

Data Storage and Organization

COSC 404 – Database System Implementation



Storage and Organization Overview



A database system relies on the operating system to store data on storage devices.

Database performance depends on:

- Properties of storage devices
- How devices are used and accessed via the operating system

We will study techniques for storing and representing data.

Review: Memory Definitions

Temporary memory retains data only while the power is on.

- Also referred to as **volatile storage**.
- e.g. dynamic random-access memory (DRAM) (main memory)

Permanent memory stores data even after the power is off.

- Also referred to as **non-volatile storage** or **secondary storage**
- e.g. flash memory, SSD, hard drive, DVD, tape drives

Cache is faster memory used to store a subset of a larger, slower memory for performance.

- processor cache (Level 1 & 2), disk cache, network cache

Research Question

In-Memory Database



Question: Does an in-memory database need a secondary storage device for persistence?

A) Yes

B) No

Review: Sequential vs. Random Access

RAM, SSDs, and flash memory allow random access. *Random access* allows retrieval of any data location in any order.

Tape drives allow sequential access. *Sequential access* requires visiting all previous locations in sequential order to retrieve a given location.

- That is, you cannot skip ahead, but must go through the tape in order until you reach the desired location.

Review: Memory Sizes

Memory size is a measure of memory storage capacity.

- Memory size is measured in **bytes**.
 - Each byte contains 8 **bits** - a bit is either a 0 or a 1.
 - A byte can store one character of text.
- Large memory sizes are measured in:

▪ kilobytes (KBs)	= 10^3	= 1,000 bytes
▪ Kibibyte (KiBs)	= 2^{10}	= 1,024 bytes
▪ megabytes (MBs)	= 10^6	= 1,000,000 bytes
▪ mebibyte (MiBs)	= 2^{20}	= 1,048,576 bytes
▪ gigabytes (GBs)	= 10^9	= 1,000,000,000 bytes
▪ gibibytes (GiBs)	= 2^{30}	= 1,073,741,824 bytes
▪ terabytes (TBs)	= 10^{12}	= 1,000,000,000,000 bytes
▪ tebibytes (TiBs)	= 2^{40}	= 1,099,511,627,776 bytes

Transfer Size, Latency, and Bandwidth

Transfer size is the unit of memory that can be individually accessed, read and written.

- DRAM, EEPROM – byte addressable
- SSD, flash – block addressable (must read/write blocks)

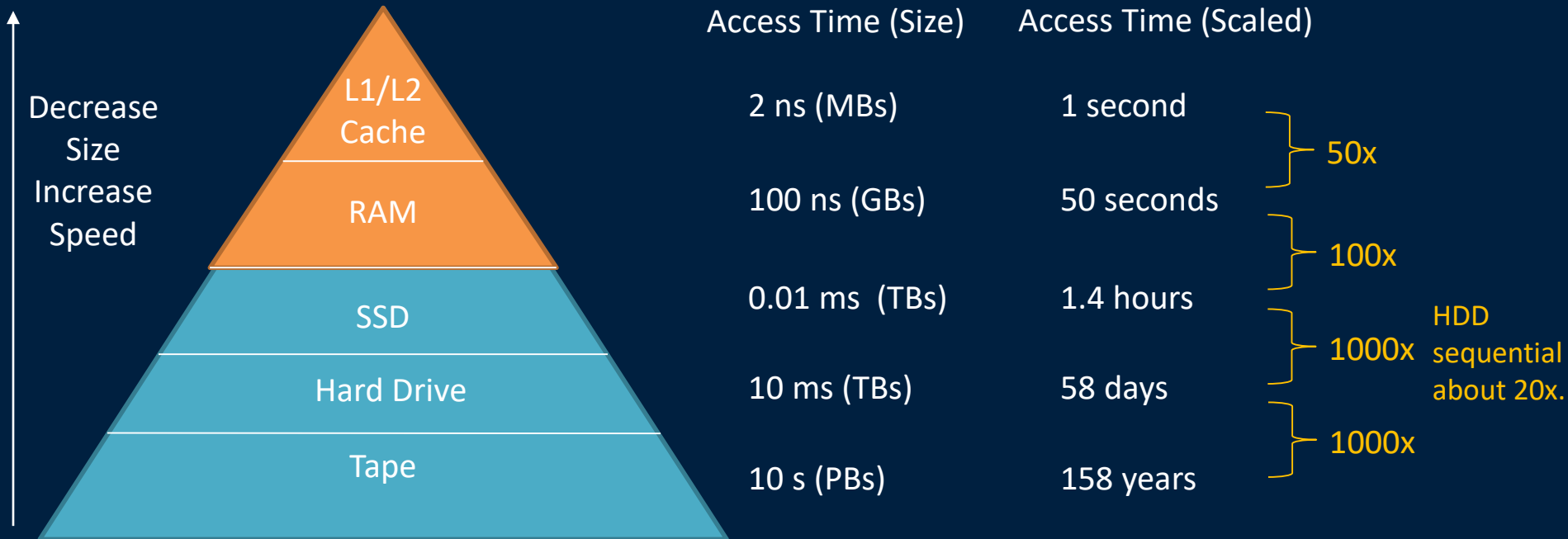
Latency is the time it takes for information to be delivered after the initial request is made.

Bandwidth is the rate at which information can be delivered.

Trends over time:

- **Capacity** grows at a fast rate. (e.g. double every few years)
- **Bandwidth** grows at a slower rate. (e.g. 5% per year)
- **Latency** does not improve very much.

Storage Hierarchy



Storage performance has dramatic effect!

Jim Gray – cache (this city), memory (nearby city), disk (Pluto)

Memory Devices

Dynamic Random Access Memory



Dynamic random access memory (DRAM) is general purpose, volatile memory.

- DRAM uses only one transistor and one capacitor per bit.
- DRAM needs periodic refreshing of the capacitor.

DRAM properties:

- low cost, high capacity
- volatile
- byte addressable
- latency ~ 10 ns
- bandwidth = 5 to 20 GB/s

Memory Devices

Processor Cache



Processor cache is faster memory storing recently used data that reduces the average memory access time.

- Cache is organized into lines/blocks of size from 64-512 bytes.
- Various levels of cache with different performance.

Cache properties:

- higher cost, very low capacity
- cache operation is hardware controlled
- byte addressable
- latency – a few clock cycles
- bandwidth – very high, limited by processor bus

Memory Devices

Flash Memory



Flash memory is used in many portable devices (cell phones, music/video players) and solid-state drives.

NAND Flash Memory properties:

- non-volatile
- low cost, high capacity
- block addressable
- asymmetric read/write performance: reads are fast, writes (which involve an erase) are slow
- erase limit of 1,000,000 cycles
- bandwidth (per chip): 40 MB/s (read), 20 MB/s (write)

Memory Devices

EEPROM



EEPROM (Electrically Erasable Programmable Read-Only Memory) is non-volatile and stores small amounts of data.

- Often available on small microprocessors.

EEPROM properties:

- non-volatile
- high cost, low capacity
- byte addressable
- erase limit of 1,000,000 cycles
- latency: 250 ns

Memory Devices

Solid State Drives



A **solid state drive** uses flash memory for storage.

Solid state drives have many benefits over hard drives:

- Increased performance (especially random reads)
- Better power utilization
- Higher reliability (no moving parts)

The performance of the solid state drive depends as much on the drive organization/controller as the underlying flash chips.

- Write performance is an issue and there is a large erase cost.

Solid state drives are non-volatile and block addressable like hard drives. The major difference is random reads are much faster (no seek time). This has a dramatic affect on the database algorithms used.

Memory Devices

Magnetic Tapes



Tape storage is non-volatile and is used primarily for backup and archiving data.

- Tapes are **sequential access** devices, so they are much slower than disks.

Since most databases can be stored in SSDs and RAID systems that support direct access, tape drives are now relegated to secondary roles as backup devices.

- Database systems no longer worry about optimizing queries for data stored on tapes.

"Tape is Dead. Disk is Tape. Flash is Disk. RAM Locality is King." – Jim Gray (2006), Microsoft/IBM, Turing Award Winner 1998 - For seminal contributions to database and transaction processing research and technical leadership in system implementation.

Device Performance Calculations

We will use simple models of devices to help understand the performance benefits and trade-offs.

These models are simplistic yet provide metrics to help determine when to use particular devices and their performance.

Memory Performance Calculations

Memory model will consider only transfer rate (determined from bus and memory speed). We will assume sequential and random transfer rates are the same.

Limitations:

- There is an advantage to sequential access compared to completely random access, especially with caching. Cache locality has a major impact as can avoid accessing memory.
- Memory alignment (4 byte/8 byte) matters.
- Memory and bus is shared by multiple processes.

Memory Performance Calculations

Example



A system has 8 GB DDR4 memory with 20 GB/sec. bandwidth.

Question 1: How long does it take to transfer 1 contiguous block of 100 MB memory?

$$\text{transfer time} = 100 \text{ MB} / 20,000 \text{ MB/sec.} = 0.005 \text{ sec} = \mathbf{5 \text{ ms}}$$

Question 2: How long does it take to transfer 1000 contiguous blocks of 100 KB memory?

$$\begin{aligned} \text{transfer time} &= 1000 * (100 \text{ KB} / 20,000,000 \text{ KB/sec.}) \\ &= 0.005 \text{ sec} = \mathbf{5 \text{ ms}} \end{aligned}$$

SSD Performance Calculations

SSD model will consider:

- **IOPS** – Input/Output Operations per Second (of given data size)
- latency
- bandwidth or transfer rate
- different performance for read and write operations.

Limitations:

- Write bandwidth is not constant. It depends on request ordering and volume, space left on drive, and SSD controller implementation.

SSD Performance Calculations

Examples



Question 1: A SSD has read bandwidth of 500 MB/sec. How long does it take to read 100 MB of data?

$$\text{read time} = 100 \text{ MB} / 500 \text{ MB/sec.} = \mathbf{0.2 \text{ sec}}$$

Question 2: The SSD IOPS for 4 KB write requests is 25,000. What is its effective write bandwidth?

$$\begin{aligned} \text{write bandwidth} &= 25,000 \text{ IOPS} * 4 \text{ KB requests} \\ &= 100,000 \text{ KB/sec.} = \mathbf{100 \text{ MB/sec.}} \end{aligned}$$

Device Performance

Question: What device would be the fastest to read 1 MB of data?

A) DRAM with bandwidth of 20 MB/sec.

B) SSD with read 400 IOPS for 100 KB data chunks.

Summary of Memory Devices

Memory Type	Volatile?	Capacity	Latency	Bandwidth	Transfer Size	Notes
DRAM	yes	High	Small	High	Byte	Best price/speed
Cache	Yes	Low	Lowest	Very high	Byte	Large reduction in memory latency
NAND Flash	No	Very high	Small	High	Block	Asymmetric read/write costs
EEPROM	No	Very low	Very small	High	Byte	High cost per bit ; On small CPUs.
Tape Drive	No	Very high	Very high	Medium	Block	Sequential access
Solid State Drive	No	Very high	Low	Medium	Block	Great random I/O ; Issue in write costs
Hard Drive	No	Very high	High	Medium	block	Beats SSDs by cost/bit but not by performance/cost.

RAID

Redundant Arrays of Independent Disks is a disk organization technique that utilizes a large number of inexpensive, mass-market disks to provide increased reliability, performance, and storage.

- Originally, the "I" stood for inexpensive as RAID systems were a cost-effective alternative to large, expensive disks. However, now performance and reliability are the two major factors.

Improvement of Reliability via Redundancy

RAID systems improve reliability by introducing **redundancy** as they store extra data that is used to recover from a disk failure.

- Redundancy occurs by duplicating data across multiple disks.
- **Mirroring** or **shadowing** duplicates an entire disk on another. Every write is performed on both disks, and if either disk fails, the other contains all the data.

By introducing more disks to the system the chance that some disk out of a set of **N** disks will fail is much higher than the chance that a specific single disk will fail.

- **Mean time to failure (MTTF)** – the average time the device is expected to run continuously without any failure.
- E.g., A system with 100 disks, each with MTTF of 100,000 hours (approx. 11 years), will have a system MTTF of 1000 hours (approx. 41 days).

Review: Parity

Parity is used for error checking. A parity bit is an extra bit added to the data. A single parity bit can detect one bit error.

In **odd parity** the number of 1 bits in the data plus the parity bit must be odd. In **even parity**, the number of 1 bits is even.

Example: What is the parity bit with *even parity* and the bit string: 01010010?

- Answer: The parity bit must be a 1, so that the # of 1's is even.

Parity Question

Question: What is the parity bit with *odd parity* and the bit string:
11111110?

A) 0

B) 1

C) 2

Improvement in Performance via Parallelism

The other advantage of RAID systems is increased **parallelism**. With multiple disks, two types of parallelism are possible:

- 1. Load balance multiple small accesses to increase throughput.
- 2. Parallelize large accesses to reduce response time.

Maximum transfer rates can be increased by allocating (**striping**) data across multiple disks then retrieving the data in parallel from the disks.

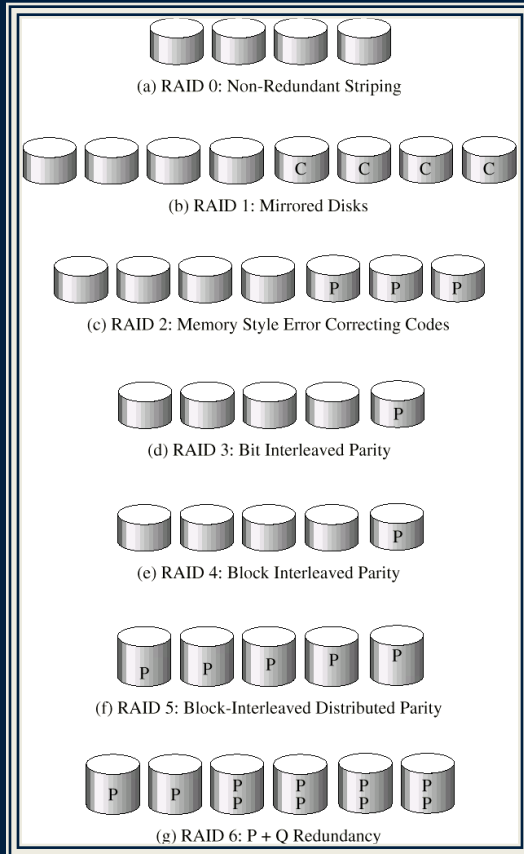
- **Bit-level striping** – split the bits of each byte across the disks
 - In an array of eight disks, write bit i of each byte to disk i .
 - Each access can read data at eight times the rate of a single disk.
 - But seek/access time worse than for a single disk.
- **Block-level striping** – with n disks, block i of a file goes to disk $(i \bmod n) + 1$

RAID Levels

There are different RAID organizations, or **RAID levels**, that have differing cost, performance and reliability characteristics:

- **Level 0:** Striping at the block level (non-redundant).
- **Level 1:** Mirrored disks (redundancy)
- **Level 2:** Memory-Style Error-Correcting-Codes with bit striping.
- **Level 3:** Bit-Interleaved Parity - a single parity bit used for error correction. Subsumes Level 2 (same benefits at a lower cost).
- **Level 4:** Block-Interleaved Parity - uses block-level striping, and keeps all parity blocks on a single disk (for all other disks).
- **Level 5:** Block-Interleaved Distributed Parity - partitions data and parity among all $N + 1$ disks, rather than storing data in N disks and parity in 1 disk. Subsumes Level 4.
- **Level 6:** P+Q Redundancy scheme - similar to Level 5, but stores extra info to guard against multiple disk failures.

RAID Levels Discussion



Level 0 is used for high-performance where data loss is not critical (parallelism).

Level 1 is for applications that require redundancy (protection from disk failures) with minimum cost.

- Level 1 requires at least two disks.

Level 5 is a common because it offers both reliability and increased performance.

- With 3 disks, the parity block for n th block is stored on disk $(n \bmod 3) + 1$. Do not have single disk bottleneck like Level 4.

Level 6 offers extra redundancy compared to Level 5 and is used to deal with multiple drive failures.

RAID Question

Question: What RAID level offers the high performance but no redundancy?

- A) RAID 0
- B) RAID 1
- C) RAID 5
- D) RAID 6

RAID Practice Question

Question: The capacity of a drive is 800 GB. Determine the capacity of the following RAID configurations:

- i) 8 drives in RAID 0 configuration
- ii) 8 drives in RAID 1 configuration
- iii) 8 drives in RAID 5 configuration

- | | | |
|----------------------|-------------|--------------|
| A) i) 6400 GB | ii) 3200 GB | iii) 5600 GB |
| B) i) 3200 GB | ii) 6400 GB | iii) 5600 GB |
| C) i) 6400 GB | ii) 3200 GB | iii) 6400 GB |
| D) i) 3200 GB | ii) 3200 GB | iii) 6400 GB |

RAID Summary

<u>Level</u>	<u>Performance</u>	<u>Protection</u>	<u>Capacity (for N disks)</u>
0	Best (parallel read/write)	Poor (lose all on 1 failure)	N
1	Good (write slower as 2x)	Good (have drive mirror)	$N / 2$
5	Good (must write parity block)	Good (one drive can fail)	$N - 1$
6	Good (must write multiple parity blocks)	Better (can have as many drives fail as dedicated to parity)	$N - X$ (where X is # of parity drives such as 2)

File Interfaces

Besides the physical characteristics of the media and device, how the data is allocated on the media affects performance (*file organization*).

The physical device is controlled by the operating system. The operating system provides one or more interfaces to accessing the device.

Block-Level Interface

A **block-level interface** allows a program to read and write a chunk of memory called a **block** (or **page**) from the device.

The page size is determined by the operating system. A page may be a multiple of the physical device's block size.

The OS maintains a mapping from logical page numbers (starting at 0) to physical blocks on the device.

Byte-Level Interface

A *byte-level interface* allows a program to read and write individually addressable bytes from the device.

A device will only directly support a byte-level interface if it is byte-addressable. However, the OS may provide a file-level byte interface to a device even if it is only block addressable.

File-Level Interface

A **file-level interface** abstracts away the device addressable characteristics and provides a standard byte-level interface for files to programs running on the OS.

A file is treated as a sequence of bytes starting from 0. File level commands allow for randomly navigating in the file and reading/writing at any location at the byte level.

Since a device may not support such access, the OS is responsible for mapping the logical byte address space in a file to physical device blocks. The OS performs buffering to hide I/O latency costs.

- Although beneficial, this level of abstraction may cause poor performance for I/O intensive operations.

Databases and File Interfaces

A database optimizes performance using device characteristics:

- 1) Favor reads over writes if device is slower for writes.
- 2) If device is block-addressable, read at the block level and know block boundaries.
- 3) If device has improved performance for large consecutive reads/writes, perform operations in sequential batches.
- 4) Put most frequently accessed data on fastest device.
- 5) If placement of data on the device matters, the database should control data placement.
- 6) The database will perform its own buffering separate from OS.

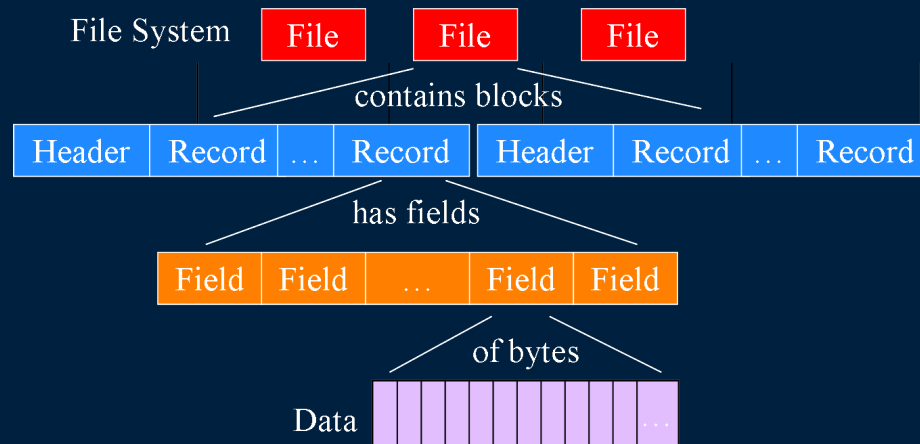
Representing Data in Databases

Overview



A **database** is made up of one or more files.

- Each **file** contains one or more blocks.
- Each **block** has a header and contains one or more records.
- Each **record** contains one or more fields.
- Each **field** is a representation of a data item in a record.



Representing Data in Memory

Every data item is translated into bytes for storage.

Example employee database:

- name : string (null-terminated or length at start)
- age : integer (binary)
- salary : double (IEEE 754 format)
- startDate : Date (string or long encoding)
- picture : BLOB (store length at start)

Representing Data in Memory

Strings and Characters



A **character** is represented by mapping the character symbol to a particular number.

- **ASCII** - maps characters/symbols to a number from 0 to 255.
- **UNICODE** - maps characters to a two-byte number (0 to 65,535) which allows for the encoding of larger alphabets.

A **string** is a sequence of characters allocated in consecutive memory bytes. A pointer indicates the location of the first byte.

- **Null-terminated string** - last byte value of 0 indicates end
- **Byte-length string** - length of string in bytes is specified (usually in the first few bytes before string starts).
- **Fixed-length string** - always the same size.

Representing Data in Memory

Dates



A **date** value can be represented in multiple ways:

- Integer representation - number of days past since a given date
 - Example: # days since Jan 1, 1900
- String representation - represent a date's components (year, month, day) as individual characters of a string
 - Example: YYYYMMDD or YYYYDDD
 - Please do not reinvent Y2K by using YYMMDD!!

A **time** value can also be represented in similar ways:

- Integer representation - number of seconds since a given time
 - Example: # of seconds since midnight
- String representation - hours, minutes, seconds, fractions
 - Example: HHMMSSFF

Representing Nothing (NULL)

The SQL concept of `NULL` (nothing, no value) may be represented as:

- Special sentinel value in a field
 - Note: Cannot use valid data for the field (such as 0 for integer or empty string)
- Boolean flag in record
- May not represent at all for records with variable schema (i.e. absence of value for field indicates `NULL`)



Storing Records in Memory

A **record** consists of one or more fields grouped together.

- Each tuple of a relation in the relational model is a record.

Two main types of records:

- **Variable-length records** - the size of the record varies.
- **Fixed-length records** - all records have the same size.

Separating Fields of a Record

The fields of a record can be separated in multiple ways:

- 1) **No separator** - store length of each field, so do not need a separate separator (fixed length field).
 - Simple but wastes space within a field.
- 2) **Length indicator** - store a length indicator at the start of the record (for the entire record) and a size in front of each field.
 - Wastes space for each length field and need to know length beforehand.
- 3) **Use offsets** – at start of record store offset to each field
- 4) **Use delimiters** - separate fields with delimiters such as a comma (comma-separated files).
 - Must make sure that delimiter character is not a valid character for field.
- 5) **Use keywords** - self-describing field names before field value (XML and JSON).
 - Wastes space by using field names.

Schemas

A *schema* is a description of the record layout.

A schema typically contains the following information:

- names and number of fields
- size and type of each field
- field ordering in record
- description or meaning of each field

Fixed versus Variable Formats

If every record has the same fields with the same types, the schema defines a *fixed record format*.

- Relational schemas generally define a fixed format structure.

It is also possible to have no schema (or a limited schema) such that not all records have the same fields or organization.

- Since each record may have its own format, the record data itself must be *self-describing* to indicate its contents.
- XML and JSON documents are considered self-describing with variable schemas (*variable record formats*).

Schemas

Fixed Format Example



Employee record is a fixed relational schema format:

<u>Field Name</u>	<u>Type</u>	<u>Size in Bytes</u>
name	char(10)	10
age	integer	4
salary	double	8
startDate	Date	8 (YYYYMMDD)

Example record:

- Joe Smith, 35, \$50,000, 1995/05/28

Memory allocation:



in ASCII? 00000023 in IEEE 754? in ASCII?

Fixed Format with Variable fields

It is possible to have a fixed format (schema), yet have variable sized records.

- In the Employee example, the picture field is a BLOB which will vary in size depending on the type and quality of the image.

It is not efficient to allocate a set memory size for large objects, so the fixed record stores a pointer to the object and the size of the object which have fixed sizes.

The object itself is stored in a separate file or location from the rest of the records.

Variable Formats

XML and JSON



XML:

```
<employees>
  <employee>
    <name>Joe Smith</name>
    <age>35</age>
    <salary>50000</salary>
    <hire>1995/05/28</hire>
  </employee>
  <employee>
    <name>CEO</name>
    <age>55</age>
    <hire>1994/06/23</hire>
  </employee>
</employees>
```

JSON:

```
{
  "employees": [
    { "name": "Joe Smith",
      "age": 35,
      "salary": 50000,
      "hire": "1995/05/28" },
    { "name": "CEO",
      "age": 55,
      "hire": "1994/06/23" }
  ]
}
```



Variable Format Discussion

Variable record formats are useful when:

- The data does not have a regular structure in most cases.
- The data values are sparse in the records.
- There are repeating fields in the records.
- The data evolves quickly so schema evolution is challenging.

Disadvantages of variable formats:

- Waste space by repeating schema information for every record.
- Allocating variable-sized records efficiently is challenging.
- Query processing is more difficult and less efficient when the structure of the data varies.

Format and Size Question

Question: JSON and XML are best described as:

- A)** fixed format, fixed size
- B)** fixed format, variable size
- C)** variable format, fixed size
- D)** variable format, variable size

Relational Format and Size Question

Question: A relational table uses a VARCHAR field for a person's name. It can be best described as:

- A) fixed format, fixed size
- B) fixed format, variable size
- C) variable format, fixed size
- D) variable format, variable size

Fixed versus Variable Formats

Discussion



We have seen fixed length/fixed format records, and variable length/variable format records.

1) Do fixed format and variable length records make sense?

Yes, you can have a fixed format schema where certain types have differing sizes. BLOBs are one example.

2) Do variable format and fixed length records make sense?

Surprisingly, Yes. Allocate a fixed size record then put as many fields with different sizes as you want and pad the rest.

320587	Joe Smith	SC	95	3	← Padding →
184923	Kathy Li	EN	92	3	← Padding →
249793	Albert Chan	SC	94	3	← Padding →

Research Question

CHAR versus VARCHAR



Question: We can represent a person's name in MySQL using either CHAR(50) or VARCHAR(50). Assume that the person's name is 'Joe'. How much space is actually used?

- A)** CHAR = 3 ; VARCHAR = 3
- B)** CHAR = 50 ; VARCHAR = 3
- C)** CHAR = 50 ; VARCHAR = 4
- D)** CHAR = 50 ; VARCHAR = 50



Storing Records in Blocks

Now that we know how to represent entire records, we must determine how to store sets of records in blocks.

There are several issues related to storing records in blocks:

- 1) **Separation** - how do we separate adjacent records?
- 2) **Spanning** - can a record cross a block boundary?
- 3) **Clustering** - can a block store multiple record types?
- 4) **Splitting** - are records allocated in multiple blocks?
- 5) **Ordering** - are the records sorted in any way?
- 6) **Addressing** - how do we reference a given record?

Storing Records in Blocks

Separation



If multiple records are allocated per block, we need to know when one record ends and another begins.

Record *separation* is easy if the records are a fixed size because we can calculate the end of the record from its start.

Variable length records can be separated by:

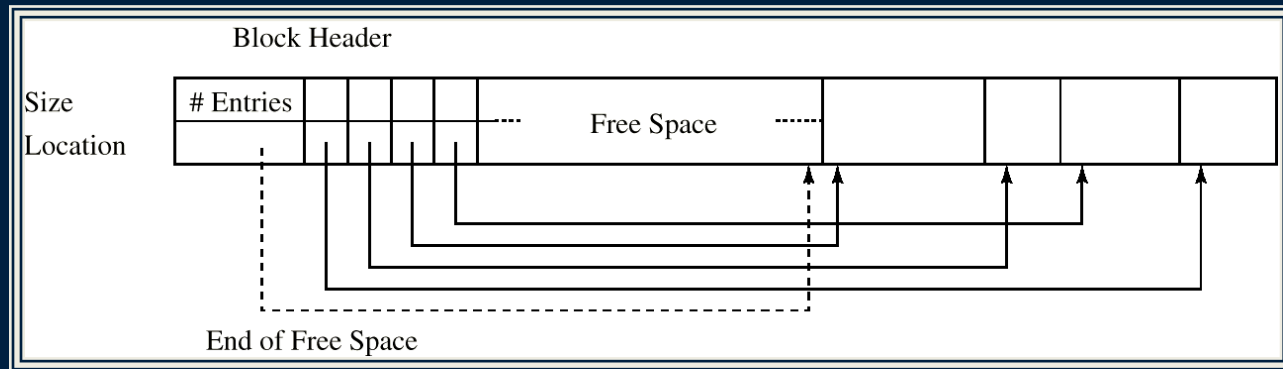
- 1) Using a special separator marker in the block.
- 2) Storing the size of the record at the start of each record.
- 3) Store the length or offset of each record in the block header.

Variable Length Records Separation and Addressing



A **block header** contains the number of records, the location and size of each record, and a pointer to block free space.

Records can be moved around within a block to keep them contiguous with no empty space between them. Header is updated accordingly.



Storing Records in Blocks

Spanning



If records do not exactly fit in a block, we have two choices:

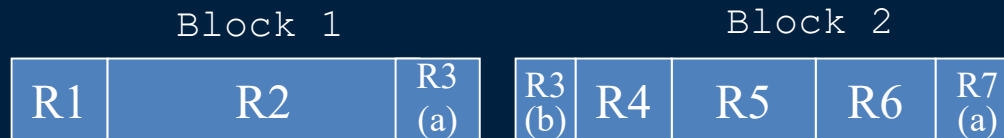
- 1) Waste the space at the end of each block.
- 2) Start a record at the end of a block and continue on the next.

Choice #1 is **unspanned** option. Simple - do not allocate across blocks.



Choice #2 is **spanned** option.

- Each piece must have a pointer to its other part.
 - Spanning is required if the record size is larger than the block size.



Storing Records in Blocks

Spanning Example



If the block size is 4096 bytes, the record size is 2050 bytes, and we have 1,000,000 records:

- How many blocks are needed for spanned/unspanned records?
- What is the block (space) utilization in both cases?

Answer:

- Unspanned
 - put one record per block implies 1,000,000 blocks
 - each block is only $2050/4096 * 100\% = 50\%$ full (utilization = 50%)
- Spanned
 - all blocks are completely full except the last one
 - # of blocks required = $1,000,000 * 2050 / 4096 = 500,489$ blocks
 - utilization is almost 100%

Storing Records in Blocks

Clustering



Clustering is allocating records of different types together on the same block (or same file) because they are frequently accessed together.

Example:

- Consider creating a block where a department record is allocated together with all employees in the department:

Block 1

DPT1	EMP1	EMP2	DEPT2	EMP3	EMP4
------	------	------	-------	------	------

Storing Records in Blocks

Clustering (2)



If the database commonly processes queries such as:

```
select * from employee, department
where employee.deptId = department.Id
```

then the clustering is beneficial because the information about the employee and department are adjacent in the same block.

However, for queries such as:

```
select * from employee
```

```
select * from department
```

clustering is harmful because the system must read in more blocks, as each block read contains information that is not needed to answer the query.

Storing Records in Blocks

Split Records



A **split record** is a record where portions of the record are allocated on multiple blocks for reasons other than spanning.

- Record splitting may be used with or without spanning.

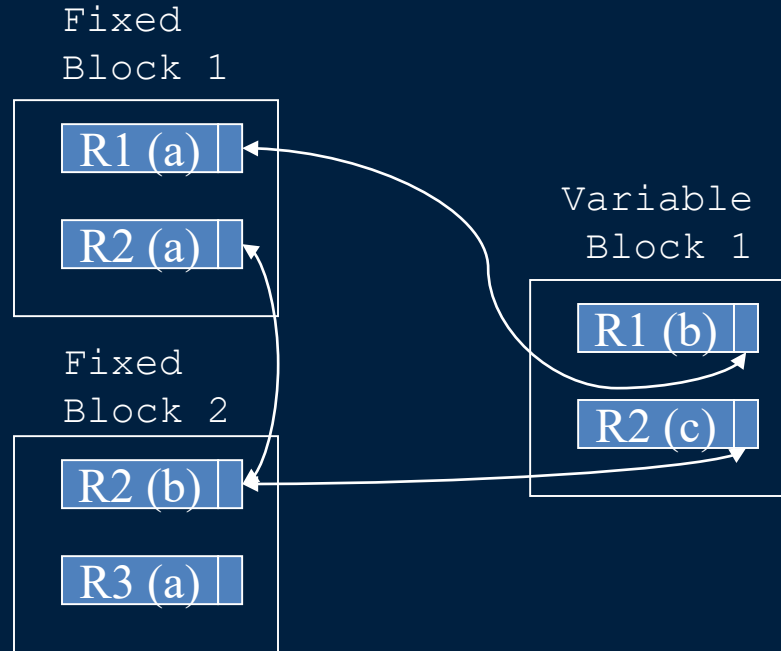
Typically, hybrid records are allocated as split records:

- The **fixed portion** of the record is allocated on one block (with other fixed record portions).
- The **variable portion** of the record is allocated on another block (with other variable record portions).

Splitting a record is done for efficiency and simplifying allocation. The fixed portion of a record is easier to allocate and optimize for access than the variable portion.

Storing Records in Blocks

Split Records with Spanning Example



Storing Records in Blocks

Ordering Records



Ordering (or sequencing) records is when the records in a file (block) are sorted based on the value of one or more fields.

Sorting records allows some query operations to be performed faster including searching for keys and performing joins.

Records can either be:

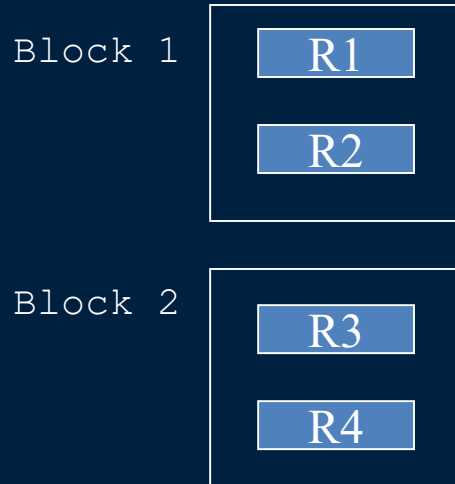
- 1) *physically ordered* - the records are allocated in blocks in sorted order.
- 2) *logically ordered* - the records are not physical sorted, but each record contains a pointer to the next record in the sorted order.

Storing Records in Blocks

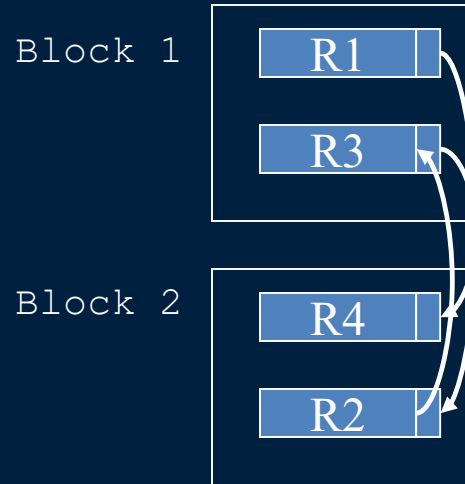
Ordering Records Example



Physical ordering



Logical Ordering



What are the tradeoffs between the two approaches?

What are the tradeoffs of any ordering versus unordered?

Storing Records in Blocks

Addressing Records



Addressing records is a method for defining a unique value or address to reference a particular record.

Records can either be:

- 1) **physically addressed** - a record has a physical address based on the device where it is stored.
 - A physical disk address may use a sector # or a physical address range exposed by the device.
- 2) **logically addressed** - a record that is logically addressed has a key value or some other identifier that can be used to lookup its physical address in a table.
 - Logical addresses are indirect addresses because they provide a mechanism for looking up the actual physical addresses. They do not provide a method for locating the record directly on the device.
 - E.g. OS provides logical block to physical sector mapping for files.

Storing Records in Blocks

Addressing Records Tradeoff

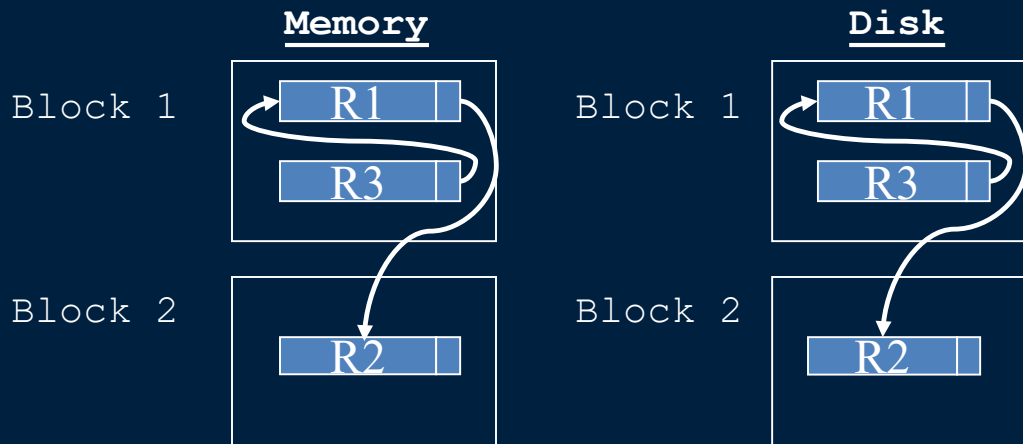


There is a tradeoff between physical and logical addressing:

- Physical addresses have better *performance* because the record can be accessed directly (no lookup cost).
- Logical addresses provide more *flexibility* because records can be moved on the physical device and only the mapping table needs to be updated.
 - The actual records or fields that use the logical address do not have to be changed.
 - Easier to move, update, and change records with logical addresses.

Pointer Swizzling

When transferring blocks between the disk and memory, we must be careful when handling pointers in the blocks. For example:



Pointer swizzling is the process for converting disk pointers to memory pointers and vice versa when blocks move between memory and disk.

Operations on Files

Once data has been stored to a file consisting of blocks of records, the database system will perform operations such as update and delete to the stored records.

How records are allocated and addressed affects the performance for update and delete operations.

Operations on Files

Record Deletion



When a record is deleted from a block, we have several options:

- 1) Reclaim deleted space
 - Move another record to the location or compress file.
- 2) Mark deleted space as available for future use

Tradeoffs:

- Reclaiming space guarantees smaller files, but may be expensive especially if the file is ordered.
- Marking space as deleted wastes space and introduces complexities in maintaining a record of the free space available.

Operations on Files

Issues with Record Deletion



We must also be careful on how to handle references to a record that has been deleted.

- If we re-use the space by storing another record in the same location, how do we know that the correct record is returned or indicate the record has been deleted?

Solutions:

- 1) Track down and update all references to the record.
- 2) Leave a "tombstone" marker at the original address indicating record deletion and not overwrite that space.
 - Tombstone is in the block for physical addressing, in the lookup table for logical addressing.
- 3) Allocate a unique record id to every record and every pointer or reference to a record must indicate the record id desired.
 - Compare record id of pointer to record id of record at address to verify correct record is returned.

Research Question

PostgreSQL VACUUM



Question: What does the **VACUUM** command do in PostgreSQL?

- A)** Cleans up your dirty house for you
- B)** Deletes records from a given table
- C)** Reclaims space used by records marked as deleted
- D)** Removes tables no longer used

Operations on Files

Record Insertion



Inserting a record into a file is simple if the file is not ordered.

- The record is **appended** to the end of the file.

If the file is physically ordered, then all records must be shifted down to perform insert.

- **Extremely costly operation!**

Inserting into a logically ordered file is simpler because the record can be inserted anywhere there is free space and linked appropriately.

- **However, a logically ordered file should be periodically re-organized to ensure that records with similar key values are in nearby blocks.**

Memory and Buffer Management

Memory management involves utilizing buffers, cache, and various levels of memory in the memory hierarchy to achieve the best performance.

- A database system seeks to minimize the number of block transfers between the disk and memory.

A **buffer** is a portion of main memory available to store copies of disk blocks.

A **buffer manager** is a subsystem responsible for allocating buffer space in main memory.

Buffer Manager Operations

All read and write operations in the database go through the buffer manager. It performs the following operations:

- **read block B** – if block B is currently in buffer, return pointer to it, otherwise allocate space in buffer and read block from disk.
- **write block B** – update block B in buffer with new data.
- **pin block B** – request that B cannot be flushed from buffer
- **unpin block B** – remove pin on block B
- **output block B** – save block B to disk (can either be requested or done by buffer manager to save space)

Key challenge: How to decide which block to remove from the buffer if space needs to be found for a new block?

Buffer Management Replacement Strategy



A *buffer replacement strategy* determines which block should be removed from the buffer when space is required.

- Note: When a block is removed from the buffer, it must be written to disk if it was modified.

Some common strategies:

- Random replacement
- Least recently used (LRU)
- Most recently used (MRU)

Buffer Replacement Strategies and Database Performance



Operating systems typically use least recently used for buffer replacement with the idea that the past pattern of block references is a good predictor of future references.

However, database queries have well-defined access patterns (such as sequential scans), and a database system can use the information to better predict future references.

- LRU can be a bad strategy for certain access patterns involving repeated scans of data!

Buffer manager can use statistical information regarding the probability that a request will reference a particular relation.

- E.g., The schema is frequently accessed, so it makes sense to keep schema blocks in the buffer.

Research Question

MySQL Buffer Management



Question: What buffer replacement policy does MySQL InnoDB use?

- A) LRU
- B) MRU
- C) 2Q

Column Storage

Each file represents all the data for a column. A file entry contains the column value and a record id. Records are rebuilt by combining columns using the record id.

Row Format

1	Ed	50	M
2	Ann	30	F
3	Fred	45	M
4	Jane	80	F

Column Format

1	Ed	50	M
2	Ann	30	F
3	Fred	45	M
4	Jane	80	F

The column format reduces the amount of data retrieved from disk (as most queries do not need all columns) and allows for better compression.

Research Question

PostgreSQL Column Layout



Question: Does PostgreSQL support column layout?

A) Yes

B) No

PAX Layout

PAX storage format combines row and column storage:

- Store all data about a row in a single disk page but
- organized data by column inside each page
- Allows for better cache locality when processing a column

Row Format

1	Ed	50	M
2	Ann	30	F
3	Fred	45	M
4	Jane	80	F

Column Format

1	Ed	50	M
2	Ann	30	F
3	Fred	45	M
4	Jane	80	F

PAX Format

1	2	Ed	Ann
50	30	M	F
3	4	Fred	Jane
45	80	M	F

Issues in Disk Organizations

There are many ways to organize information on a disk. Key concepts:

- **Locality**: items tend to be accessed together
- **Searchability**: how fast can records be found

The "best" disk organization will be determined by a variety of factors such as: **flexibility**, **complexity**, **space utilization**, and **performance**.

Performance measures to evaluate a given strategy include:

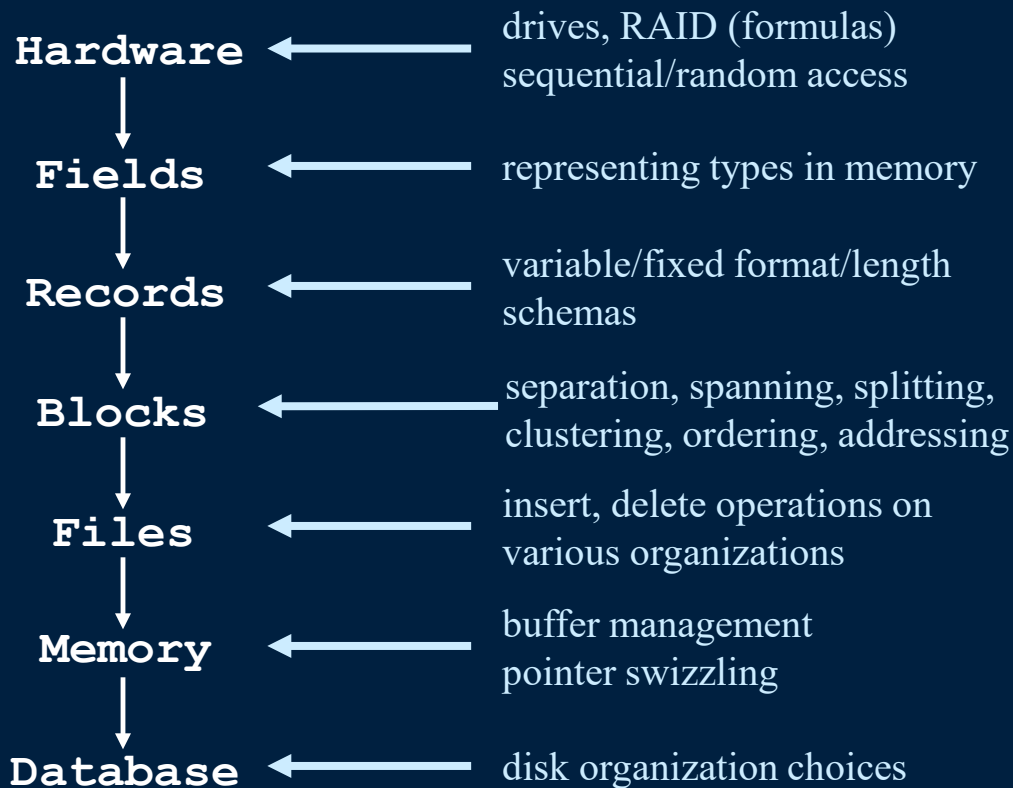
- space utilization
- expected times to search for a record given a key, search for the next record, insert/append/delete/update records, reorganize the file, read the entire file.

Key terms:

- **Storage structure** is a particular organization of data.
- **Access mechanism** is an algorithm for manipulating the data in a storage structure.

Summary

Storage and Organization



Major Objectives

The "One Things":

- Perform device calculations such as computing transfer times.
- Explain the differences between fixed and variable schemas.
- List and briefly explain the six record placement issues in blocks.

Major Theme:

- There is no single correct organization of data. The "best" organization will be determined by a variety of factors such as: *flexibility, complexity, space utilization, and performance.*

Objectives

- Compare/contrast volatile versus non-volatile memory.
- Compare/contrast random access versus sequential access.
- Use both metric and binary units for memory sizes.
- List the benefits of RAID and common RAID levels.
- List different ways for representing strings in memory.
- List different ways for representing date/times in memory.
- Explain the difference between fixed and variable length records.
- Compare/contrast the ways of separating fields in a record.
- Define and explain the role of schemas.
- Compare/contrast variable and fixed formats.
- List and briefly explain the six record placement issues in blocks.

Objectives (2)

- Explain the tradeoffs for physical/logical ordering and addressing.
- List the methods for handling record insertion/deletion in a file.
- List some buffer replacement strategies.
- Explain the need for pointer swizzling.
- Define storage structure and access mechanism.



THE UNIVERSITY OF BRITISH COLUMBIA

