



- ① Conceptual Design (ER Diagram)
- 2) Logical Design (eg Relational Model)
- 3) Physical Design

Topic 4 – Part 1 (File Organization)

Chapter 13, ①④

Indexing & Hashing

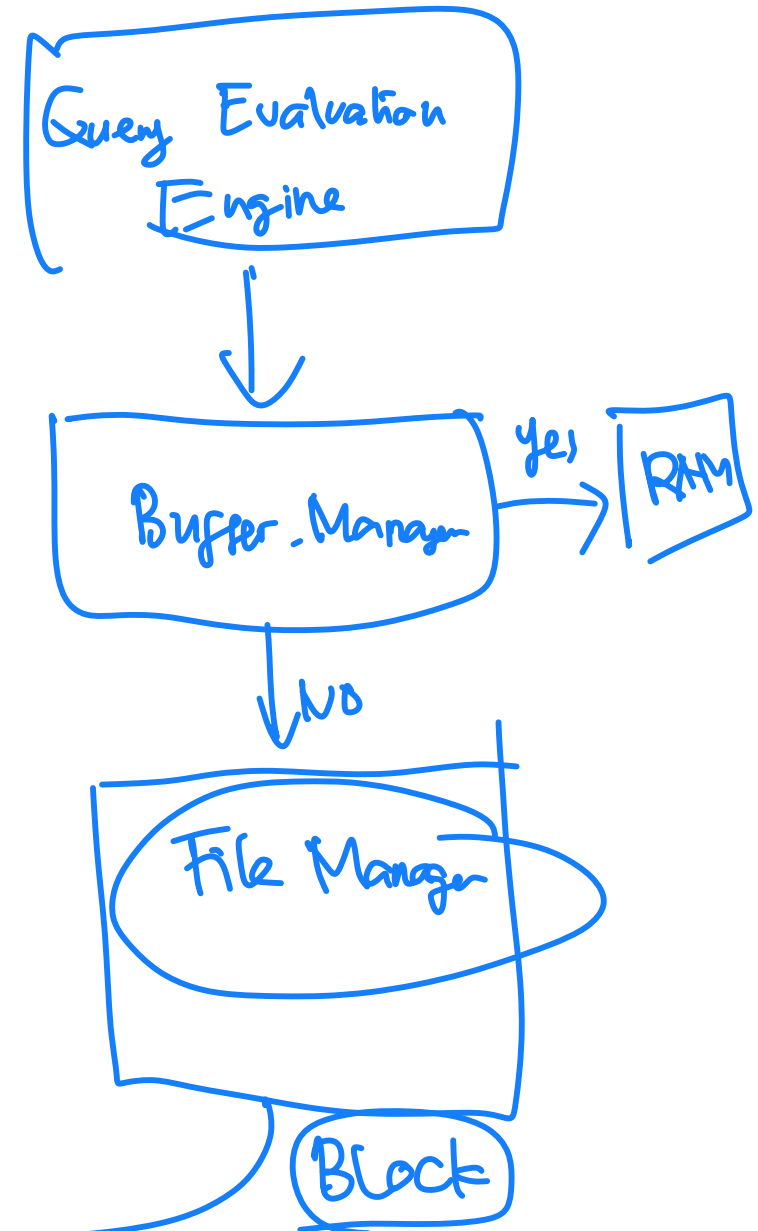
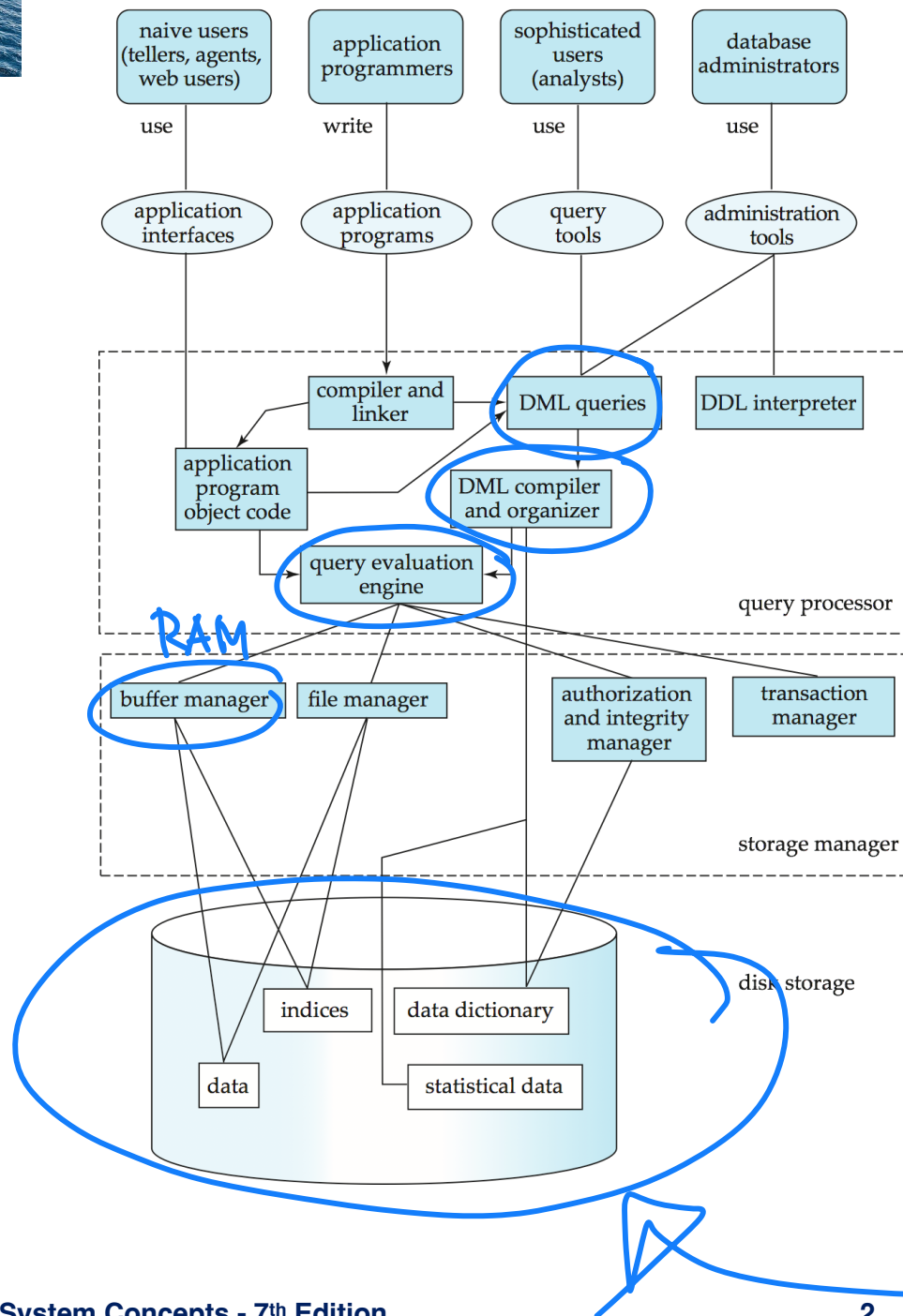
Database System Concepts, 7th Ed.

©Silberschatz, Korth and Sudarshan

See www.db-book.com for conditions on re-use

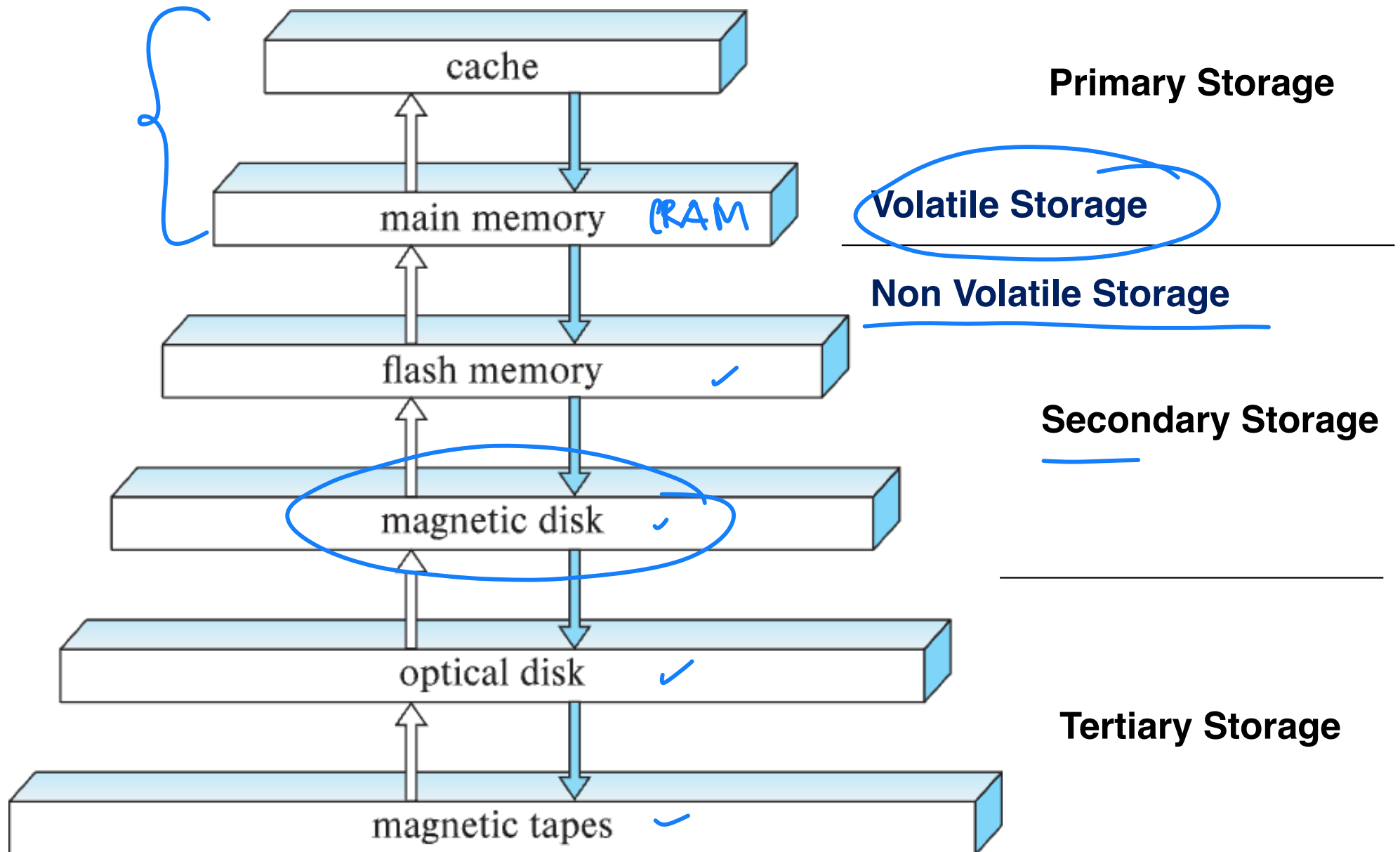


Database System Internals





Storage Hierarchy





Magnetic Hard Disk Mechanism

sector = block = 4KB

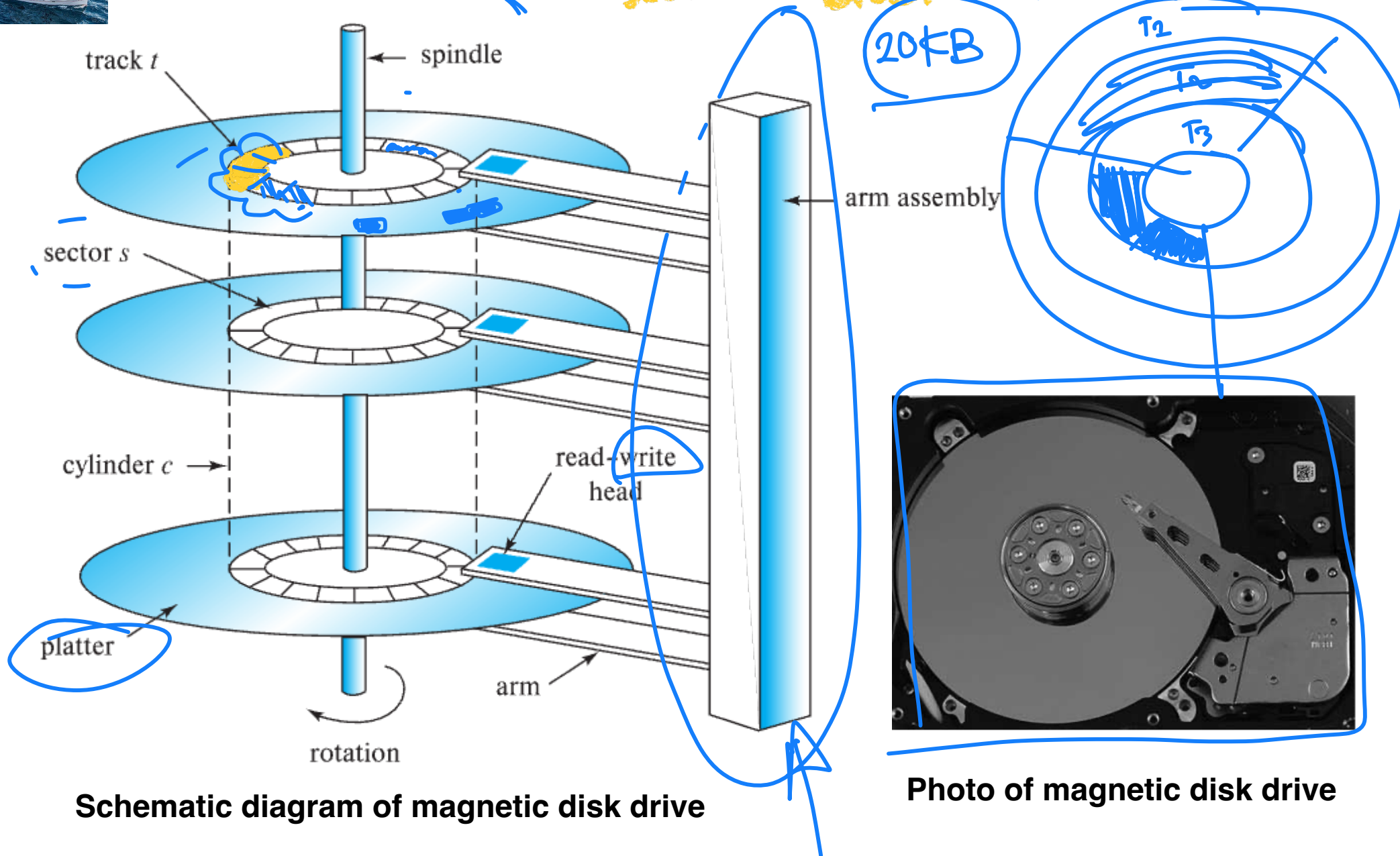


Photo of magnetic disk drive



Magnetic Hard Disk (Cont.)

20 KBytes

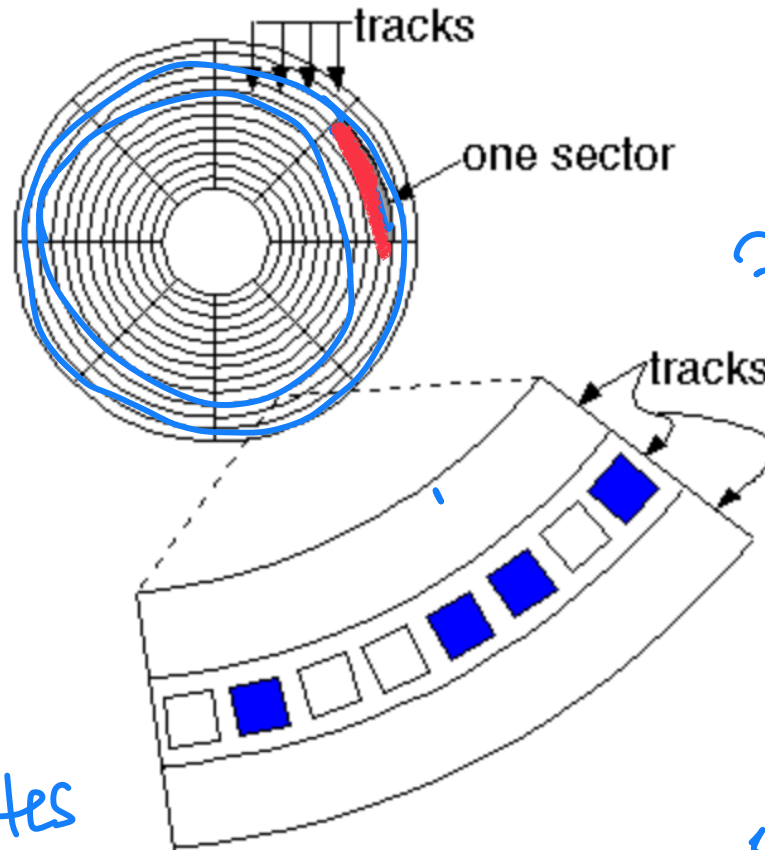
Sector is the smallest physical storage unit on a disk

512 bytes

4 Kbytes

⇒ 4096 bytes

1 KB = 1024 bytes



⇒ block ⇒ 4 KBytes

Logical unit for storage allocation & retrieval

4 V0

Tracks and Sectors

Tracks are concentric circles that are broken up into storage units called "sectors," typically 4,096 bytes long. The sector is the smallest unit that can be read or written. Tracks are only 75 nanometers wide today, and bit density is greater than one terabit per square inch. See areal density.

Source: <https://www.pcmag.com/encyclopedia/term/magnetic-disk>



Performance Measures of Disks

I/O operation

- **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.
 - Average seek time is 1/2 the worst case seek time.
 - Would be 1/3 if all tracks had the same number of sectors, and we ignore the time to start and stop arm movement
 - 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.)
 - Average latency is 1/2 of the above latency.
 - Overall latency is 5 to 20 msec depending on disk model
- **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.
 - 25 to 200 MB per second max rate, lower for inner tracks



Performance Measures (Cont.)

4 KB

- **Disk block** is a logical unit for storage allocation and retrieval (data transfer)

- 4 to 16 kilobytes typically

- Smaller blocks: more transfers from disk
- Larger blocks: more space wasted due to partially filled blocks

1 KB

20 KB

5 I/O

2 I/O

- **Sequential access pattern**

- Successive requests are for successive disk blocks
- Disk seek required only for first block

- **Random access pattern**

- Successive requests are for blocks that can be anywhere on disk
- Each access requires a seek
- Transfer rates are low since a lot of time is wasted in seeks

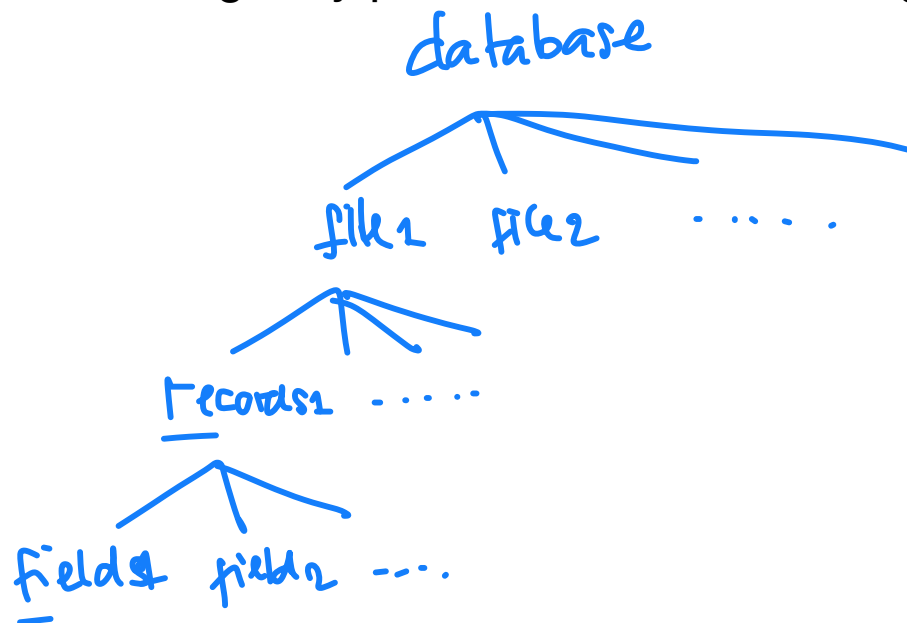
- **I/O operations per second (IOPS)**

- Number of random block reads that a disk can support per second
- 50 to 200 IOPS on current generation magnetic disks



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.
- Each file is also logically partitioned into fixed-length storage units called block.





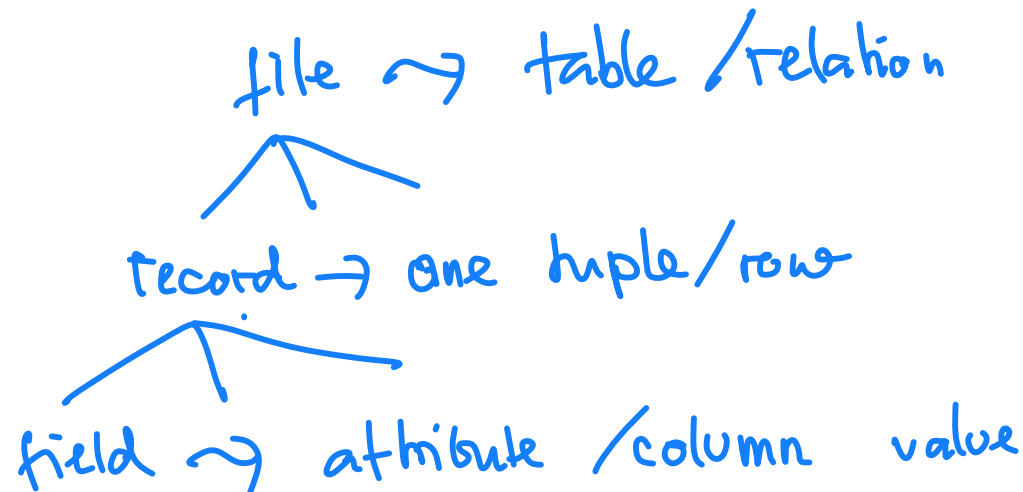
File Organization (Cont.)

- One approach of organizing data:
 - Assume record size is fixed
 - Each file has records of one particular type only
 - Different files are used for different relations

This case is easiest to implement; will consider variable length records later

- We assume that records are smaller than a disk block

4KB
 $2 \ll 4$
record size





Fixed-Length Records

■ Simple approach:

- Store record i starting from byte $n * (i - 1)$, where n is the size of each record.

$$100 * (40 - 1) = 3900$$

- Record access is simple but records may cross blocks

- Modification: do not allow records to cross block boundaries

Block size : 4KB

= 4096 Bytes

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
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

instructor (ID, name, dept. name, salary)
100 bytes

$$4096 / 100 = 40 \text{ records}$$

96 bytes

i: 1 = 0	
i: 2 = 100	
i: 3 = 200	
i: 40 = 3900	

i = 41

96

4 byte 10 another block



Fixed-Length Records

- Deletion of record i : alternatives:
 - move records $i + 1, \dots, n$ to $i, \dots, n - 1$
 - move record n to i
 - do not move records, but link all free records on a *free list*

Record 3 deleted

i	i''			
record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
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

$O(n)$



Fixed-Length Records

- Deletion of record i : alternatives:
 - move records $i + 1, \dots, n$ to $i, \dots, n - 1$
 - **move record n to i**
 - do not move records, but link all free records on a *free list*

Record 3 deleted and replaced by record 11

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 11	98345	Kim	Elec. Eng.	80000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
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



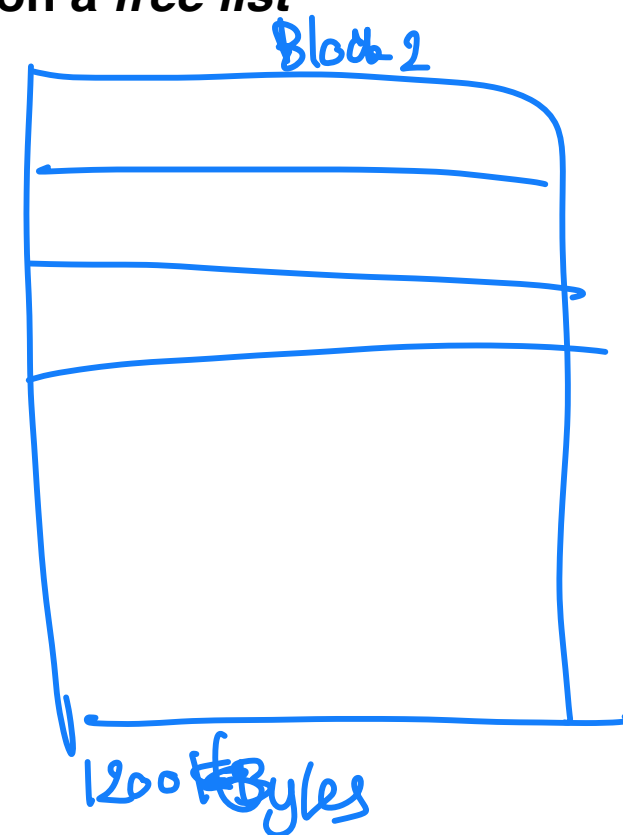
Fixed-Length Records

- Deletion of record i : alternatives:
 - move records $i + 1, \dots, n$ to $i, \dots, n - 1$
 - move record n to i
 - **do not move records, but link all free records on a *free list***

Block 1

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

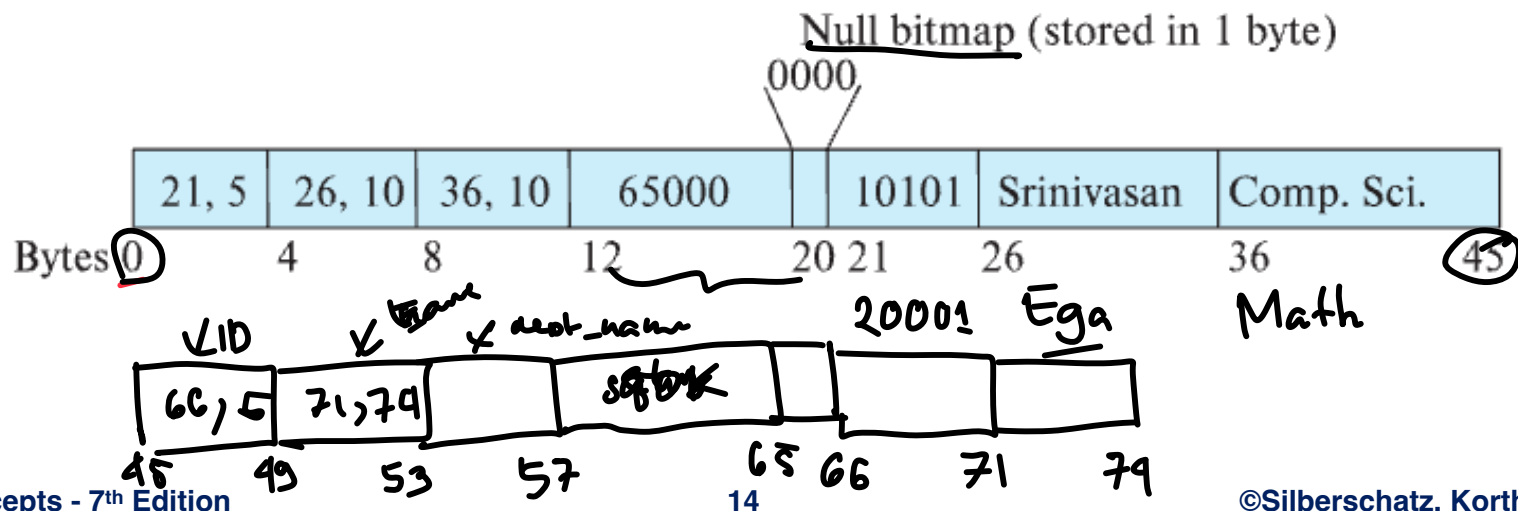
100 bytes





Variable-Length Records

- Variable-length records arise in database systems in several ways:
 - Storage of multiple record types in a file.
 - Record types that allow variable lengths for one or more fields such as strings (**varchar**)
 - Record types that allow repeating fields (used in some older data models).
- Attributes are stored in order
- Variable length attributes represented by fixed size (**offset, length**), with actual data stored after all fixed length attributes
- Null values represented by null-value bitmap



me
ardar (:)



```
create table instructor (  
  ID varchar(5),  
  name varchar(50),  
  dept_name varchar(25),  
  salary int,  
  primary key (ID),  
  foreign key (dept_name) references (department)  
)
```

4 char → 1 byte



Storing Large Objects

- Databases often store data that can be much larger than a disk block.
- E.g. blob/clob types
- Records must be smaller than pages
- Alternatives:
 - Store as files in file systems → OS
 - Store as files managed by database
 - Break into pieces and store in multiple tuples in separate relation
 - PostgreSQL TOAST

→ binary large object

images / videos

→ character large object

4KB



5MB

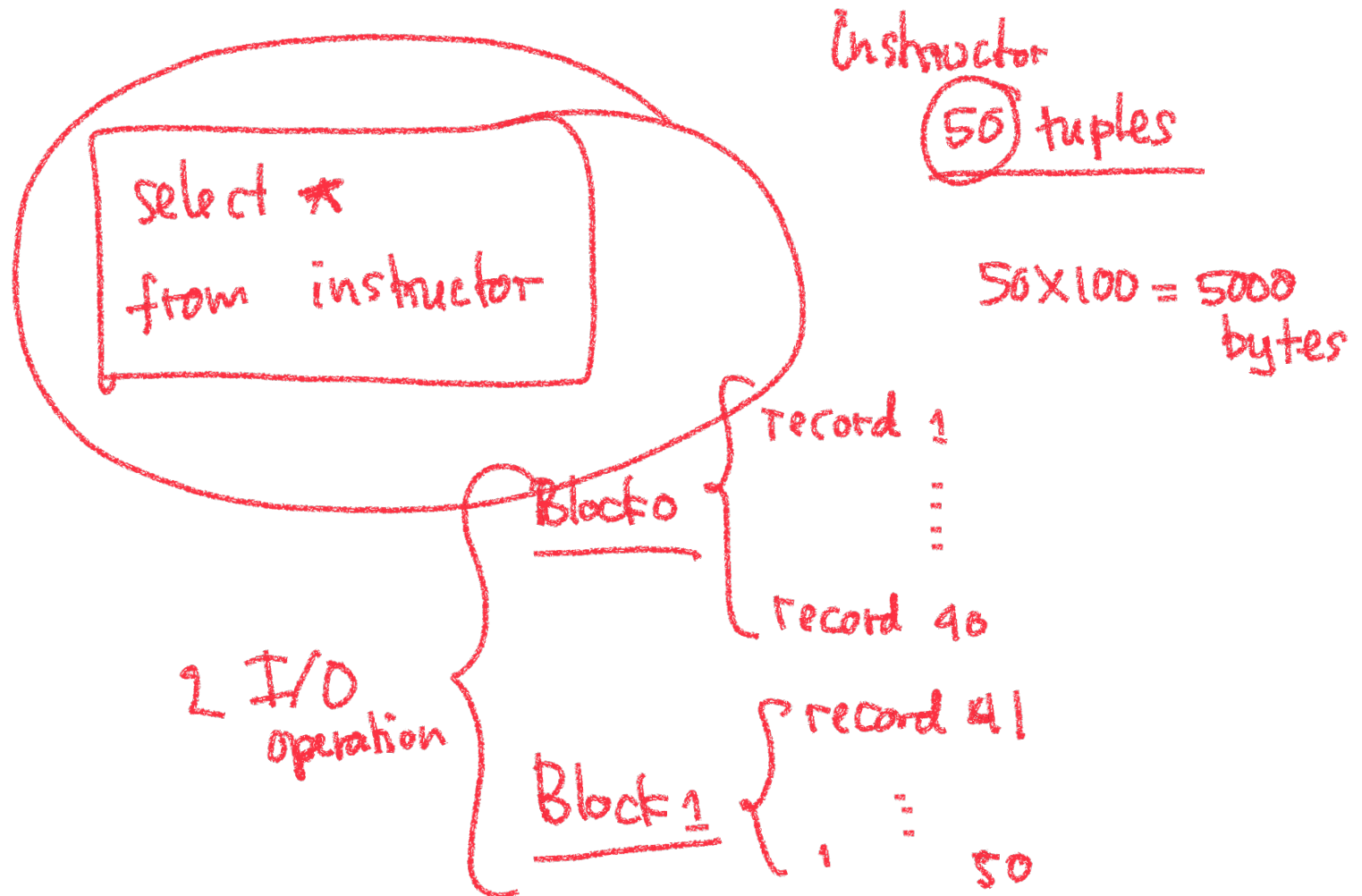
1 file ~> 1 relation



Topic 4 Poll 1

$$4096/100 = 40 \text{ records} \\ 96 \text{ bytes}$$

- Suppose a disk has a block size of 4 KB (4096 bytes). A file contains a database relation with 50 fixed-size records, each 100 bytes. How many I/O operations are required to retrieve all data from the relation?





Organization of Records in Files

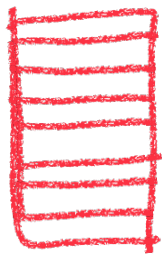
- **Heap** – record can be placed anywhere in the file where there is space
- **Sequential** – store records in sequential order, based on the value of the search key of each record
- In a **multitable clustering file organization** records of several different relations can be stored in the same file
 - Motivation: store related records on the same block to minimize I/O
- **B+-tree file organization**
 - Ordered storage even with inserts/deletes
 - More on this in Chapter 14 (Topic 4 Part 2)
- **Hashing** – a hash function computed on search key; the result specifies in which block of the file the record should be placed
 - More on this in Chapter 14 (Topic 4 Part 2)



Heap File Organization

- Records can be placed anywhere in the file where there is free space
- Records usually do not move once allocated
- Important to be able to efficiently find free space within file
- Free-space map** → **data structure that tracks the fraction of the free space in each block**

- Array with 1 entry per block. Each entry is a few bits to a byte, and records fraction of block that is free
- In example below, 3 bits per block, value divided by 8 indicates fraction of block that is free



- Can have :

4	2	1	4	7	3	6	5	1	2	0	1	1	0	5	6
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---
- In example below, each entry stores maximum from 4 entries of first-level free-space map

1st level

2nd level

4	7	2	6
---	---	---	---

Block is $\frac{6}{8}$
6 blocks

- Free space map written to disk periodically, OK to have wrong (old) values for some entries (will be detected and fixed)



Sequential File Organization

- Suitable for applications that require sequential processing of the entire file → *relation / table*
- The records in the file are ordered by a search-key

IPV

0

1

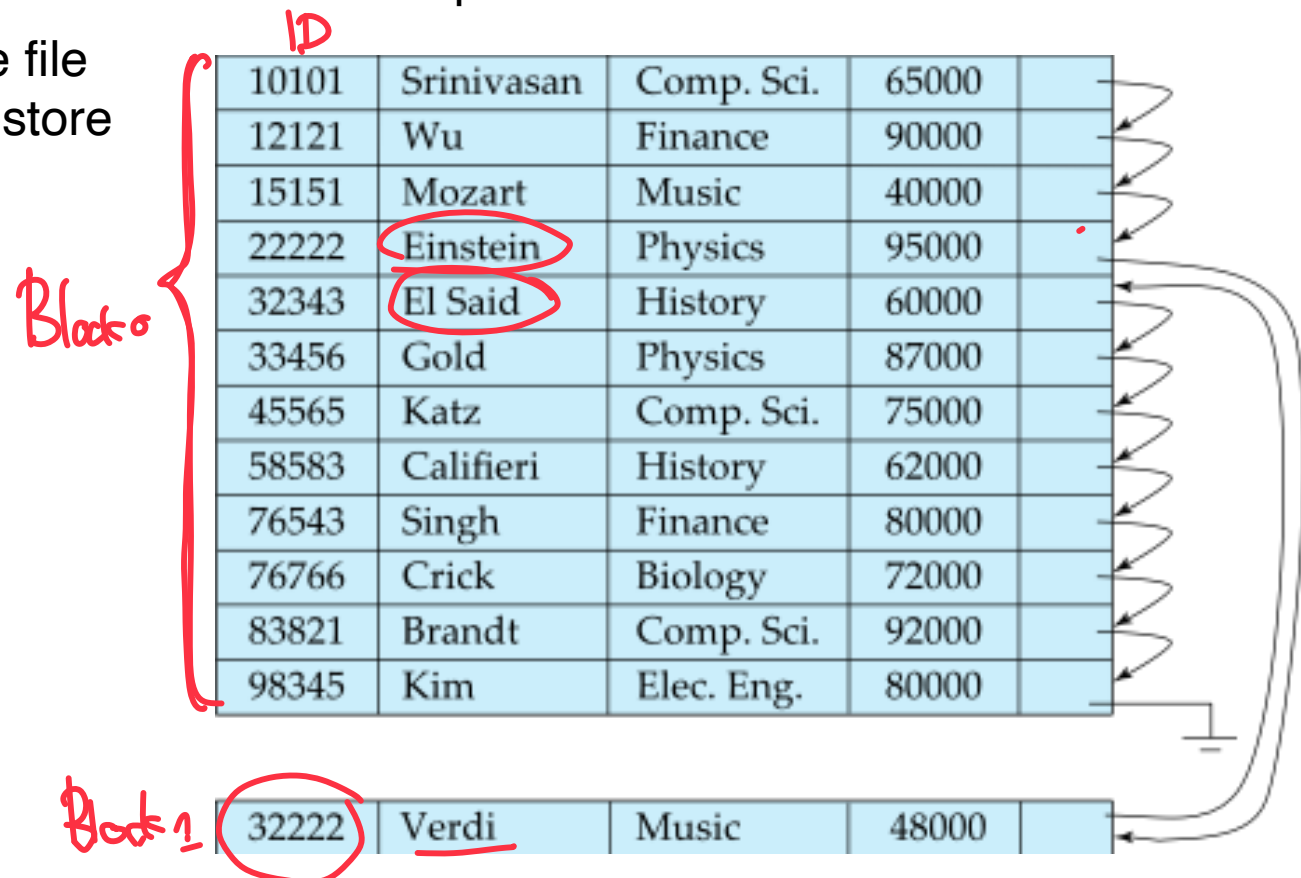
10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califieri	History	62000	
76543	Singh	Finance	80000	
76766	Crick	Biology	72000	
83821	Brandt	Comp.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

instructor



Sequential File Organization (Cont.)

- Deletion – use pointer chains
- 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





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

Select *

from department,

instructor

where department.

dept_name

=

instructor.

dept_name

multitable clustering
of *department* and
instructor

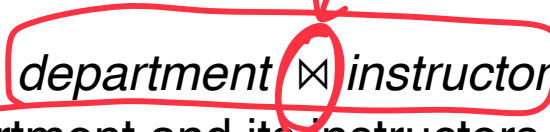
Comp. Sci.	Taylor	100000	dept
10101	Srinivasan	Comp. Sci.	65000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000
Physics	Watson	70000	
33456	Gold	Physics	87000



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

join





Partitioning

- **Table partitioning:** Records in a relation can be partitioned into smaller relations that are stored separately
- E.g., *transaction* relation may be partitioned into *transaction_2018*, *transaction_2019*, etc. credit_card_usage
- Queries written on *transaction* must access records in all partitions
 - Unless query has a selection such as *year=2019*, in which case only one partition is needed
- Partitioning
 - Reduces costs of some operations such as free space management
 - Allows different partitions to be stored on different storage devices
 - E.g., *transaction* partition for current year on SSD, for older years on magnetic disk



Data Dictionary Storage

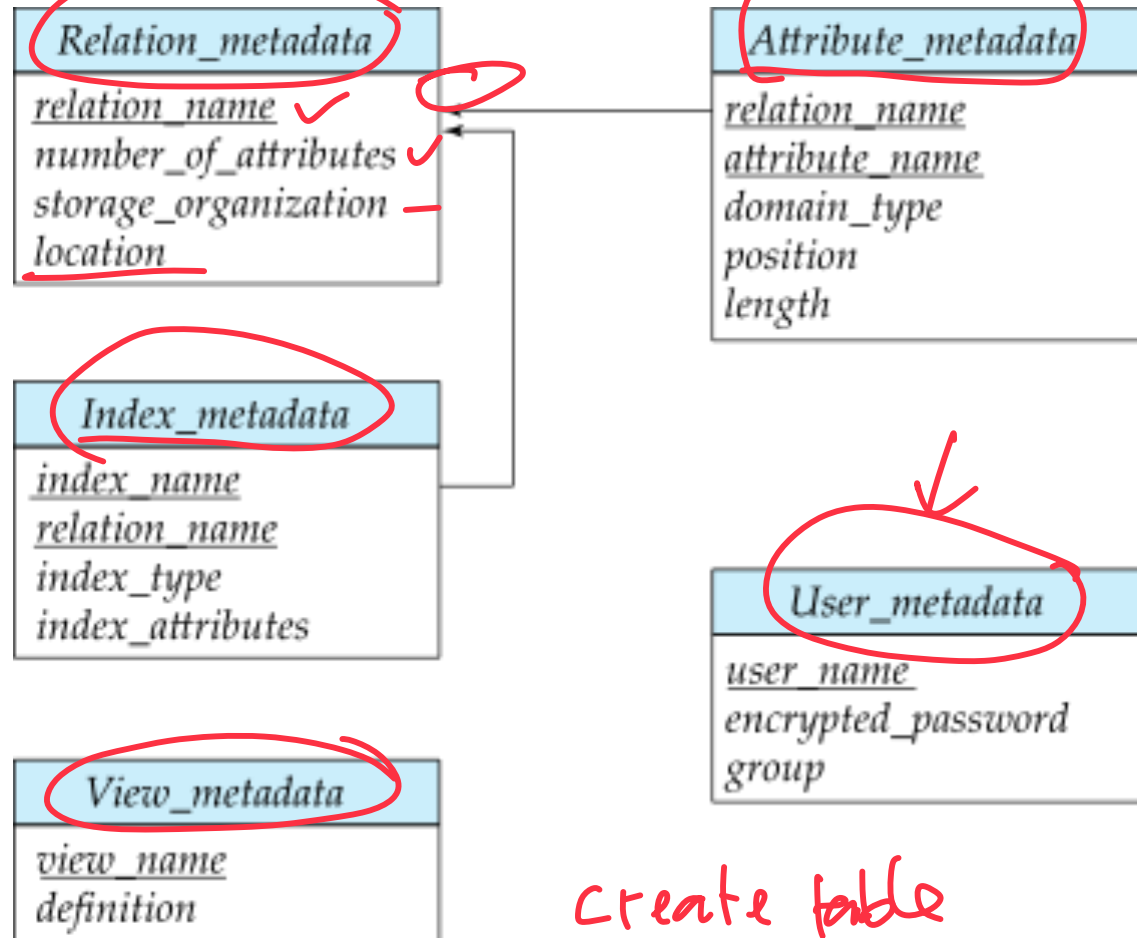
The **Data dictionary** (also called **system catalog**) stores **metadata**; that is, data about data, such as

- Information about relations
 - names of relations
 - names, types and lengths of attributes of each relation
 - names and definitions of views
 - integrity constraints
- User and accounting information, including passwords
- Statistical and descriptive data
 - number of tuples in each relation
- Physical file organization information
 - How relation is stored (sequential/hash/...)
 - Physical location of relation
- Information about indices (Chapter 14)



Relational Representation of System Metadata

- Relational representation on disk
- Specialized data structures designed for efficient access, in memory





End of Topic 4 Part 1