# **Storing Data**

CSC443, Winter 2019 Sayyed Nezhadi

#### **Computer Memory Hierarchy**



# Levels of Memory Hierarchy

- <u>Cache:</u> Volatile, fastest accessible memory (a few nanoseconds), expensive, typical size in megabytes
- Main Memory: Volatile, most used memory, fast access (around 100 nanoseconds), typical size in order of 10GB
- <u>Secondary Storage</u>: Hard drive or SSD, nonvolatile, very large storage capacity (terabytes), much slower than memory (a few milliseconds for hard drives).
- <u>Tertiary Storage</u>: tape drives or robotic disk arrays, nonvolatile, mostly used for data backup, huge storage capacity (petabytes), very slow access (seconds or minutes).

► *In this course, we are mainly concerned with the main memory and the secondary storage* 

#### Move Data to & from Secondary Storage

- The <u>mechanical design</u> of the hard drive naturally leads itself to <u>superior sequential access</u>.
- A hard drive consists of an (array of) rotating disk.
  Each disk has a head driven by a mechanical arm which performs the bitwise read and write. If a sequential data block is accessed, all bits can be read or written during a single rotation.
- <u>Random access</u> of the hard drive data may require <u>multiple rotations</u> for the head to reach all the regions on the disk.

### Data Transfer Rate

• <u>Seek Time</u>: measures *average* time it takes the head to travel to the region of interest on disk. Typically, the seek time is a few milliseconds.

Roational speed (rpm) Seek time (ms)

15,000	2
10,000	3
7,200	4.16

 <u>Data Transfer Rate</u>: is the steady state speed the head can read or write data between the disk buffer and the disk. The data transfer rate can reach 1 terabytes / second.

# **Block Size**

- The data transfer speed depends on:
  - contiguousness of the accessed physical locations on the disk
  - the size of the accessed data per request (block)



When data is requested in large enough blocks, we see a drastic improvement in data transfer rate.

# **Buffer Manager**

- Another technique to further combat disk latency.
- All disk I/O issued by the database are done in units of <u>pages</u>.
- A section of the memory is used as the **<u>buffer pool</u>**.
- The buffer pool is partitioned into <u>frames</u>. Each frame is the size of a page.
- The buffer manager performs lazy disk I/O:
  - Disk read and write are first performed as frame read and write in the buffer pool
  - Data is transferred between the buffer pool and disk only when a requested page is not cached by the buffer pool, or new pages need to be cached but the buffer pool is full.

### **Relational Data on Disk**

# Heap File for Pages

- A <u>heap file</u> is a data structure that holds pages of data.
- Functionally, a heap file is similar to the main memory data structure <u>linked list</u>.
- It serves as a container of generic data items in an <u>unordered</u> fashion
- Similar to linked list, heap file:
  - offers efficient append of new data,
  - supports sequential scanning,
  - does not offer random access.

## **Properties of Heap File**

- <u>Data exists on secondary storage</u>: it's designed to be efficient for large data volume, with capacity only limited by the available secondary storage.
- <u>A heap file can span over multiple disks or machines:</u> heap file can use large addressing schemes to span across multiple disks or ever networks.
- Data is organized into pages.
- <u>Pages can be partially empty</u>.

#### Heap File as a Linked List



• <u>Problem:</u> finding free spaces is very expensive. To locate a page with free space (*page3*), we have to scan through the pages (*page1* and *page2*).

# **Directory-based Heap File**

 <u>Idea:</u> use one or more blocks to store the directory of all the page addresses and their respective free space.



► Since page1 and page2 are both full, their free space entries n1 and n2 will be both zero



If there are too many pages that the directory itself needs to span over multiple pages, we may use a linked list to store the directory.

# **Record Serialization**

- How to convert tuple records into byte arrays to be stored on disk (inside pages)?
- <u>Serializing fixed length records</u>: all records *must* be exactly the same length.
- Example:



\* sizeof(CHAR) + sizeof(INTEGER) + sizeof(DATETIME) 100

### **Record Serialization**

- <u>Serializing variable length records</u>: records have variable size. A popular method is to encode the <u>offsets of the field boundaries</u> in the header.
- Example: CREATE TABLE Person2 ( name VARCHAR(1024) NOT NULL, age INTEGER NOT NULL, birthdate DATETIME



 Having these boundary offsets provides a natural way of encoding NULL valued fields -- they can be viewed as fields with zero bytes:



➡ Here, we have assumed all value in a record are stored together. This is sometimes called <u>row-oriented storage</u>. Alternatively, all values of an attribute can be stored together in what is called <u>column-oriented storage</u>.

## Page Format

- How data is formatted in a each page?
- <u>Simplification</u>: Every page stores tuples with the same schema:
  - Every relation (table) occupies at least one page.
  - Page contains either fixed length tuples or variable length records.
- Important properties of page formats:
  - efficiently utilize disk space.
  - support scan, insert, delete and update with minimal disk I/O.
  - efficient way of assigning unique ID's to records (invariant during insert, delete and update )
  - each record fits into a page (will be relaxed later)

#### Page Format for Fixed Length Records

- M slots of equal size
- Each slot one record
- Last segment:
  - Value of M
  - M bits to indicate if the respective slot is occupied or free.



#### Page Format for Fixed Length Records

- <u>Record ID:</u> concat of <page ID, record offset>
- The record ID will never be affected by changes to the database.

INSERT	UPDATE	DELETE
1. Get <b>M</b>	Simply overwrite	Set the
2. Scan the M bits	the old record (at	corresponding
backward	the same location)	bit in the bit
3. All the bits are 1 ->		index to 0
no room in the page		
4. K-th bit is 0 -> write		
at offset <b>nk</b>		

Variable length records are <u>much more difficult</u> to maintain in a page for the following reasons:

- The record boundaries and field boundaries are dependent on the data content.
- When the data is modified, <u>record size may</u> <u>change</u>, resulting in physical relocation of the record (requires efficient methods to move records in a page or between pages)
- If record ID needs to be maintained, <u>the record ID</u> <u>assignment must be invariant</u> even if the record is to be moved.

One efficient design is to maintain a <u>page header</u> which tracks the free space and a <u>directory of records</u>.



- <u>Record ID</u>: concat of <page ID, directory index>
- The record ID is independent of changes to the database.

	INSERT		DELETE
1.	Check free space can hold the new record and the new entry in the directory	1.	If last entry in the directory, delete the entry and update free space pointer
2.	Write the record at the start of free space	2.	Else, make the offset = -1 in the entry (mark sit as deleted)
3.	Update directory (offset & size)		
4.	Update the free space pointer		

► Deletion leaves <u>holes</u> in the page. The space occupied by deleted records cannot be reclaimed unless we reorganize the page such that all space before the free space pointer are contiguously occupied. This is called <u>compacting</u> the page.

#### UPDATE

- 1. If new record smaller than the old record, update the record and the record size in the directory
- 2. Else, if enough free space, insert in the free space and update directory (offset & size)
- 3. Else, try to compact the page
- 4. If still not enough space, allocate a new page and write the new record in the new page. <u>But record ID needs to</u> <u>be changed</u>. We will make the original record as a <u>forwarding</u> record. The forwarding record only contains header information indicating that it is a forwarding record, and the record ID of the new record.