



# Chapter 10: Storage and File Structure

**Edited by Radhika Sukapuram**

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use

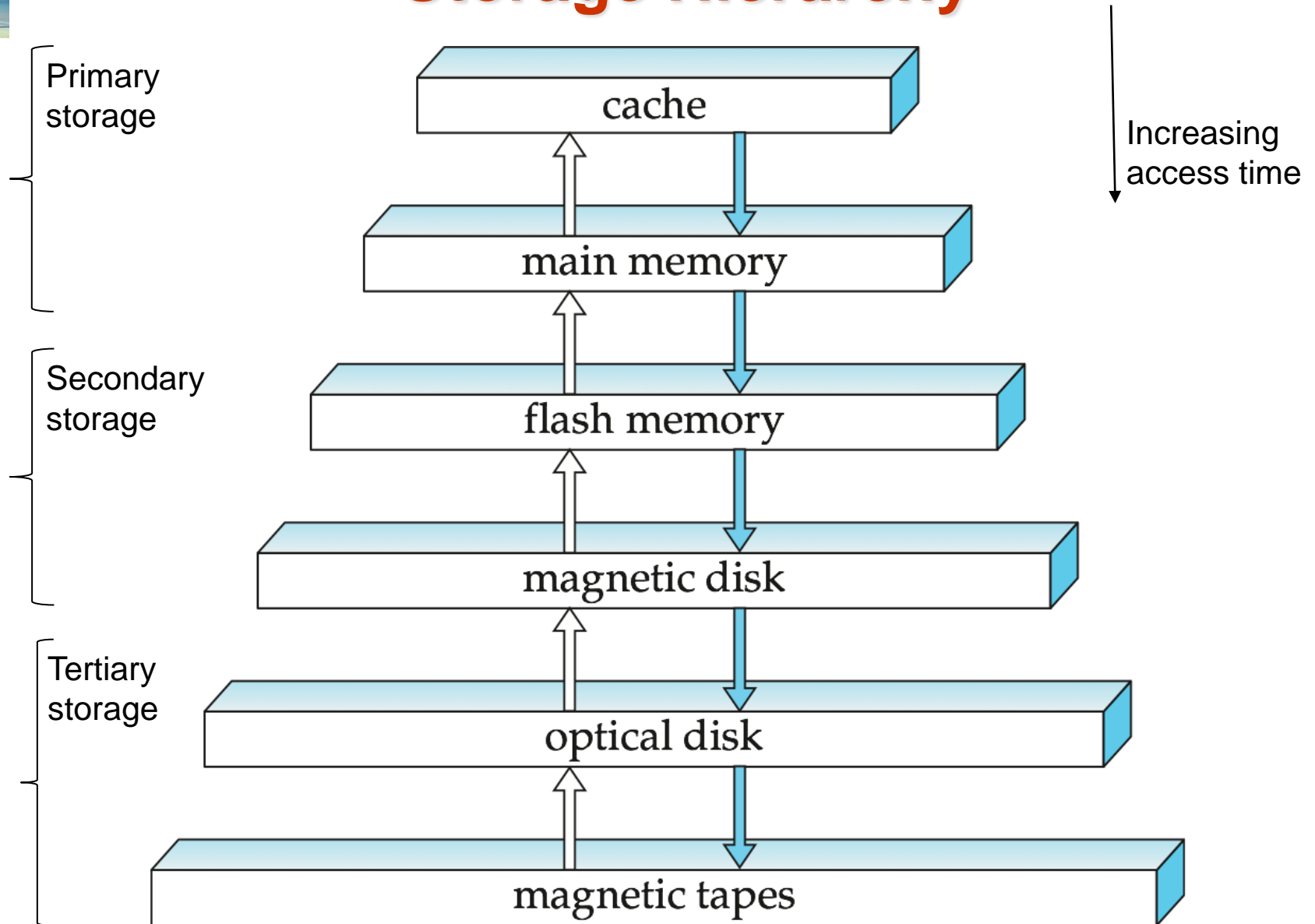


# Classification of Physical Storage Media

- Speed with which data can be accessed
- Cost per unit of data
- Reliability
  - data loss on power failure or system crash
  - physical failure of the storage device
- Can differentiate storage into:
  - **volatile storage**: loses contents when power is switched off
  - **non-volatile storage**:
    - ▶ Contents persist even when power is switched off.
    - ▶ Includes secondary and tertiary storage, as well as battery-backed up main-memory.

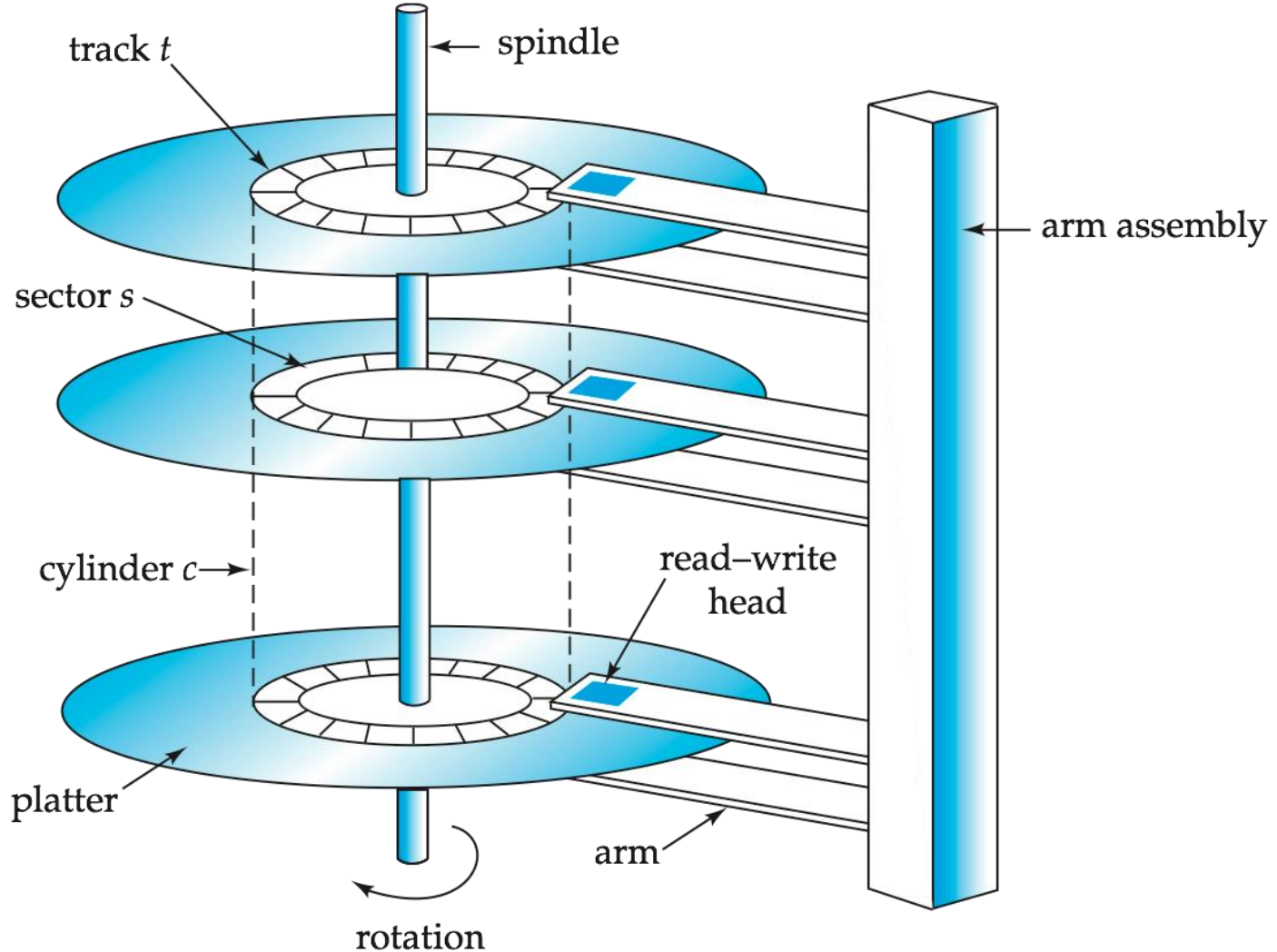


# Storage Hierarchy





# Magnetic Hard Disk Mechanism



**NOTE: Diagram is schematic, and simplifies the structure of actual disk drives**



# Magnetic Disks

- **Read-write head**
  - Positioned very close to the platter surface (almost touching it)
  - Reads or writes magnetically encoded information.
- Surface of platter divided into circular **tracks**
  - Over 50K-100K tracks per platter on typical hard disks
- Each track is divided into **sectors**.
  - A sector is the smallest unit of data that can be read or written.
  - Sector size typically 512 bytes
  - Typical sectors per track: 500 to 1000 (on inner tracks) to 1000 to 2000 (on outer tracks)
- To read/write a sector
  - disk arm swings to position head on right track
  - platter spins continually; data is read/written as sector passes under head
- Head-disk assemblies
  - multiple disk platters on a single spindle (1 to 5 usually)
  - one head per platter, mounted on a common arm.
- **Cylinder**  $i$  consists of  $i^{\text{th}}$  track of all the platters



# Magnetic Disks (Cont.)

- **Disk controller** – interfaces between the computer system and the disk drive hardware.
  - accepts high-level commands to read or write a sector
  - initiates actions such as moving the disk arm to the right track and actually reading or writing the data
  - Computes and attaches **checksums** to each sector to verify that data is read back correctly
    - ▶ If data is corrupted, with very high probability stored checksum won't match recomputed checksum
  - Ensures successful writing by reading back sector after writing it
  - Performs **remapping of bad sectors**



# Performance Measures of Disks

- **Access time** – the time it takes from when a read or write request is issued to when data transfer begins. Consists of:
  - **Seek time** – time it takes to reposition the arm over the correct track.
    - ▶ 4 to 10 milliseconds on typical disks
  - **Rotational latency** – time it takes for the sector to be accessed to appear under the head.
    - ▶ 4 to 11 milliseconds on typical disks (5400 to 15000 r.p.m.)
- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
  - 25 to 100 MB per second max rate, lower for inner tracks



# Performance Measures (Cont.)

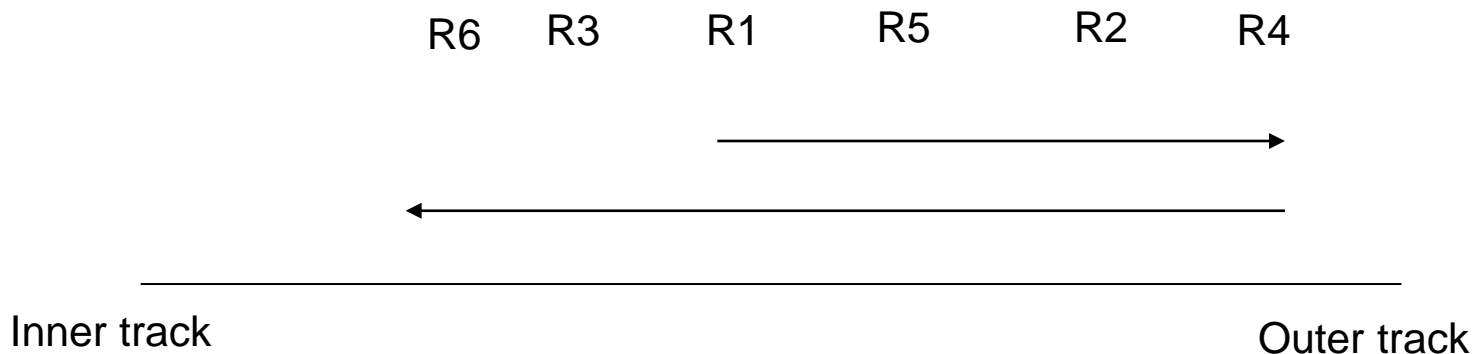
- **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.
  - Typically 3 to 5 years
  - Probability of failure of new disks is quite low, corresponding to a “theoretical MTTF” of 500,000 to 1,200,000 hours for a new disk
    - ▶ E.g., an MTTF of 1,200,000 hours for a new disk means that given 1000 relatively new drives, on an average one will fail every 1200 hours
  - MTTF decreases as disk ages





# Optimization of Disk-Block Access

- **Block** – a contiguous sequence of sectors from a single track
  - data is transferred between disk and main memory in blocks
  - sizes range from 512 bytes to several kilobytes
    - ▶ Smaller blocks: more transfers from disk
    - ▶ Larger blocks: more space wasted due to partially filled blocks
    - ▶ Typical block sizes today range from 4 to 16 kilobytes
- **Disk-arm-scheduling** algorithms order pending accesses to tracks so that disk arm movement is minimized
  - **elevator algorithm**:  $R_n$ : Read request





# Optimization of Disk Block Access (Cont.)

- **File organization** – optimize block access time by organizing the blocks to correspond to how data will be accessed
  - E.g. Store related information on the same or nearby cylinders.
  - Files may get **fragmented** over time
    - ▶ E.g. if data is inserted to/deleted from the file
    - ▶ Or free blocks on disk are scattered, and newly created file has its blocks scattered over the disk
    - ▶ Sequential access to a fragmented file results in increased disk arm movement
  - Some systems have utilities to **defragment** the file system, in order to speed up file access

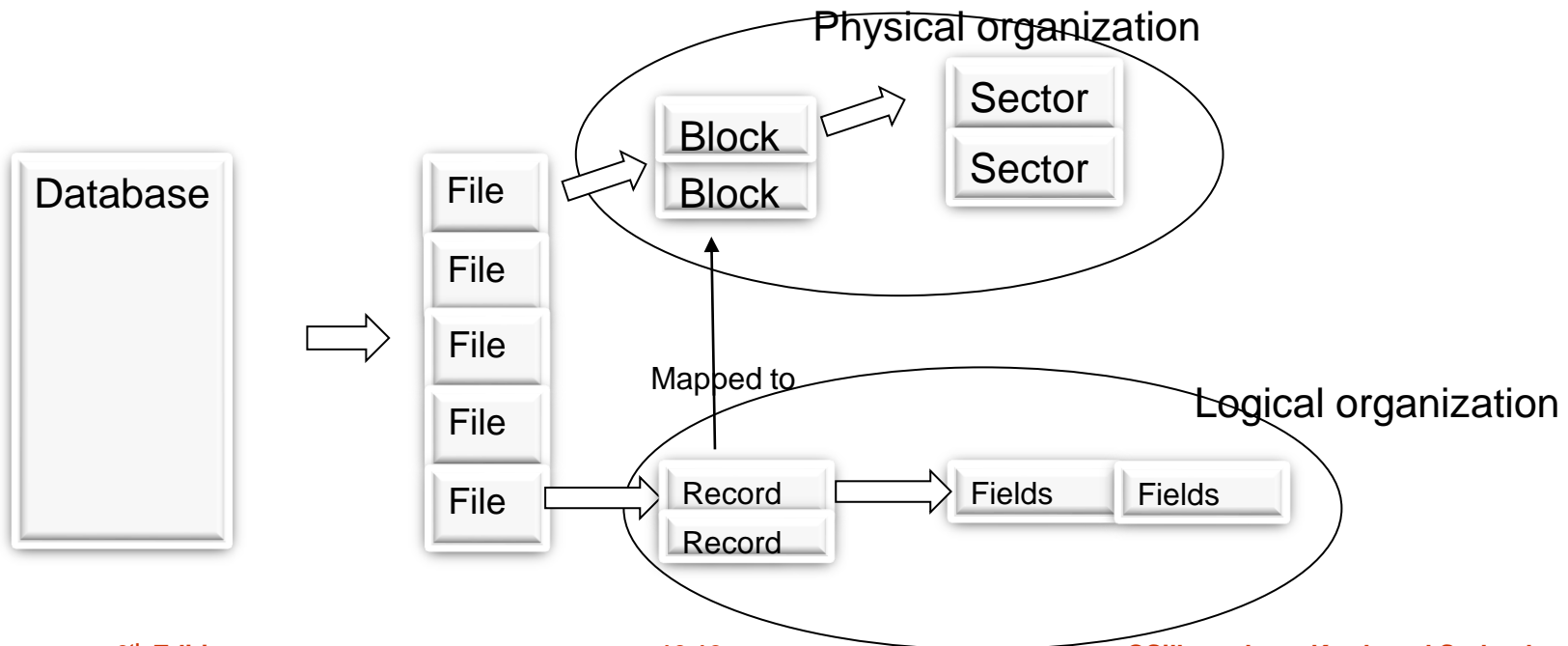


# **File Organization, Record Organization and Storage Access**



# File Organization

- The database is stored as a collection of *files*.
  - Each file is a sequence of *records*.
  - A record is a sequence of fields.
  - A database file is partitioned into fixed-length *storage units* called **blocks**.
    - ▶ Blocks are units of both storage allocation and data transfer.
    - ▶ A block may contain several records





# File organization cont.

- How to represent records in a file ?
- One approach:
  - assume record size is fixed
  - each file has records of only one particular type
  - do not allow records to cross block boundaries
  - different files are used for different relations

This case is easiest to implement; variable length records are also possible.



# Organization of Records in Files

- We discussed how to represent records in a file. But how to organize them in a file ?
- **Heap** – a record can be placed anywhere in the file where there is space
- **Sequential** – store records in sequential order, based on the *value of the search key* of each record
- **Hashing** – a hash function computed on *some attribute* of each record; the result specifies in which block of the file the record should be placed
- Records of each relation may be stored in a separate file.
  - In a **multitable clustering file organization** records of several different relations can be stored in the same file
  - Motivation: store related records on the same block to minimize I/O



# Sequential File Organization

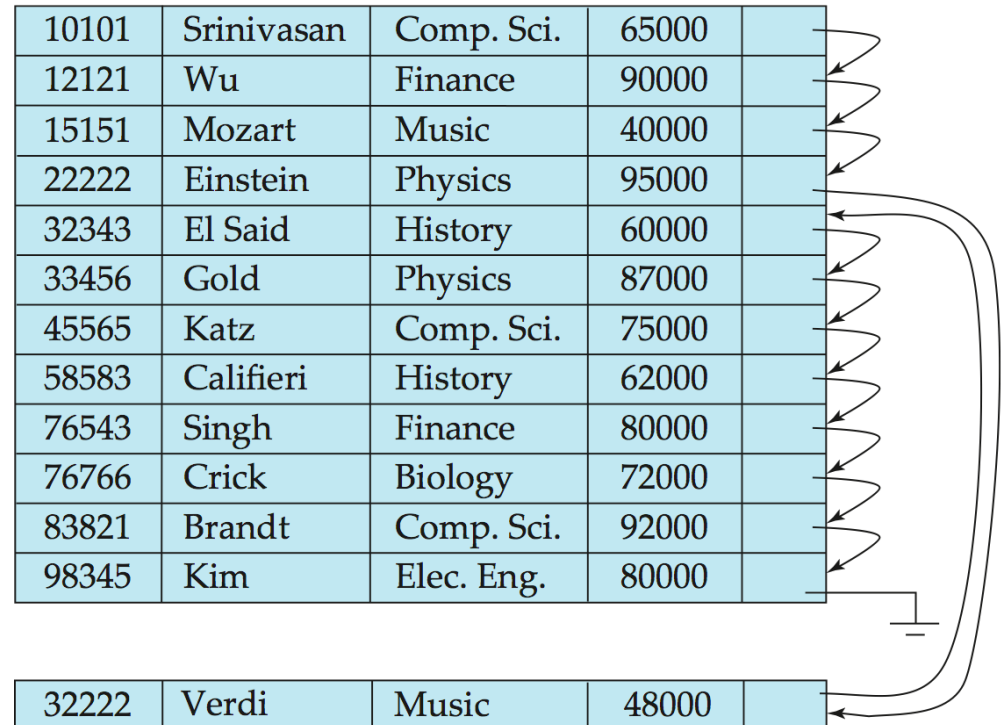
- ❑ Suitable for applications that require sequential processing of the entire file
- ❑ The records in the file are ordered by a **search-key** (ID)
- ❑ For fast retrieval records are chained

|       |            |            |       |  |  |
|-------|------------|------------|-------|--|--|
| 10101 | Srinivasan | Comp. Sci. | 65000 |  |  |
| 12121 | Wu         | Finance    | 90000 |  |  |
| 15151 | Mozart     | Music      | 40000 |  |  |
| 22222 | Einstein   | Physics    | 95000 |  |  |
| 32343 | El Said    | History    | 60000 |  |  |
| 33456 | Gold       | Physics    | 87000 |  |  |
| 45565 | Katz       | Comp. Sci. | 75000 |  |  |
| 58583 | Califieri  | History    | 62000 |  |  |
| 76543 | Singh      | Finance    | 80000 |  |  |
| 76766 | Crick      | Biology    | 72000 |  |  |
| 83821 | Brandt     | Comp. Sci. | 92000 |  |  |
| 98345 | Kim        | Elec. Eng. | 80000 |  |  |



# Sequential File Organization (Cont.)

- ❑ Insertion –locate the position where the record is to be inserted
  - ❑ if there is free space insert there
  - ❑ if no free space, insert the record in an **overflow block**
  - ❑ In either case, pointer chain must be updated
- ❑ Need to reorganize the file from time to time to restore sequential order







# Deletion – pointer chains

- ❑ Store the address of the first deleted record in the file header.
- ❑ Use this first record to store the address of the second deleted record, and so on
- ❑ Can think of these stored addresses as **pointers** since they “point” to the location of a record.
- ❑ More space efficient representation: reuse space for normal attributes of free records to store pointers. (No pointers stored in in-use records.)

|           |       |            |            |       |
|-----------|-------|------------|------------|-------|
| header    |       |            |            |       |
| record 0  | 10101 | Srinivasan | Comp. Sci. | 65000 |
| record 1  |       |            |            |       |
| record 2  | 15151 | Mozart     | Music      | 40000 |
| record 3  | 22222 | Einstein   | Physics    | 95000 |
| record 4  |       |            |            |       |
| record 5  | 33456 | Gold       | Physics    | 87000 |
| record 6  |       |            |            |       |
| record 7  | 58583 | Califieri  | History    | 62000 |
| record 8  | 76543 | Singh      | Finance    | 80000 |
| record 9  | 76766 | Crick      | Biology    | 72000 |
| record 10 | 83821 | Brandt     | Comp. Sci. | 92000 |
| record 11 | 98345 | Kim        | Elec. Eng. | 80000 |



# Multitable Clustering File Organization

Store several relations in one file using a **multitable clustering** file organization

*department*

| <i>dept_name</i> | <i>building</i> | <i>budget</i> |
|------------------|-----------------|---------------|
| Comp. Sci.       | Taylor          | 100000        |
| Physics          | Watson          | 70000         |

*instructor*

| <i>ID</i> | <i>name</i> | <i>dept_name</i> | <i>salary</i> |
|-----------|-------------|------------------|---------------|
| 10101     | Srinivasan  | Comp. Sci.       | 65000         |
| 33456     | Gold        | Physics          | 87000         |
| 45565     | Katz        | Comp. Sci.       | 75000         |
| 83821     | Brandt      | Comp. Sci.       | 92000         |

multitable clustering  
of *department* and  
*instructor*

|            |            |        |
|------------|------------|--------|
| Comp. Sci. | Taylor     | 100000 |
| 45564      | Katz       | 75000  |
| 10101      | Srinivasan | 65000  |
| 83821      | Brandt     | 92000  |
| Physics    | Watson     | 70000  |
| 33456      | Gold       | 87000  |



# Multitable Clustering File Organization (cont.)

- ❑ good for queries involving *department* ⋈ *instructor*, and for queries involving one single department and its instructors
- ❑ bad for queries involving only *department*
- ❑ results in variable size records
- ❑ Can add pointer chains to link records of a particular relation

|            |            |        |  |
|------------|------------|--------|--|
| Comp. Sci. | Taylor     | 100000 |  |
| 45564      | Katz       | 75000  |  |
| 10101      | Srinivasan | 65000  |  |
| 83821      | Brandt     | 92000  |  |
| Physics    | Watson     | 70000  |  |
| 33456      | Gold       | 87000  |  |



# Storage Access

- A database file is partitioned into fixed-length storage units called **blocks**.
  - Blocks are units of both storage allocation and data transfer.
- Database system seeks to
  - minimize the number of block transfers between the disk and memory.
  - We can reduce the number of disk accesses by
    - ▶ keeping as many blocks as possible in main memory.
- **Buffer** – portion of main memory available to store copies of disk blocks.
- **Buffer manager** – subsystem responsible for allocating buffer space in main memory
  - Not the same as the virtual memory manager of the OS



# End of Chapter 10

**Edited by Radhika Sukapuram**

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use