



## — Chapter 13: Data Storage Structures

Database System Concepts, 7<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use



## 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.
- One approach
  - ① Assume record size is fixed
  - ② Each file has records of one particular type only
  - ③ 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

```
type instructor = record
    ID varchar (5);
    name varchar(20);
    dept_name varchar (20);
    salary numeric (8,2);
end
```



## Fixed-Length Records

- Simple approach:

- Store record  $i$  starting from byte  $n * (i - 1)$ , where  $n$  is the size of each record.
- Record access is simple but records may cross disk blocks
  - Modification: do not allow records to cross block boundaries
  - Disk block** is a logical unit for storage allocation and retrieval
    - 4 to 16 kilobytes typically
    - Smaller blocks: more transfers from disk
    - Larger blocks: more space wasted due to partially filled blocks

|           |       |            |            |       |
|-----------|-------|------------|------------|-------|
| 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 |



## 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

➤

|           |       |            |            |       |
|-----------|-------|------------|------------|-------|
| 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 |



# 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 |

to move the final record into space occupied by the deleted record



# Fixed-Length Records

- Deletion of record  $i$ : alternatives:
  - move records  $i + 1, \dots, n$  to  $i, \dots, n - 1$
  - move record  $n$  to  $i$

linked list  
://.

do not move records, but link all free records on a *free list*

change the header pointer  
to point to next available record.  
If no space is available  
⇒ add the new record  
to the end of the file

| 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 |

there is the address  
of the first record whose contents are deleted.

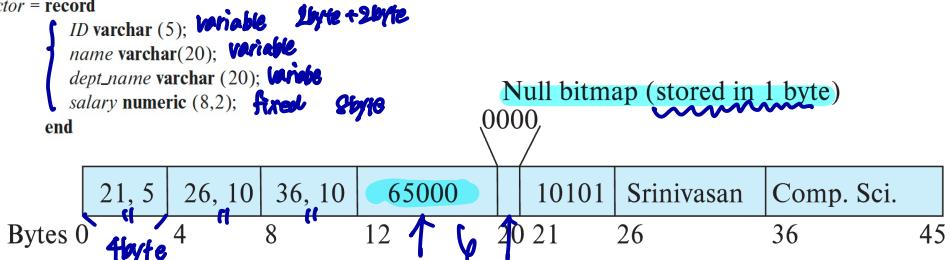


# Variable-Length Records

⇒ いじくのいじき

- 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

`type instructor = record`



Database System Concepts - 7<sup>th</sup> Edition

13.7

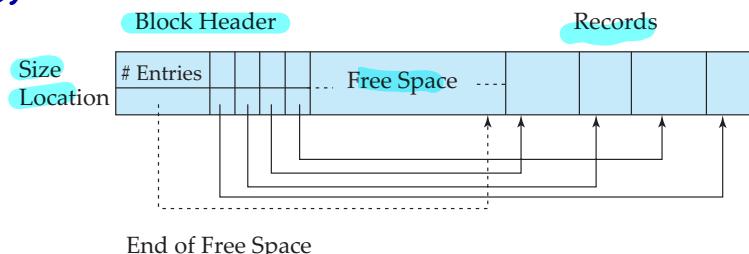
©Silberschatz, Korth and Sudarshan

and 13~19 bytes were ignored



## Variable-Length Records: Slotted Page Structure

⇒ record いじき



End of Free Space

- **Slotted page (i.e., block) header** contains:
  - number of record entries
  - end of free space in the block
  - location and size of each record
- If a record is deleted, the space that it occupies is freed, and its entry is set to deleted (its size is set to -1, for example).
- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record — instead they should point to the entry for the record in header.

Database System Concepts - 7<sup>th</sup> Edition

13.8

©Silberschatz, Korth and Sudarshan



## Storing Large Objects

- e.g., blob/clob types (binary large objects / character large objects)
- Records must be smaller than pages  
*(block)*
- Alternatives:
  - Store as files in file systems
  - Store as files managed by database
  - Break into pieces and store in multiple tuples in separate relation
    - PostgreSQL TOAST



## 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 *it is ordered.*
  - More on this in Chapter 14
- **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



# 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**
    - 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
- |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4 | 2 | 1 | 4 | 7 | 3 | 6 | 5 | 1 | 2 | 0 | 1 | 1 | 0 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
- e.g., If a block size is 512 bytes and there is 8 partitions, then the size of each partition is  $512/8 = 64$  bytes.  
1<sup>st</sup> entry indicates at least 4/8<sup>th</sup> of the space (i.e.,  $4 \times 64$  bytes) in the block is free.
  - Can have second-level free-space map
  - In example below, each entry stores the maximum from 4 entries of first-level free-space map
- |   |   |   |   |
|---|---|---|---|
| 4 | 7 | 2 | 6 |
|---|---|---|---|



# Sequential File Organization

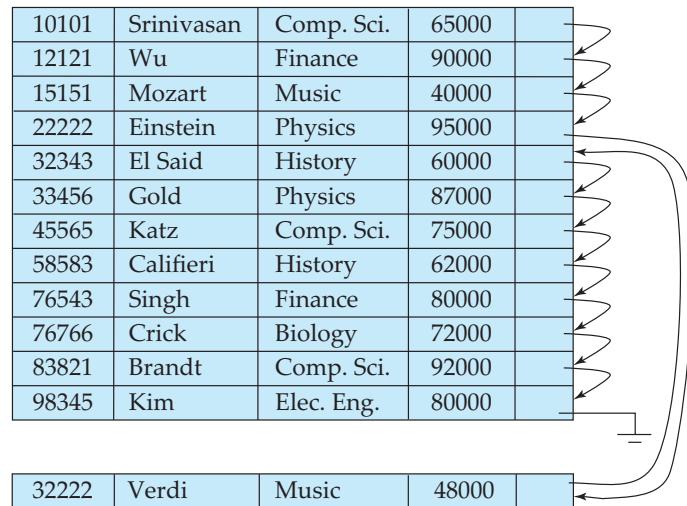
- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a search-key

|       |            |            |       |  |
|-------|------------|------------|-------|--|
| 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.)

- 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*

| dept_name  | building | budget |
|------------|----------|--------|
| Comp. Sci. | Taylor   | 100000 |
| Physics    | Watson   | 70000  |

*instructor*

| ID    | name       | dept_name  | salary |
|-------|------------|------------|--------|
| 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     |       |
| 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*  $\bowtie$  *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
- Multitable clustering is supported by the Oracle database system



## 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.
- 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

*F2P4*  
2019 year is SSD interface  
⇒ cheap storage.



# Data Dictionary Storage

## Data of Data

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

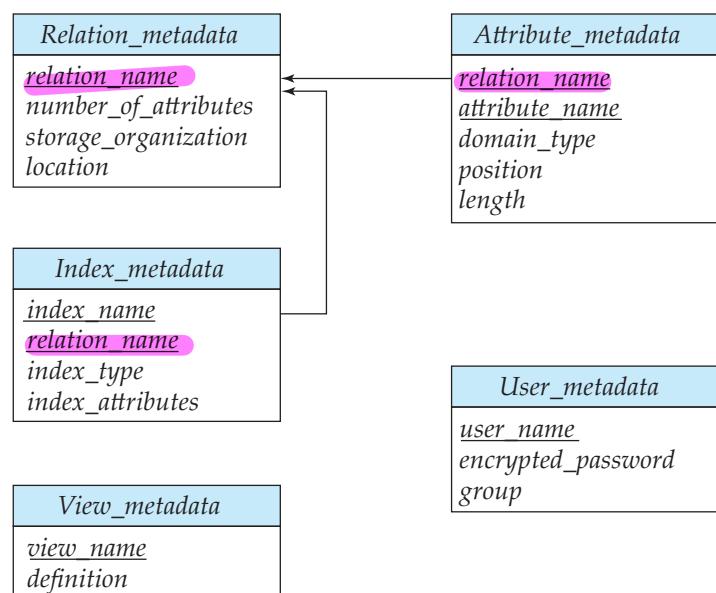
- Information about relations
  - names of relations
  - names, types and lengths of attributes of each relation
  - names and definitions of views
  - integrity constraints
- User and accounting information, including passwords
- Statistical and descriptive data
  - number of tuples in each relation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



# Relational Representation of System Metadata

## table 表의 구조

- Relational representation on disk
- Specialized data structures designed for efficient access, in memory
- Whenever the database system needs to retrieve records from a relation, it must first consult the Relation metadata relation to find the location and storage organization of the relation, and then fetch records using this information.

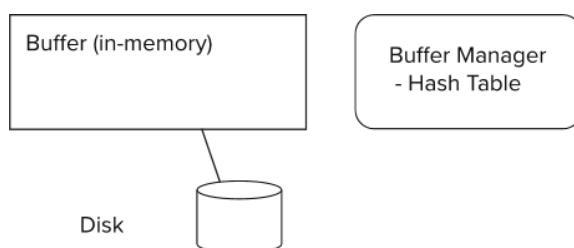




# Storage Access

- 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. *Buffer manager is a table*.

Disk I/O Flow



↓ Disk is block manager.

## Buffer Manager

\* physical level    user level.

- Programs call on the buffer manager when they need a block from disk.
  - If the block is already in the buffer, buffer manager returns the address of the block in main memory
  - If the block is not in the buffer, the buffer manager
    - Allocates space in the buffer for the block
      - Replacing (throwing out) some other block, if required, to make space for the new block.
      - Replaced block written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk.
    - Reads the block from the disk to the buffer, and returns the address of the block in main memory to requester.

\* Block of data  
descriptor of block.



# Buffer Manager

- Buffer replacement strategy (details coming up!)
- ➡ Pinned block: memory block that is not allowed to be written back to disk
  - Pin done (before) reading/writing data from a block ← the buffer manager never evicts pinned block
  - Unpin done when read /write is complete ← allowing the block to be evicted
- Keep a pin count, buffer block can be evicted only if pin count = 0
- Shared and exclusive locks on buffer
  - Needed to prevent concurrent operations from reading page contents as they are moved/reorganized, and to ensure only one move/reorganize at a time
  - Readers get shared lock, updates to a block require exclusive lock
  - Locking rules: DB process can buffer block if & only if it has lock on it
- Only one process can get exclusive lock at a time
- Shared lock cannot be concurrently with exclusive lock
- Multiple processes may be given shared lock concurrently

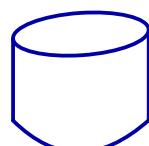


## Buffer-Replacement Policies

가장 먼저 액세스한 것을 만날 때

- X Most operating systems replace the block **least recently used (LRU strategy)**
  - Idea behind LRU – use past pattern of block references as a predictor of future references
  - LRU can be bad for some queries
- Queries have well-defined access patterns (such as sequential scans), and a database system can use the information in a user's query to predict future references
- Mixed strategy with hints (on replacement strategy provided by the query optimizer is preferable)
- Example of bad access pattern for LRU: when computing the join of 2 relations r and s by a nested loops

for each tuple  $tr$  of  $r$  do  
 for each tuple  $ts$  of  $s$  do  
 if the tuples  $tr$  and  $ts$  match ...





## Buffer-Replacement Policies (Cont.)

- **Toss-immediate** strategy – frees the space occupied by a block as soon as **the final tuple of that block** has been processed
- **Most recently used (MRU) strategy** – system must pin the block currently being processed. After the final tuple of that block has been processed, the block is unpinned, and it becomes the most recently used block. *가장 최근에 사용된 블록은 버퍼에서 해제된다.*
- Buffer manager can use statistical information regarding the probability that a request will reference a particular relation
  - E.g., the data dictionary is frequently accessed. Heuristic: keep data-dictionary blocks in main memory buffer
- Operating system or buffer manager may reorder writes
  - Can lead to corruption of data structures on disk *in the event of a system crash*
    - e.g., linked list of blocks with missing block on disk
    - File systems perform consistency check to detect such situations
  - Careful ordering of writes can avoid many such problems



## Optimization of Disk Block Access (Cont.)

- Buffer managers support **forced output** of blocks for the purpose of recovery (more in Chapter 19)
- **Nonvolatile write buffers** speed up disk writes by writing blocks to a non-volatile RAM or flash buffer immediately
  - *Writes can be reordered to minimize disk arm movement*
- **Log disk** – a disk devoted to writing a sequential log of block updates
  - Used exactly like nonvolatile RAM *로그 디스크는 이진 형식으로 데이터를 차례대로 기록하는 방식*
    - Write to log disk is very fast since no seeks are required
- **Journaling file systems** write data in-order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data



# Column-Oriented Storage

• 한 행 / 행

- Also known as **columnar representation**
- Store each attribute of a relation separately
- Example *instructor 가 한 행에 걸친 열을 갖는 columnar 형태.*

|       |            |            |       |
|-------|------------|------------|-------|
| 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 |



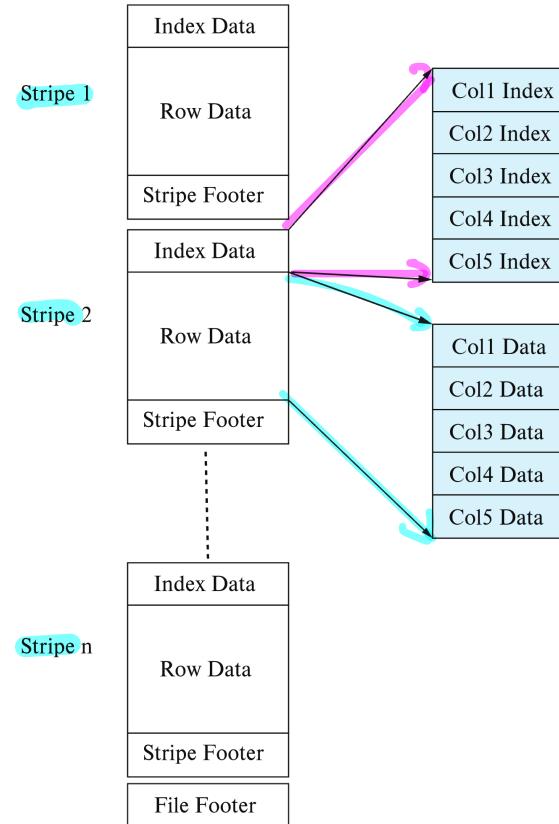
# Columnar Representation

- *\* 한 행 푸시 ⇒ column.*
- Benefits:
  - Reduced IO if only some attributes are accessed
  - Improved CPU cache performance
  - Improved compression
  - **Vector processing** on modern CPU architectures
- Drawbacks
  - ⌚ Cost of tuple reconstruction from columnar representation
  - ⌚ Cost of tuple deletion and update
  - ⌚ Cost of decompression
- Columnar representation found to be more efficient for decision support than row-oriented representation  
*• 예제.*
- Traditional row-oriented representation preferable for transaction processing
- Some databases support both representations
  - Called **hybrid row/column stores**



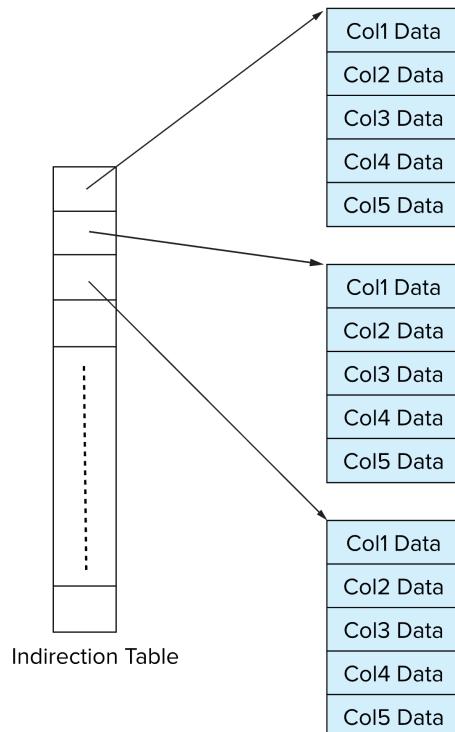
# Columnar File Representation

- ORC(Optimized Row Columnar) and Parquet: file formats with columnar storage inside file
- Very popular for big-data applications
- Orc file format shown on right:



# Storage Organization in Main-Memory Databases

- Can store records directly in memory without a buffer manager
- Column-oriented storage can be used in-memory for decision support applications
  - Compression reduces memory requirement





## End of Chapter 13