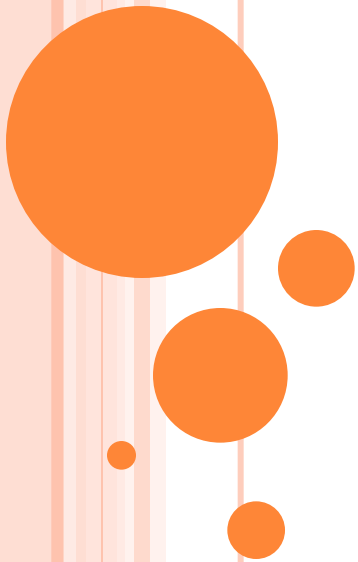# Overview of Storage and Indexing

# 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. Page: typically 4 Kbytes.

# 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**.
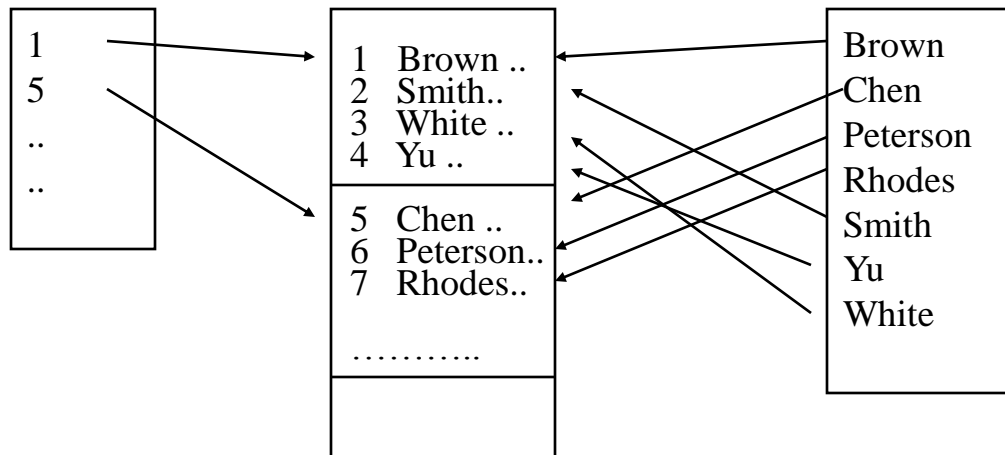
# INDEX CLASSIFICATION

- *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 order of data entries, then called **clustered** index.
    - 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!
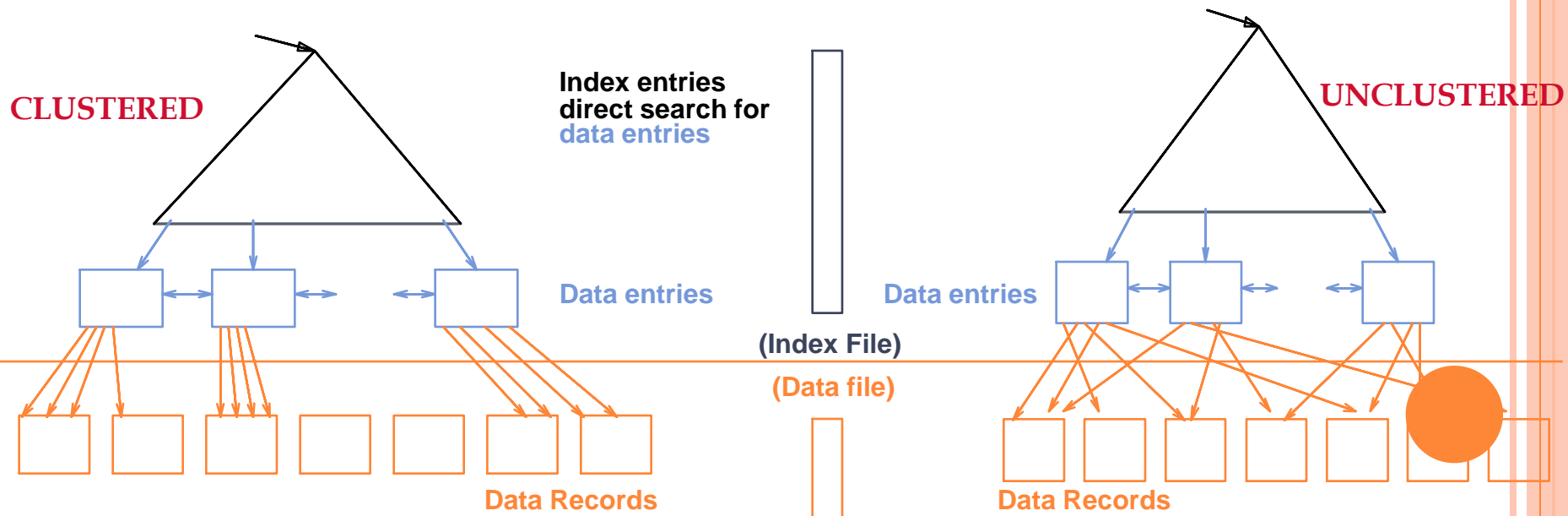
# INDEX CLASSIFICATION

- *Dense* vs *Sparse*:  If there is an entry in the index for each key value -> **dense index** (unclustered indices are dense). If there is an entry for each page -> **sparse index**.

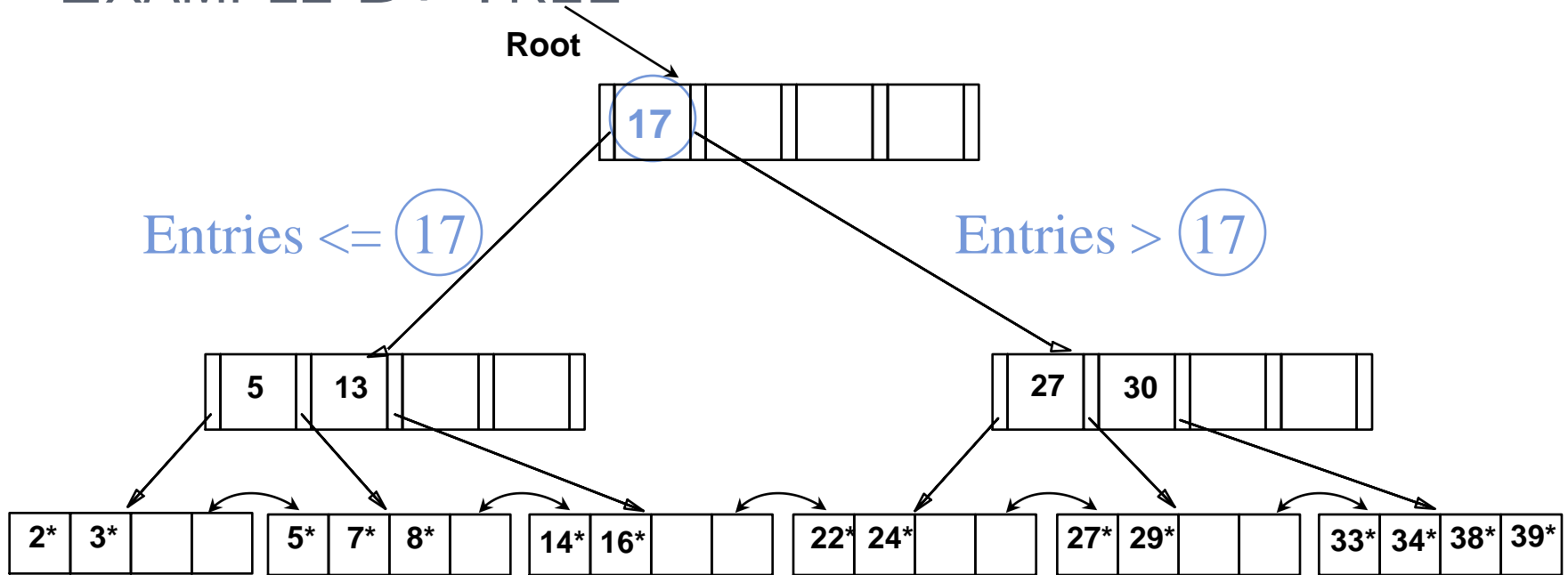| 1 | | 1 Brown .. | | Brown |
| 5 | | 2 Smith.. | | Chen |
| .. | | 3 White .. | | Peterson |
| .. | | 4 Yu .. | | Rhodes |
| | | 5 Chen .. | | Smith |
| | | 6 Peterson.. | | Yu |
| | | 7 Rhodes.. | | White |
| | | ……….. | | |

# CLUSTERED VS. UNCLUSTERED INDEX

- 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.)

**CLUSTERED**

**Index entries direct search for data entries**

**UNCLUSTERED**

**Data entries**

**Data entries**

**(Index File)**

**(Data file)**

**Data Records**

**Data Records**

# EXAMPLE B+ TREE

**Root**

17

Entries <= 17                    Entries > 17

| 5 | 13 |    |    |

| 27 | 30 |    |    |

| 2* | 3* |    |    |

| 5* | 7* | 8* |    |

| 14* | 16* |    |    |

| 22* | 24* |    |    |

| 27* | 29* |    |    |

| 33* | 34* | 38* | 39* |

- Good for range queries.
- Insert/delete:  Find data entry in leaf, then change it. Need to adjust parent sometimes. All leaves at he same height.
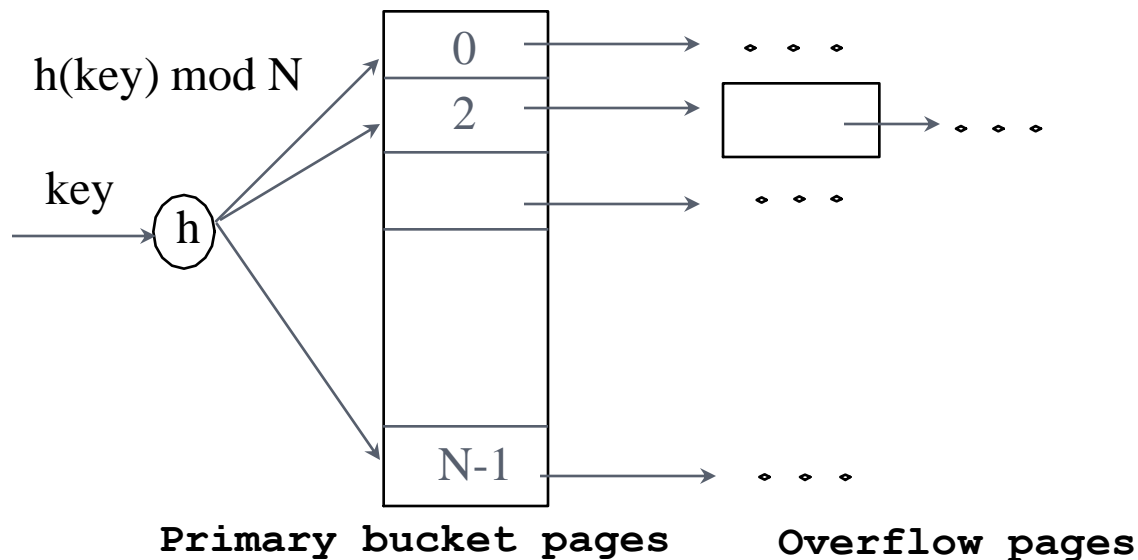
# Hash-Based Indexes

- 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 $r$.
- Buckets may contain the data records or just the rids.
- _Hash-based_ indexes are best for _equality selections_. **Cannot** support range searches

# STATIC HASHING

- # primary pages fixed, allocated sequentially, never de-allocated; overflow pages if needed.

- **h**(*k*) mod N = bucket to which data entry with key *k* belongs. (N = # of buckets)

- Long overflow chains can develop and degrade performance.
  - *Extendible* and *Linear Hashing*: Dynamic techniques to fix this.

h(key) mod N

key

h

0

2

N-1

**Primary bucket pages**      **Overflow pages**

# Static Hashing (Contd.)

- Buckets contain *data entries*.
- Hash fn works on *search key* field of record *r.* Must distribute values over range 0 ... M-1.
  - **h**(*key*) = (a * *key* + b) usually works well.
  - a and b are constants; lots known about how to tune **h**.

# 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.

⊠ *Good enough to show the overall trends!*

# COMPARING FILE ORGANIZATIONS

- Heap files (random order; insert at eof)
- Sorted files, sorted on *<age, sal>*
- Clustered B+ tree file, Alternative (1), search key *<age, sal>*
- Heap file with unclustered B + tree index on search key *<age, sal>*
- Heap file with unclustered hash index on search key *<age, sal>*

# Operations to Compare

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

# 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** | $D\log_2 B$ | $D\log_2 B +$ **# matches** | **Search + BD** | **Search +BD** |
| (3) Clustered | **1.5BD** | $D\log_F 1.5B$ | $D\log_2 1.5B +$ **# matches** | **Search + D** | **Search +D** |
| (4) Unclustered Tree index | **BD(R+0.15)** | $D(1 +\log_F 0.15B)$ | $D\log_F 0.15B +$ **# matches** | $D(3 +\log_F 0.15B)$ | **Search + 2D** |
| (5) Unclustered Hash index | **BD(R+0.125)** | **2D** | **BD** | **4D** | **Search + 2D** |

**B:** The number of data pages
**R:** Number of records per page
**D:** (Average) time to read or write disk page

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

# CHOICE OF INDEXES

- What indexes should we create?

- 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.
- Multi-attribute search keys should be considered when a WHERE clause contains several conditions.
- 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?
- 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

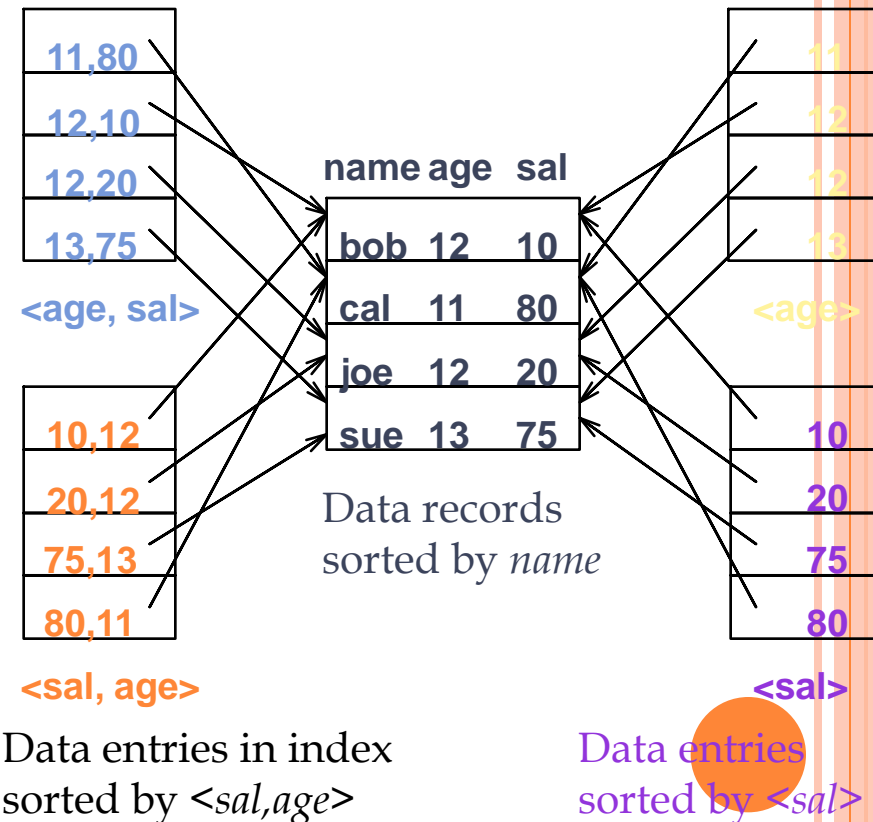SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age>10
GROUP BY E.dno

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

# 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.
- Order or attributes is relevant.

Examples of composite key



name age  sal

| bob | 12 | 10 |
| cal | 11 | 80 |
| joe | 12 | 20 |
| sue | 13 | 75 |

<age, sal>

10,12
20,12
75,13
80,11

<sal, age>

11
12
12
13

<age>

10
20
75
80

<sal>

Data records sorted by *name*

Data entries in index sorted by <*sal,age*>

Data entries sorted by <*sal*>

# 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.

# 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.

# 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.