed):



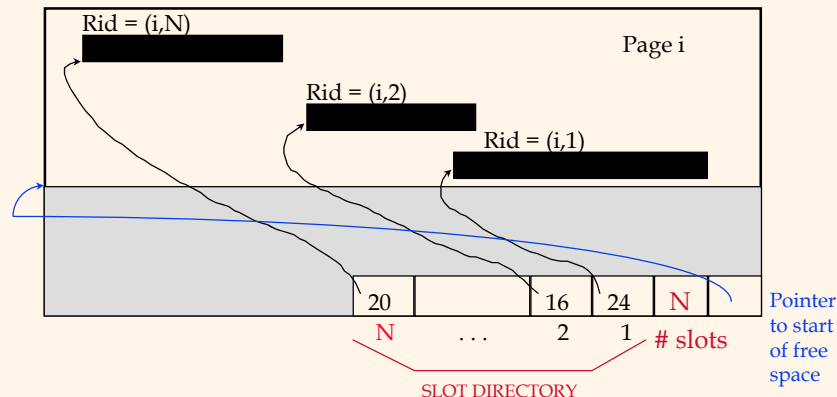
- * Second offers direct access to *i*'th field, efficient storage of *nulls* (special *don't know* value); small directory overhead.

Page Formats: Fixed Length Records



- * **Record id** = *<page id, slot #>*. In first alternative, moving records for free space management changes *rid*; may not be acceptable.

Page Formats: Variable Length Records

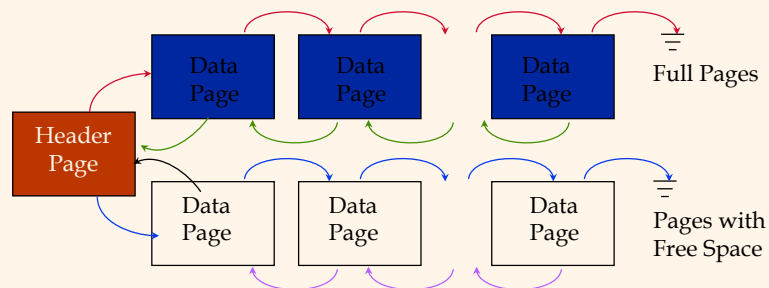


- * Can move records on page without changing rid; so, attractive for fixed-length records too.

Unordered (Heap) Files

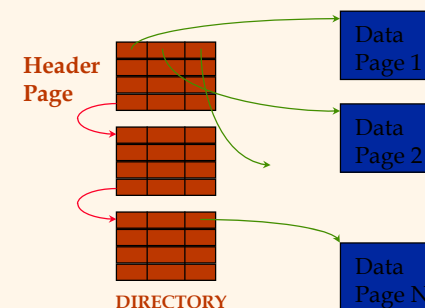
- ❖ Simplest file structure contains records in no particular order.
- ❖ As file grows and shrinks, disk pages are allocated and de-allocated.
- ❖ To support record level operations, we must:
 - keep track of the *pages* in a file
 - keep track of *free space* on pages
 - keep track of the *records* on a page
- ❖ There are many alternatives for keeping track of this.

Heap File Implemented as a List



- ❖ The header page id and Heap file name must be stored someplace.
- ❖ Each page contains 2 'pointers' plus data.

Heap File Using a Page Directory



- ❖ The entry for a page can include the number of free bytes on the page.
- ❖ The directory is a collection of pages; linked list implementation is just one alternative.
 - *Much smaller than linked list of all HF pages!*

Indexes

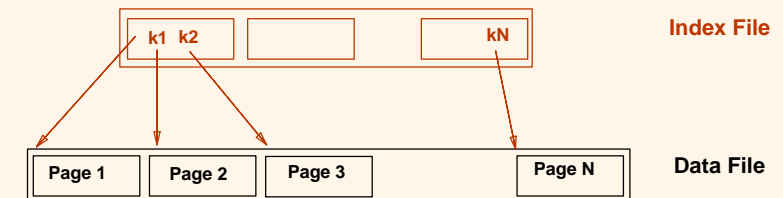


- ❖ An **index** on a file speeds up selections on the **search key fields** for the index.
 - Any subset of the fields of a relation can be the search key for an index on the relation.
 - **Search key** is **not** the same as **key** (minimal set of fields that uniquely identify a record in a relation).
- ❖ 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. (Details soon ...)

Range Searches

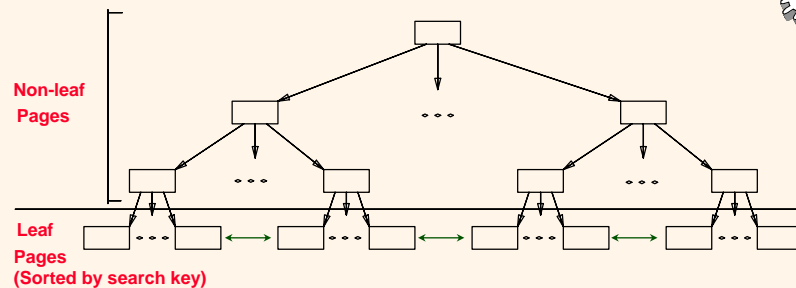


- ❖ ``Find all students with $gpa > 3.0$ ''
 - If data is in sorted file, do binary search to find first such student, then scan to find others.
 - Cost of binary search can be quite high.
- ❖ Simple idea: Create an 'index' file.

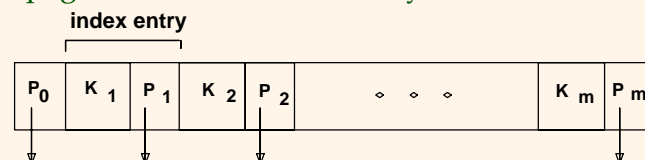


* Can do binary search on (smaller) index file!

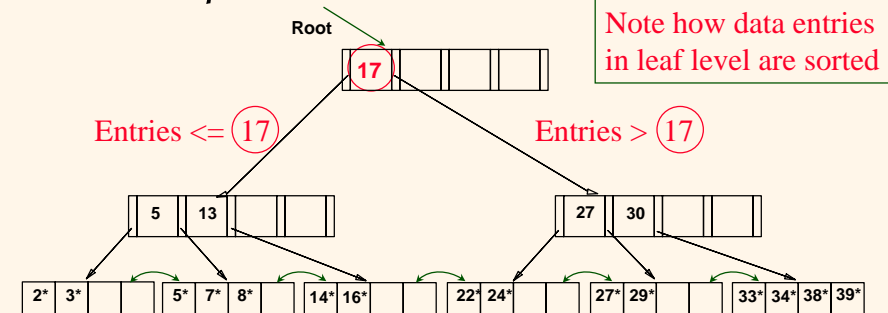
B+ Tree Indexes



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



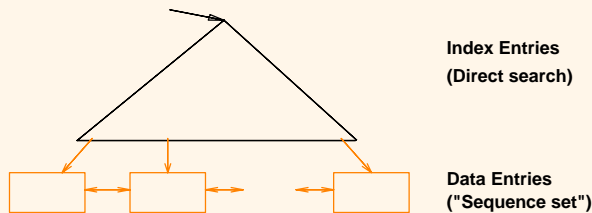
Example B+ Tree



- ❖ 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

B+ Tree: Most Widely Used Index

- ❖ Insert/delete at $\log_F N$ cost; keep tree *height-balanced*. (F = fanout, N = # leaf pages)
- ❖ Minimum 50% occupancy (except for root). Each node contains $d \leq m \leq 2d$ entries. The parameter d is called the *order* of the tree.
- ❖ Supports equality and range-searches efficiently.

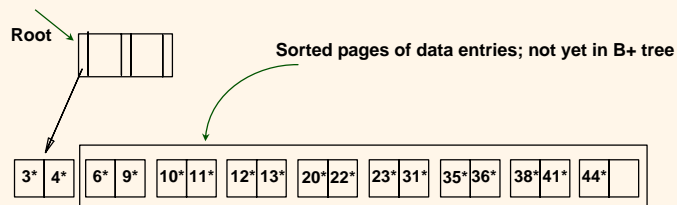


Prefix Key Compression

- ❖ Important to increase fan-out. (Why?)
- ❖ Key values in index entries only 'direct traffic'; can often compress them.
 - E.g., If we have adjacent index entries with search key values *Dannon Yogurt*, *David Smith* and *Devarakonda Murthy*, we can abbreviate *David Smith* to *Dav*. (The other keys can be compressed too ...)
 - Is this correct? Not quite! What if there is a data entry *Davey Jones*? (Can only compress *David Smith* to *Davi*)
 - In general, while compressing, must leave each index entry greater than every key value (in any subtree) to its left.
- ❖ Insert/delete must be suitably modified.

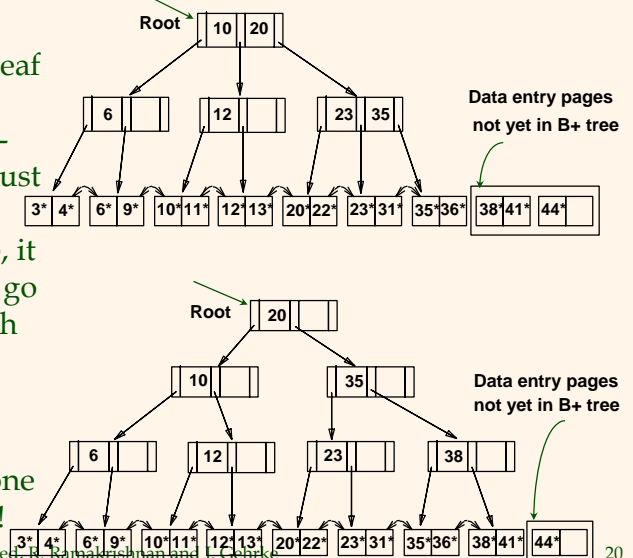
Bulk Loading of a B+ Tree

- ❖ If we have a large collection of records, and we want to create a B+ tree on some field, doing so by repeatedly inserting records is very slow.
- ❖ Bulk Loading can be done much more efficiently.
- ❖ Initialization: Sort all data entries, insert pointer to first (leaf) page in a new (root) page.



Bulk Loading (Contd.)

- ❖ Index entries for leaf pages always entered into right-most index page just above leaf level. When this fills up, it splits. (Split may go up right-most path to the root.)
- ❖ Much faster than repeated inserts, especially when one considers locking!



Summary of Bulk Loading



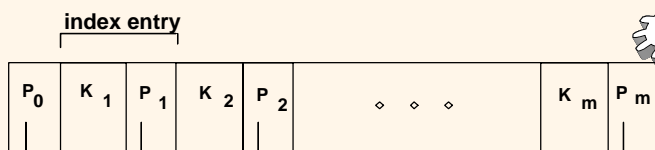
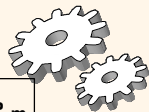
- ❖ Option 1: multiple inserts.
 - Slow.
 - Does not give sequential storage of leaves.
- ❖ Option 2: Bulk Loading
 - Has advantages for concurrency control.
 - Fewer I/Os during build.
 - Leaves will be stored sequentially (and linked, of course).
 - Can control “fill factor” on pages.

A Note on ‘Order’

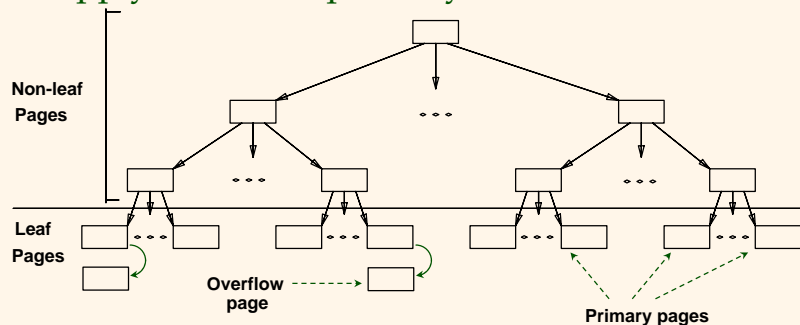


- ❖ *Order (d)* concept replaced by physical space criterion in practice (*‘at least half-full’*).
 - Index pages can typically hold many more entries than leaf pages.
 - Variable sized records and search keys mean different nodes will contain different numbers of entries.
 - Even with fixed length fields, multiple records with the same search key value (*duplicates*) can lead to variable-sized data entries (if we use Alternative (3)).

ISAM

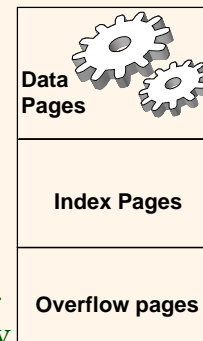


- ❖ Index file may still be quite large. But we can apply the idea repeatedly!



* Leaf pages contain *data entries*.

Comments on ISAM

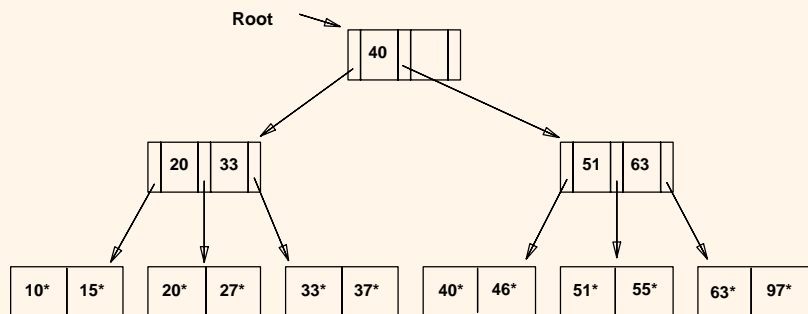


- ❖ *File creation*: Leaf (data) pages allocated sequentially, sorted by search key; then index pages allocated, then space for overflow pages.
- ❖ *Index entries*: *<search key value, page id>*; they ‘direct’ search for *data entries*, which are in leaf pages.
- ❖ *Search*: Start at root; use key comparisons to go to leaf. Cost $\propto \log_F N$; $F = \# \text{ entries/index pg}$, $N = \# \text{ leaf pgs}$
- ❖ *Insert*: Find leaf data entry belongs to, and put it there.
- ❖ *Delete*: Find and remove from leaf; if empty overflow page, de-allocate.
- * **Static tree structure**: *inserts/deletes affect only leaf pages.*

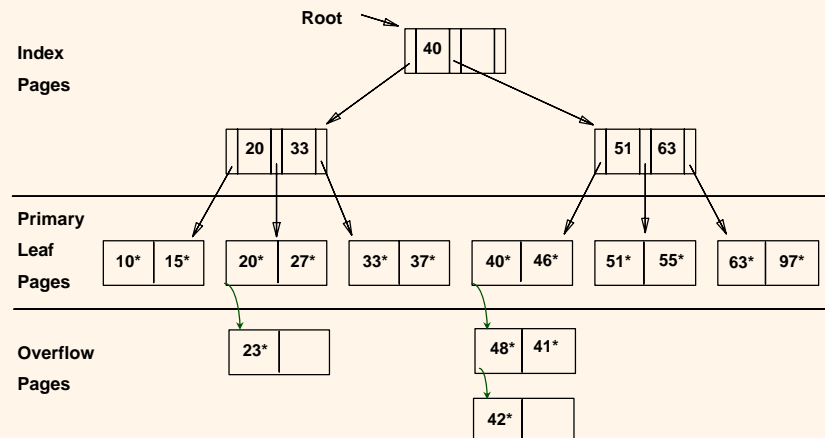
Example ISAM Tree



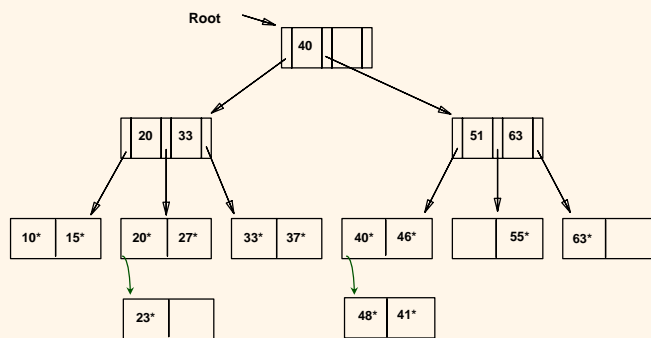
- ❖ Each node can hold 2 entries; no need for 'next-leaf-page' pointers. (Why?)



After Inserting 23*, 48*, 41*, 42* ...



... Then Deleting 42*, 51*, 97*



* Note that 51* appears in index levels, but not in leaf!

Hash-Based Indexes

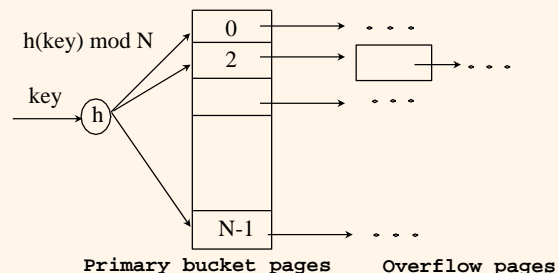


- ❖ Good for equality selections.
- ❖ Index is a collection of buckets.
 - Bucket = *primary page* plus zero or more *overflow pages*.
 - Buckets contain data entries.
- ❖ **Hashing function h :** $h(r)$ = bucket in which (data entry for) record r belongs. h looks at the *search key* fields of r .
 -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

Static Hashing



- ❖ # primary pages fixed, allocated sequentially, never de-allocated; overflow pages if needed.
- ❖ $h(k) \bmod M$ = bucket to which data entry with key k belongs. (M = # of buckets)



Static Hashing (Contd.)



- ❖ Buckets contain *data entries*.
- ❖ Hash fn works on *search key* field of record r . Must distribute values over range $0 \dots M-1$.
 - $h(\text{key}) = (a * \text{key} + b)$ usually works well.
 - a and b are constants; lots known about how to tune h .
- ❖ **Long overflow chains** can develop and degrade performance.
 - *Extendible* and *Linear Hashing*: Dynamic techniques to fix this problem.

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:

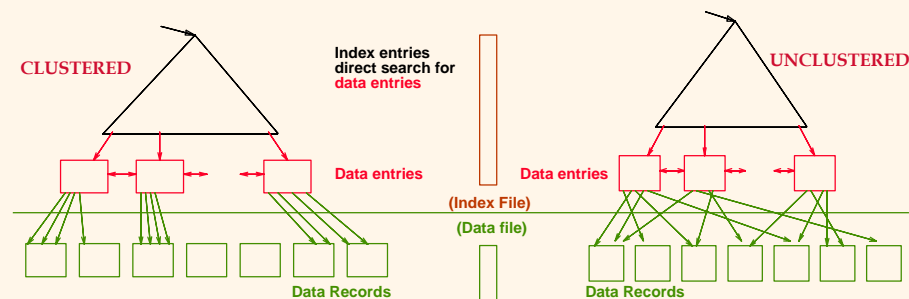
- Data entries typically much smaller than data records. So, better than Alternative 1 with large data records, especially if search keys are small. (Portion of index structure used to direct search, which depends on size of data entries, is much smaller than with Alternative 1.)
- Alternative 3 more 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

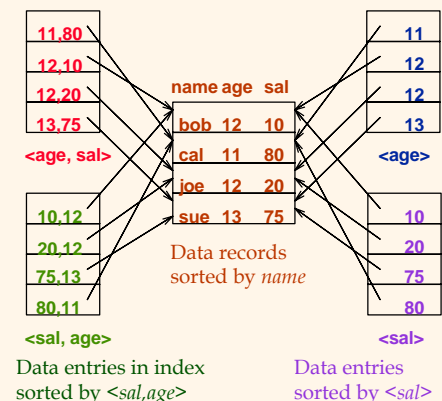
- ❖ Suppose that Alternative (2) is used for data entries, and that the data records are stored in a Heap file.
 - To build clustered index, first sort the Heap file (with some free space on each page for future inserts).
 - Overflow pages may be needed for inserts. (Thus, order of data recs is 'close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



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:
 - $\text{age}=20$ and $\text{sal}=75$
 - **Range query:** Some field value is not a constant. E.g.:
 - $\text{age}=20$; or $\text{age}=20$ and $\text{sal} > 10$
- ❖ 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.



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 $\langle age, sal \rangle$ index much better than $\langle sal, age \rangle$ index!
- ❖ Composite indexes are larger, updated more often.