# Storage, Disks, and File Organization

EECS 339

Lecture 10

# Data on External Storage

- <u>Disks:</u> Can retrieve random page at fixed cost
  - But reading several consecutive pages is much cheaper than reading them in random order
- <u>Tapes:</u> Can only read pages in sequence
  - Cheaper than disks; used for archival storage
- <u>File organization:</u> 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
- <u>Architecture:</u> Buffer manager stages pages from external storage to main memory buffer pool. File and index layers make calls to the buffer manager.
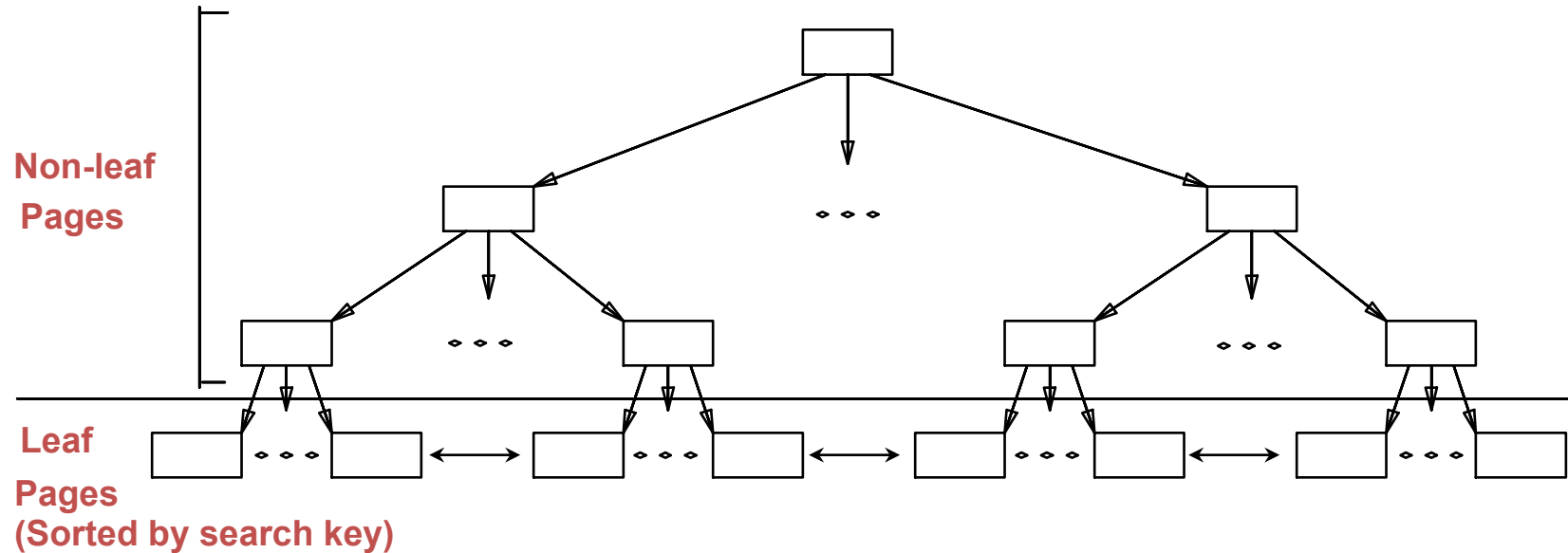
# File Organization Choices

Many choices 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
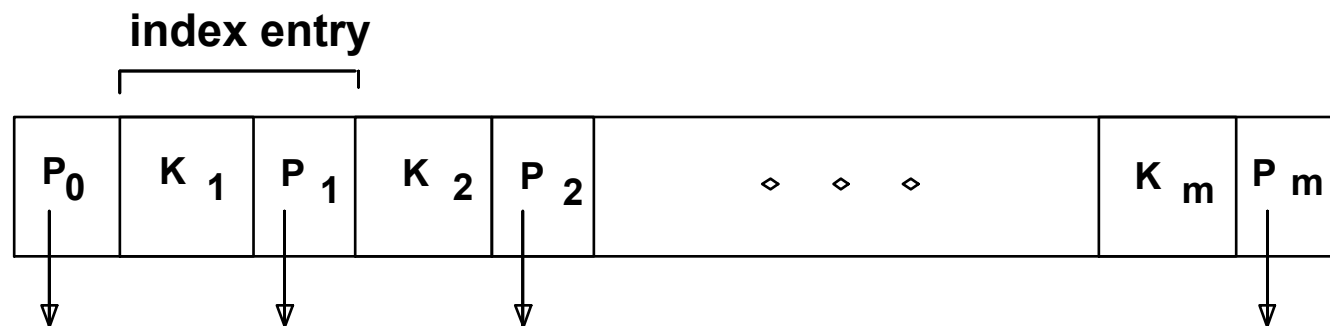    - Updates are much faster than in sorted files.

# 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.  (Details soon …)
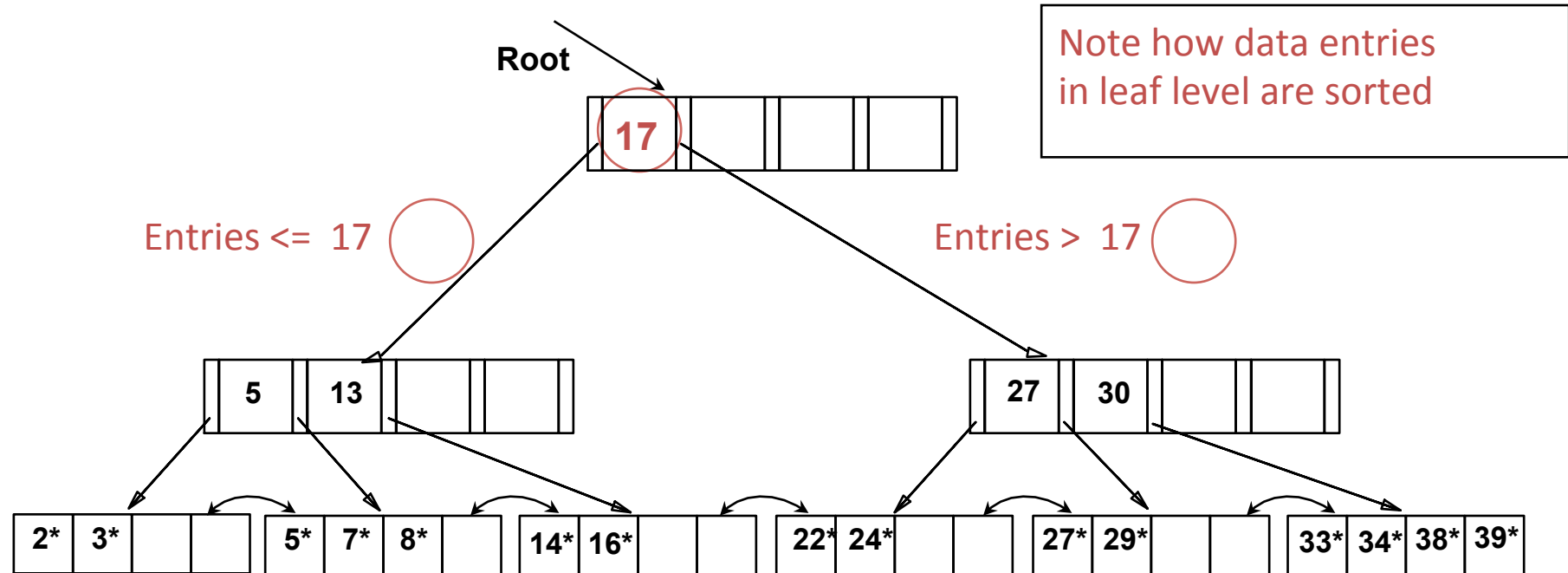
# B+ Tree Indexes

**Non-leaf Pages**

**Leaf Pages (Sorted by search key)**

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

**index entry**

| $P_0$ | $K_1$ | $P_1$ | $K_2$ | $P_2$ | $\diamond \quad \diamond \quad \diamond$ | $K_m$ | $P_m$ |
|---|---|---|---|---|---|---|---|

# 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* |

- 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.*

# Options for Data Entry k* in Index

- In a data entry k* we can store:
  - Data record with key value **k,** or
  - <**k**, rid of data record with search key value **k**>, or
  - <**k**, list of rids of data records with search key **k**>
- Choice of representation 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

# Options for Data Entries (Contd.)

- Choice 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 Choice 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.
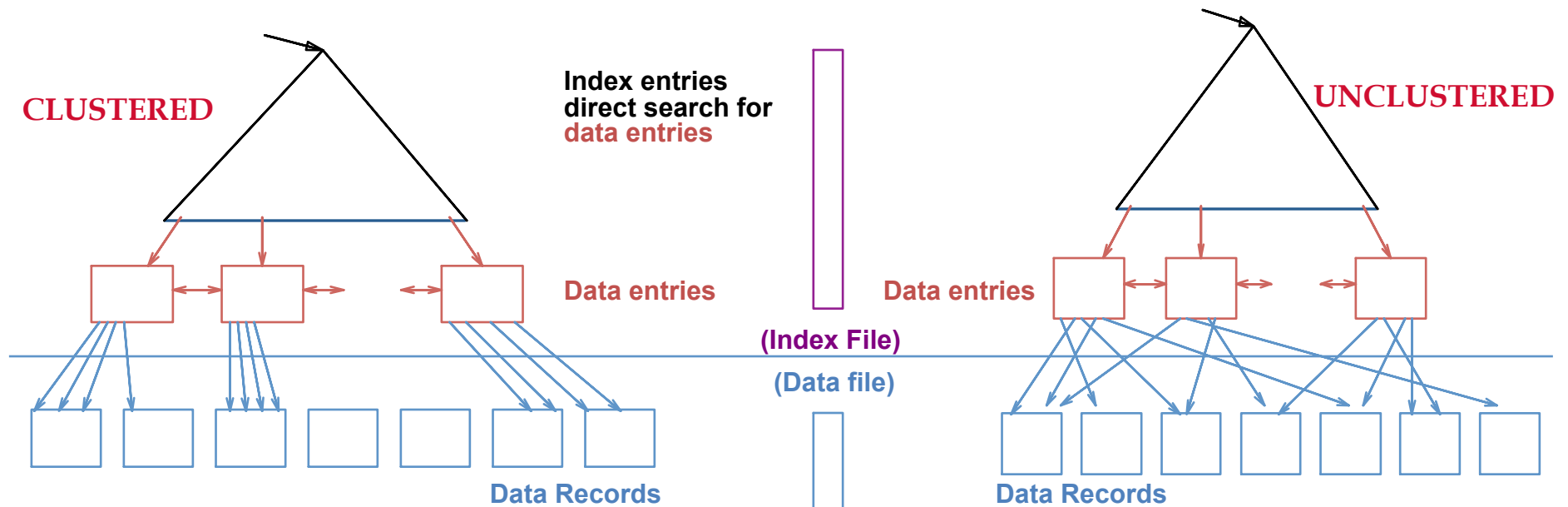
# Alternatives for Data Entries (Contd.)

- Choices 2 and 3:
  - Data entries typically much smaller than data records. So, better than Choice 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.
  - Choice 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 Choice (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.)

**CLUSTERED**

**Index entries
direct search for
data entries**

**UNCLUSTERED**

Data entries

Data entries

**(Index File)**

**(Data file)**

**Data Records**

**Data Records**

# 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, Choice (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

# Assumptions in Our Analysis

- Heap Files:
  - Equality selection on key; exactly one match.
- Sorted Files:
  - Files compacted after deletions.
- Indexes:
  - Choice (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

# Assumptions (contd.)

- Scans:
  - Leaf levels of a tree-index are chained.
  - Index data-entries plus actual file scanned for unclustered indexes.
- Range searches:
  - We use tree indexes to restrict the set of data records fetched, but ignore hash indexes.

# Cost of Operations

|  | (a) Scan | (b) Equality | (c ) Range | (d) Insert | (e) Delete |
|---|---|---|---|---|---|
| (1) Heap |  |  |  |  |  |
| (2) Sorted |  |  |  |  |  |
| (3) Clustered |  |  |  |  |  |
| (4) Unclustered Tree index |  |  |  |  |  |
| (5) Unclustered Hash index |  |  |  |  |  |

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

# 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 | $D\log_2 B$ | $D(\log_2 B +$ # pgs with match recs) | Search + BD | Search +BD |
| (3) Clustered | 1.5BD | $D\log_F 1.5B$ | $D(\log_F 1.5B$ + # pgs w. match recs) | Search + D | Search +D |
| (4) Unclust. Tree index | BD(R+0.15) | $D(1 + \log_F 0.15B)$ | $D(\log_F 0.15B$ + # pgs w. 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!*

# Study Break

- Consider a delete specified using an equality condition. For each of the five file organizations, what is the cost if no record qualifies? What is the cost if the condition is not on a key?

# Study Break Solution

- If the search key is not a candidate key, there may be several qualifying records. In a heap file, this means we have to search the entire file to be sure that we've found all qualifying records; the cost is B(D + RC). In a sorted file, we find the first record (cost is that of equality search; Dlog2B + Clog2R) and then retrieve and delete successive records until the key value changes. The cost of the deletions is C per deleted record, and D per page containing such a record. In a hashed file, we hash to find the appropriate bucket (cost H), then retrieve the page (cost D; let's assume no overflow pages), then write the page back if we find a qualifying record and delete it (cost D).

- If no record qualifies, in a heap file, we have to search the entire file. So the cost is B(D + RC). In a sorted file, even if no record qualifies, we have to do equality search to verify that no qualifying record exists. So the cost is the same as equality search, Dlog2B + Clog2R. In a hashed file, if no record qualifies, assuming no overflow page, we compute the hash value to find the bucket that would contain such a record (cost is H), bring that page in (cost is D), and search the entire page to verify that the record is not there(cost is RC). So the total cost is H+D+RC.

- In all three file organizations, if the condition is not on the search key we have to search the entire file. There is an additional cost of C for each record that is deleted, and an additional D for each page containing such a record.

# 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.
- 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

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

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

- 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.hobby=Stamps
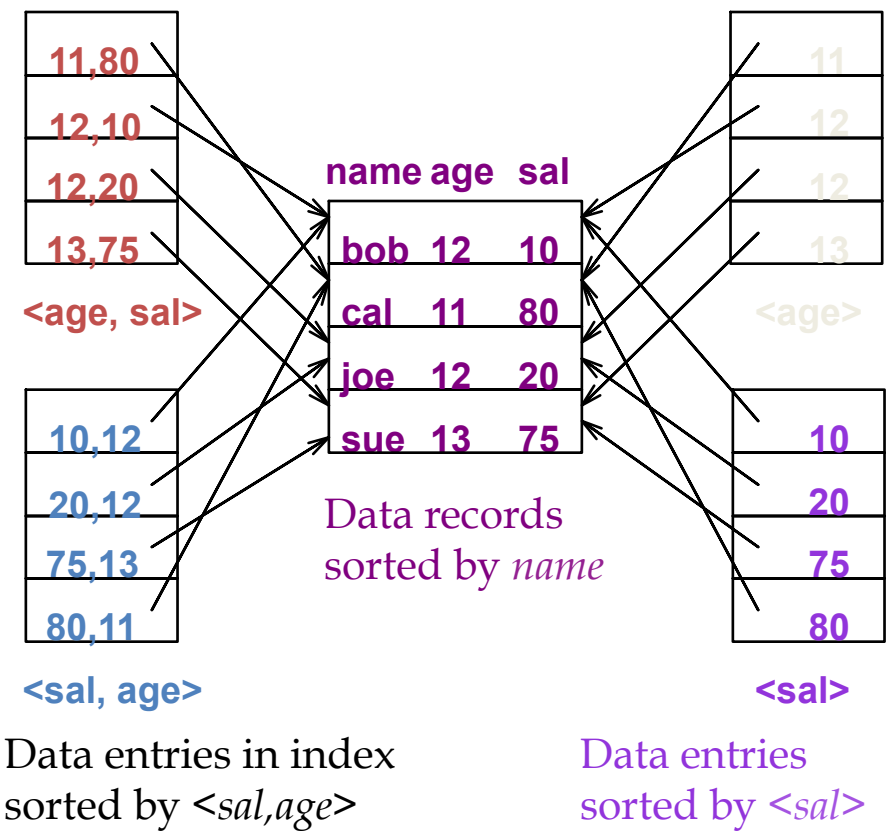
# Indexes with Composite Search Keys

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

Examples of composite key indexes using lexicographic order.



| | name | age | sal |
|---|---|---|---|
| | bob | 12 | 10 |
| | cal | 11 | 80 |
| | joe | 12 | 20 |
| | sue | 13 | 75 |

11,80
12,10
12,20
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.

# 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>*

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

*<E.dno,E.sal>*

*Tree index!*

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

*<E. age,E.sal>*
or
*<E.sal, E.age>*

*Tree index!*

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

# Index-Only Plans (Contd.)

- Index-only plans are possible if the key is <dno,age> or we have a tree index with key <age,dno>
  – Which is better?
  – What if we consider the second query?

```
SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age=30
GROUP BY E.dno
```

```
SELECT  E.dno,  COUNT (*)
FROM  Emp E
WHERE  E.age>30
GROUP BY E.dno
```

# Index-Only Plans (Contd.)

- Index-only plans can also be found for queries involving more than one table; more on this later.

*<E.dno>*

```
SELECT  D.mgr
FROM  Dept D, Emp E
WHERE  D.dno=E.dno
```

*<E.dno,E.eid>*

```
SELECT  D.mgr, E.eid
FROM  Dept D, Emp E
WHERE  D.dno=E.dno
```

# Recap

- 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.

# Recap (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.

# Recap (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.

# 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!
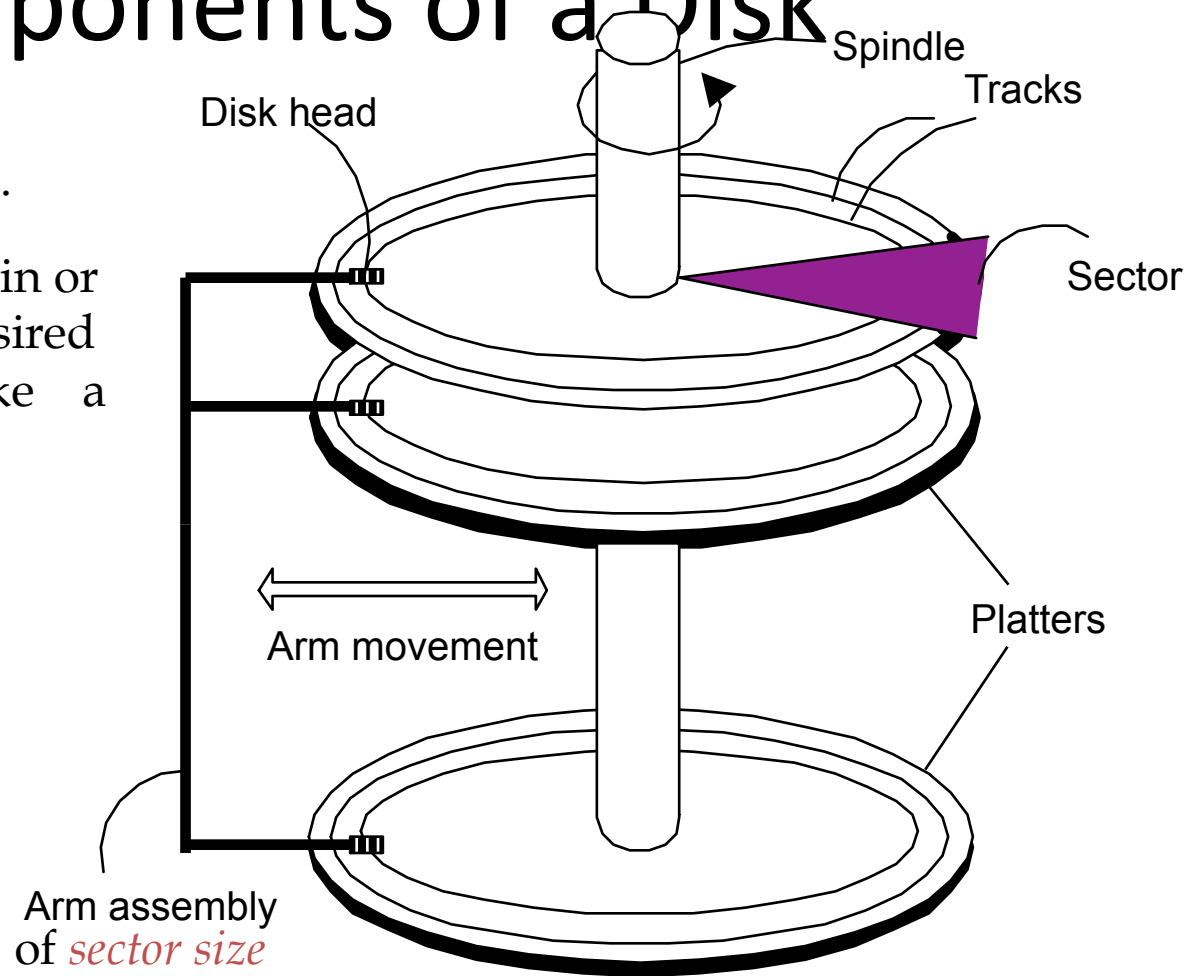
# Why Not Store Everything in Main Memory?

- *Costs too much.* $1000 will buy you either 128 GB of RAM or 16 TB of disk today.
- *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).

# 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!

# 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 (which is fixed).

Spindle

Tracks

Disk head

Sector

Arm movement

Platters

Arm assembly of *sector size*

# 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?

# 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!

# Study Break: Disk Costs

Consider a disk with a sector size of 512 bytes, 2000 tracks per surface, 50 sectors per track, five double-sided platters, and average seek time of 10 msec.

- What is the capacity of a track in bytes? What is the capacity of each surface? What is the capacity of the disk?

- How many cylinders does the disk have?

- Give examples of valid block sizes. Is 256 bytes a valid block size? 2048? 51200?

- If the disk platters rotate at 5400 rpm (revolutions per minute), what is the maximum rotational delay?

- If one track of data can be transferred per revolution, what is the transfer rate?

# Disk Costs: Solution

1.
- bytes/track = bytes/sector × sectors/track = 512 × 50 = 25K bytes/surface = bytes/track × tracks/surf ace = 25K × 2000 = 50, 000K

- bytes/disk = bytes/surf ace × surf aces/disk = 50, 000K × 5 × 2 = 500, 000K

2. The number of cylinders is the same as the number of tracks on each platter, which is 2000.

3. The block size should be a multiple of the sector size. We can see that 256 is not a valid block size while 2048 is. 51200 is not a valid block size in this case because block size cannot exceed the size of a track, which is 25600 bytes.

4. If the disk platters rotate at 5400rpm, the time required for one complete rotation, which is the maximum rotational delay, is

- $1 × \frac{60}{5400} = 0.011$ seconds

- . The average rotational delay is half of the rotation time, 0.006 seconds.

5. The capacity of a track is 25K bytes. Since one track of data can be transferred per revolution, the data transfer rate is

- $\frac{25K}{0.011} = 2,250K$ bytes/second

# 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.

# 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

# 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

# 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.

# Buffer Management in a DBMS

Page Requests from Higher Levels

BUFFER POOL

disk page

free frame

MAIN MEMORY

DISK

DB

choice of frame dictated by **replacement policy**

- *Data must be in RAM for DBMS to operate on it!*
- *Table of <frame#, pageid> pairs is maintained.*

# When a Page is Requested …

- If requested page is not in pool:
  - Choose a frame for *replacement*
  - If  frame is dirty, write it to disk
  - Read requested page into chosen frame
- *Pin* the page and return its address.

☛ *If requests can be predicted (e.g., sequential scans) pages can be* <u>pre-fetched</u> *several pages at a time!*

# More on Buffer Management

- Requestor of page must unpin it, and indicate whether page has been modified:
  - *dirty* bit is used for this.
- Page in pool may be requested many times,
  - a *pin count* is used.  A page is a candidate for replacement iff *pin count* = 0.
- CC & recovery may entail additional I/O when a frame is chosen for replacement. (*Write-Ahead Log* protocol; more later.)

# DBMS vs. OS File System

OS does disk space & buffer mgmt: why not let OS manage these tasks?

- Differences in OS support: portability issues
- Some limitations, e.g., files can't span disks.
- Buffer management in DBMS requires ability to:
  - pin a page in buffer pool, force a page to disk (important for implementing CC & recovery),
  - adjust *replacement policy,* and pre-fetch pages based on access patterns in typical DB operations.

# Buffer Replacement Policy

- Frame is chosen for replacement by a *replacement policy:*
  - Goal: minimize *cache misses*
- Policy can have big impact on # of I/Os; depends on the *access pattern*.
- *Sequential flooding*:  Nasty situation caused by LRU + repeated sequential scans.
  - # buffer frames < # pages in file means each page request causes an I/O.  MRU much better in this situation (but not in all situations, of course).

# Least Recently Used Replacement

- When the buffer pool is full, evict the oldest data
  - Data must not be pinned
  - If dirty, needs to be flushed to disk
- Implemented with a queue of pages
- Each page added when pin count goes to 0

# LRU w/ 3 pages

| A | B | C | Read | Hit/miss? |
|---|---|---|---|---|
| 1 | | | 1 | m |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# LRU w/ 3 pages

| A | B | C | Read | Hit/miss? |
|---|---|---|---|---|
| 1 | | | 1 | m |
| 1 | 2 | | 2 | m |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

54

# LRU w/ 3 pages

| A | B | C | Read | Hit/miss? |
|---|---|---|------|-----------|
| 1 |   |   | 1    | m         |
| 1 | 2 |   | 2    | m         |
| 1 | 2 | 3 | 3    | m         |
|   |   |   |      |           |
|   |   |   |      |           |
|   |   |   |      |           |

# LRU w/ 3 pages

| A | B | C | Read | Hit/miss? |
|---|---|---|------|-----------|
| 1 | | | 1 | m |
| 1 | 2 | | 2 | m |
| 1 | 2 | 3 | 3 | m |
| ~~1~~ 4 | 2 | 3 | 4 | m |
| | | | | |
| | | | | |

# LRU w/ 3 pages

| A | B | C | Read | Hit/miss? |
| --- | --- | --- | --- | --- |
| 1 | | | 1 | m |
| 1 | 2 | | 2 | m |
| 1 | 2 | 3 | 3 | m |
| ~~1~~ 4 | 2 | 3 | 4 | m |
| ~~1~~ 4 | ~~2~~ 1 | 3 | 1 | m |
| | | | | |

# Clock Replacement

- Rather than maintain a queue, treat pages like a circular buffer (or clock face)

- Make sweep over pages in buffer pool to identify ones that have not been accessed recently

- Lower overhead than LRU

# Clock Replacement

- c = *current* page in buffer pool (1…N pages), circular buffer
- r = *referenced* bit, 1 / page, buffer pool sets to 1 when a page's pin count goes to 0

- c = 1
- while(c.r != 0)

```
if(c.pinCount > 0) ++c;
if(c.r == 1) c.r = 0; ++c;
if(c > N) c = 0;
```

# Memory Management

- When preparing to execute a query, the database calculates how many pages it will need in the buffer pool

- We use this figure to make sure that memory is not overcommitted
  - Otherwise the system thrashes

- Do not run queries that won't fit.

# Query Memory Footprint (LRU)

- Sequential table scan = 1 page

- Random index access = 1 page

- Join = sizeof(inner relation)

- Index lookup – 1 page

- Repeated index lookup (e.g., inner join, B+ Tree) – top few pages of tree + footprint of bottom level

# Admission Control

- Queries arrive in a queue
- Analyze each, determining # of pages it needs to run well
- Greedily admit first query that will fit.
  - In practice accompanied by anti-starvation strategies for large queries.
    - Won't consider that here.

# Memory Management Study Break

- You have a queue of queries with the following working set sizes (in pages):

- 13, 16, 11, 18, 20, 5, 2, 15, 7, 3, 10, 1, 8, 14, 4

- Your buffer pool contains 30 pages.  If the queries all have the same duration, in what order will they be executed?

# Memory Management Solution

- Schedules:
- 13, 16, 1
- 11, 18
- 20, 5, 2, 3
- 15, 7, 8
- 10, 14, 4