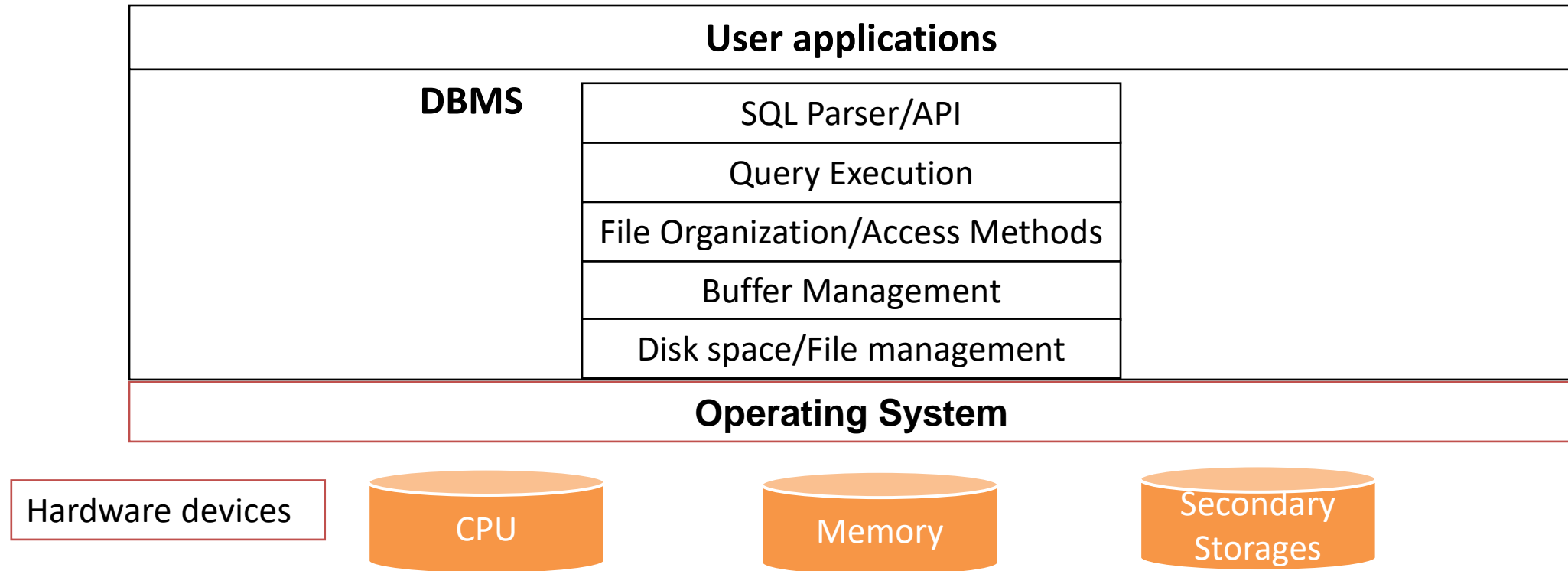


CSE462/562: Database Systems (Spring 22)

Lecture 2: Physical Storage

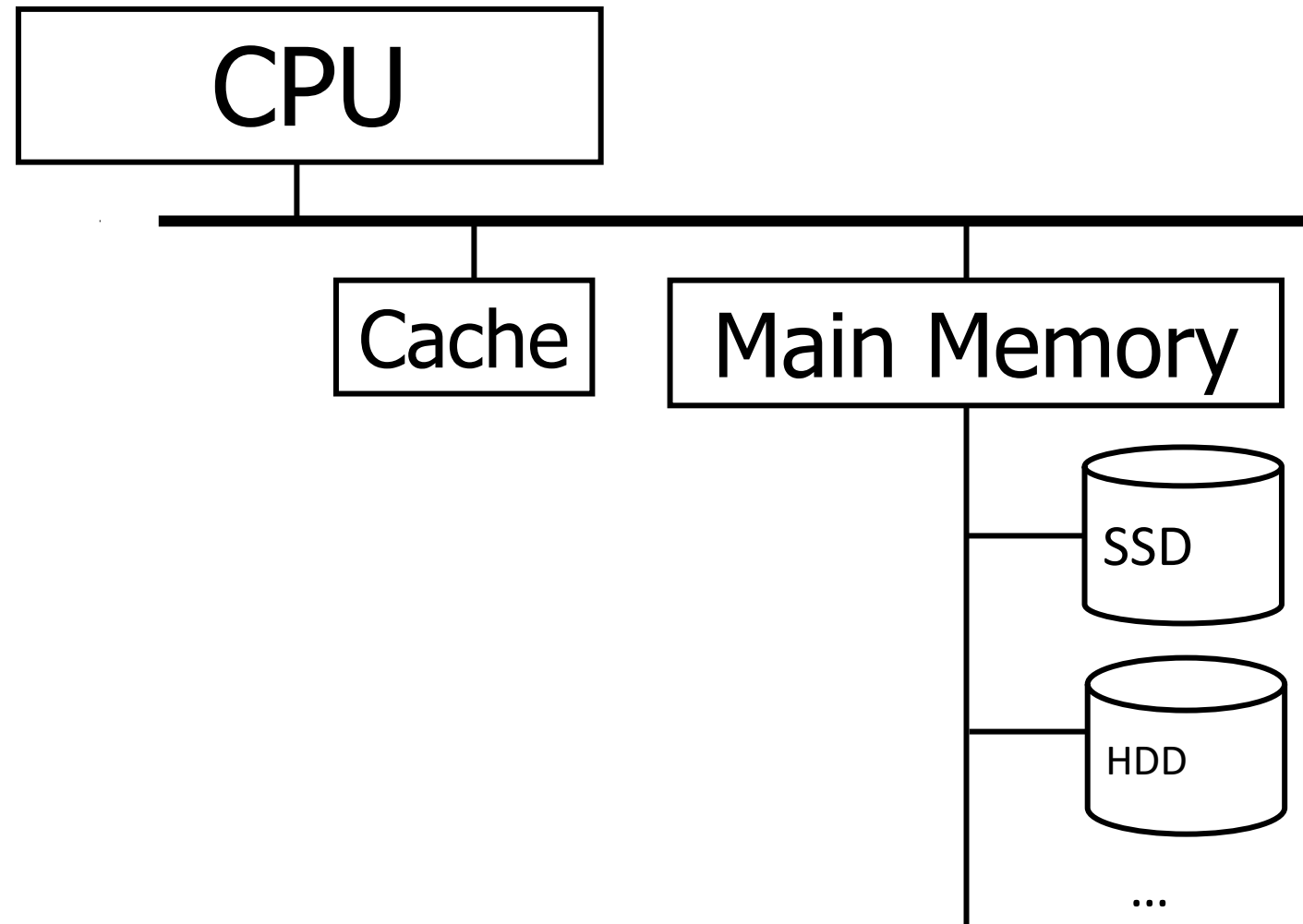
2/2/2022

Big Picture



Typical (& oversimplified) computer architecture

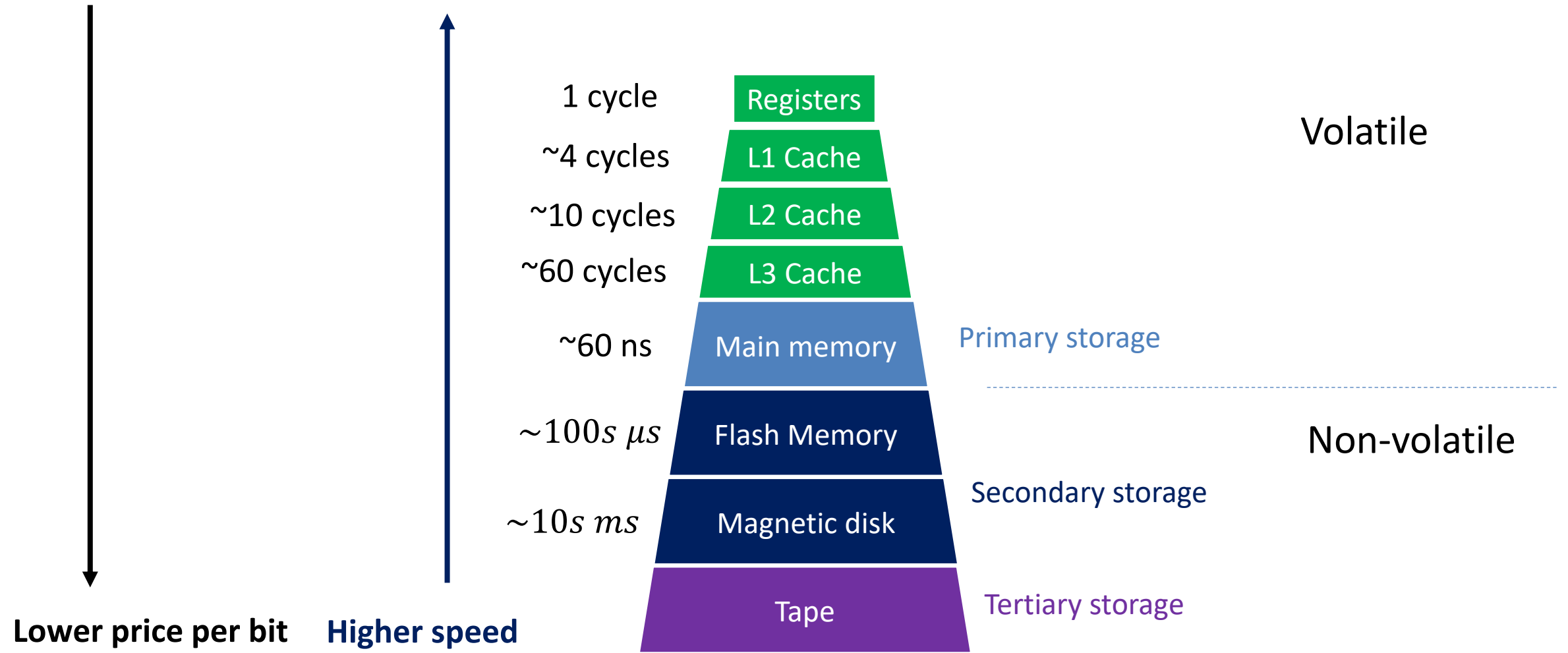
- A simplistic view of a computer



Typical
Computer

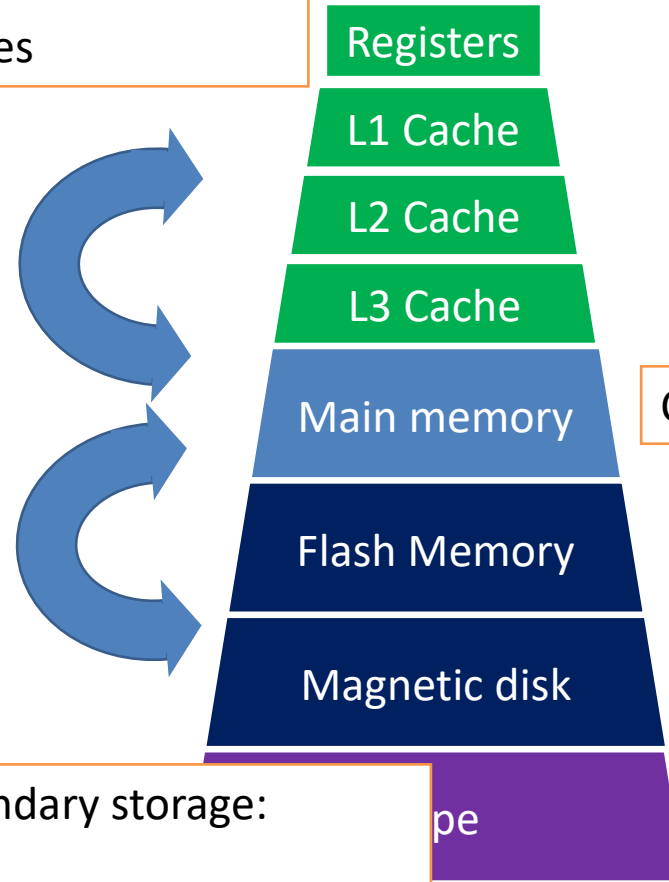
Secondary
Storage

Storage Hierarchy



Data Transfers

Between cache and main memory:
hardware/OS controlled
usually in small units of cache lines



Volatile

CPU operates on main memory (byte addressable)

Non-volatile

Between main memory and secondary storage:
DBMS controlled (read/write)
usually with large block I/O

Volatile storage

- Register
 - Very fast but very limited amount
 - CPU directly operates on registers
- Cache
 - Faster than main memory but takes multiple cycles to access
 - Stores cache lines that are likely to be read/write again
 - Usually managed by CPU
- Main memory
 - Still quite fast albeit it takes hundreds of cycles
 - CPU instructions can read/write byte addressable data into/from registers

Why not store everything in memory?

- Too expensive
 - Data growth is much faster than what you can afford
- Volatile
 - Power loss -> data loss
- Typical storage hierarchy in (traditional) DBMS
 - Main memory as buffer/working space
 - Disk as the main database storage
 - Tape for archiving old data
- Main memory DB actually uses memory for main database storage
 - Persistency of data? [Logging & checkpointing \(later lectures\)](#)

Non-volatile storage

- Common non-volatile (secondary) storage
 - Flash memory (e.g., SSD)
 - Magnetic disk
- Advantages
 - Cheaper -- can store much more data than memory with the same cost
 - Non-volatile – data are saved in server shutdown/power failure
- Disadvantages
 - Block device: read/write in the units of sectors (usually 512B/4096B)
 - Higher latency: usually $\geq 1 - 2$ orders of magnitude slower than main memory
- Tertiary storage: tape (sequential I/O only)
 - Very slow but inexpensive; good for archiving data

Closer look at non-volatile storage

- We need to know the performance characteristics of non-volatile storage
 - to optimize database storage design



Magnetic disk (HDD)

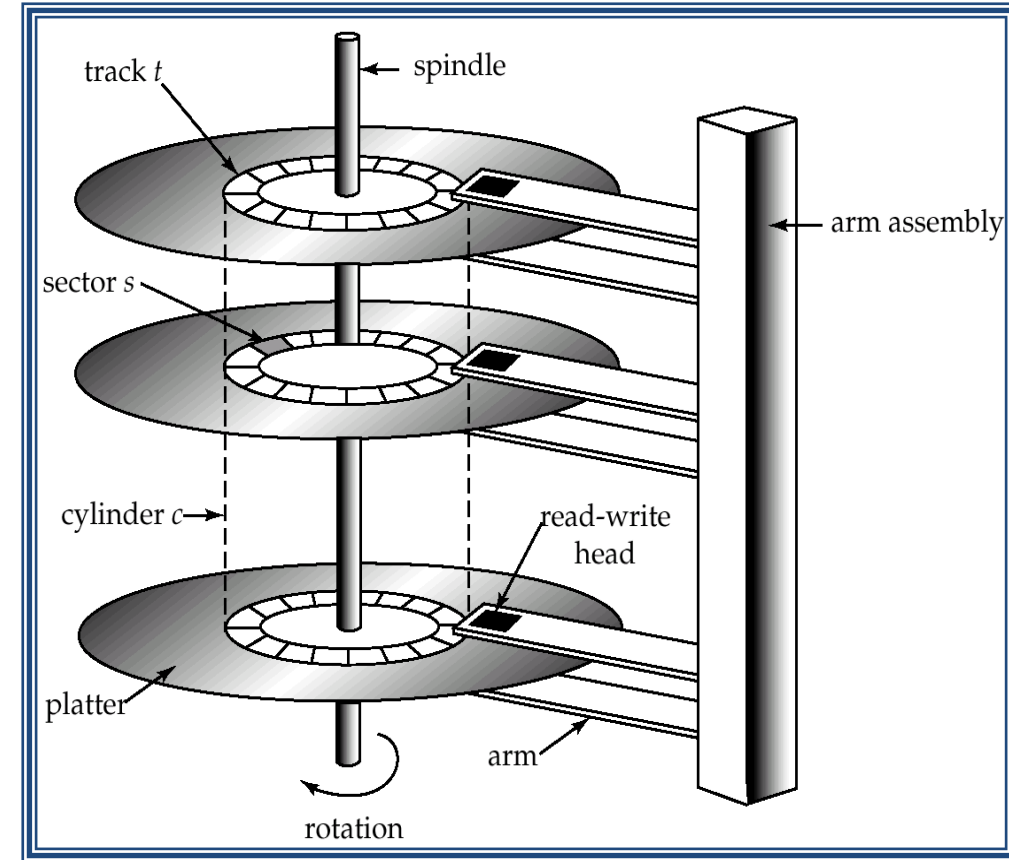


[This Photo](#) by Unknown Author is licensed under [CC BY](#)

Solid State Drive (SSD)

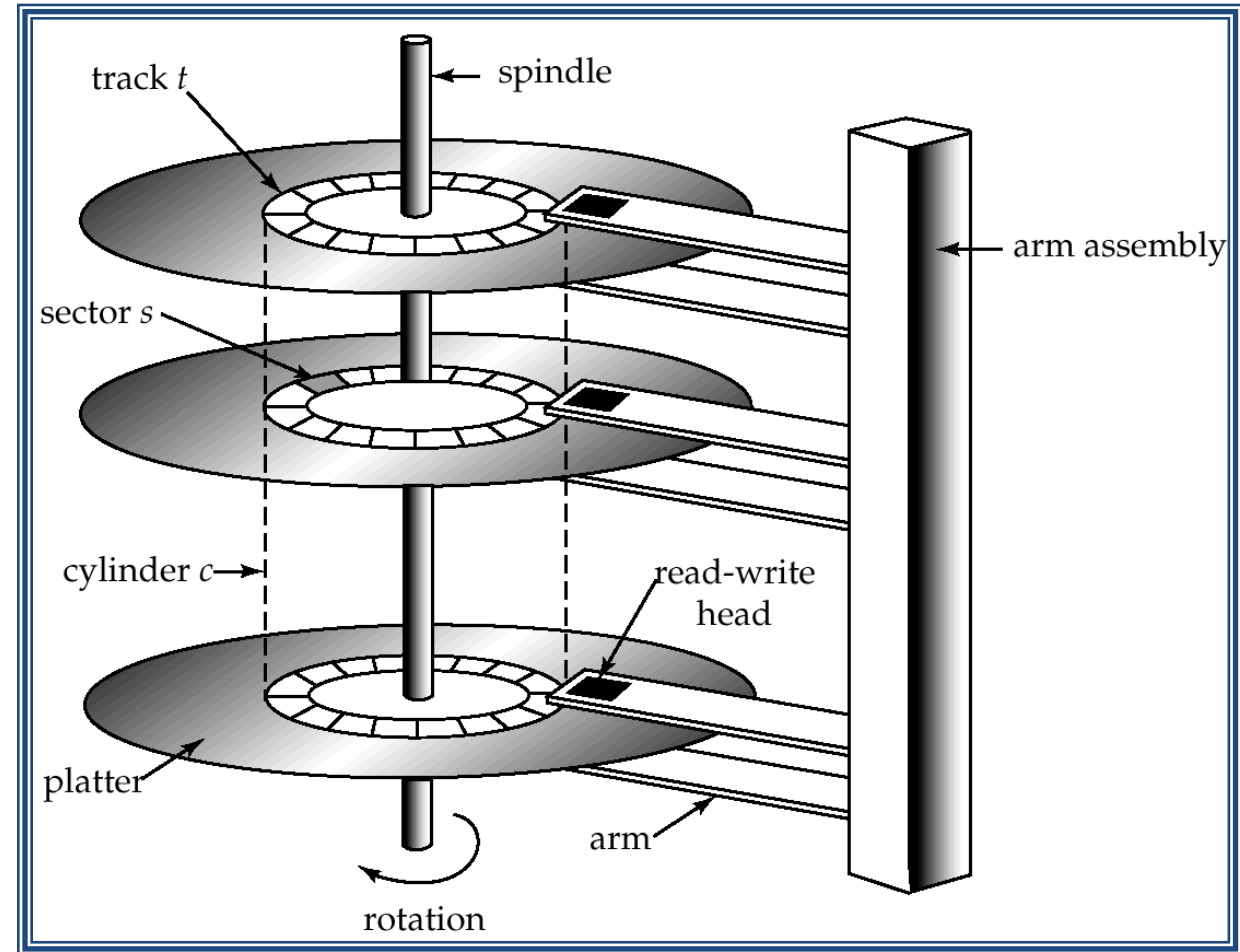
Magnetic disk organization

- Multiple platters
 - Each platter has **two** surfaces for data storage
 - Platters spin at the **same** rate (e.g., 7200 rpm)
 - A ring on a surface is called a **track**
 - A track is divided into many **sectors** of fixed size (512 B)
 - A sector is the **smallest** unit of I/O
- A single arm assembly with multiple disk heads
 - Can only move inward/outward **together**
 - The vertical stack of tracks is called a **cylinder**
 - Disk heads can be over the tracks of the **same** cylinder at the **same** time
 - Usually one read/writes at the same time
- Address of a sector: **cylinder - head - sector**
 - (0, 0, 0) : first sector; (0, 0, 1): second sector, ...
 - (0, 1, 0) : the S^{th} sector, (1, 0, 0) the $(SH)^{th}$ where S is the max # of sectors/track and H is the # of heads
 - Reality: today's disks use logical block addressing (linear **block #**)
 - Translated to the actual geometry by disk controller
 - Nevertheless, this is still a good model for understanding HDD performance.



Magnetic disk I/O latency

- File systems perform I/O in units of multiple sector (page)
 - 4KB~16KB are most common
- Break-down of I/O latency of a page
 - **Seek time:** moving arms to the cylinder
 - 2 ~ 20 ms per seek
 - 4 ~ 10 ms on average
 - **Rotation delay:**
wait for the sector to be under a head
 - Depending on rotation speed (5400 rpm - 15000 rpm)
 - E.g, 7200 rpm = 120 rotations/second
 $\Rightarrow 1/120 = 8.33 \text{ ms / rotation}$
on average it needs a half rotation
 $\Rightarrow 8.33 / 2 = 4.17 \text{ ms on average}$
 - **Transfer time:** time for reading/writing data
 - Data transfer rate: 50 - 200 MB/s
 - $\Leftrightarrow 0.02 \sim 0.08 \text{ ms for 4KB pages}$
- **Average access time**
 - 4KB page, 7200 rpm: roughly 8 ~ 15 ms

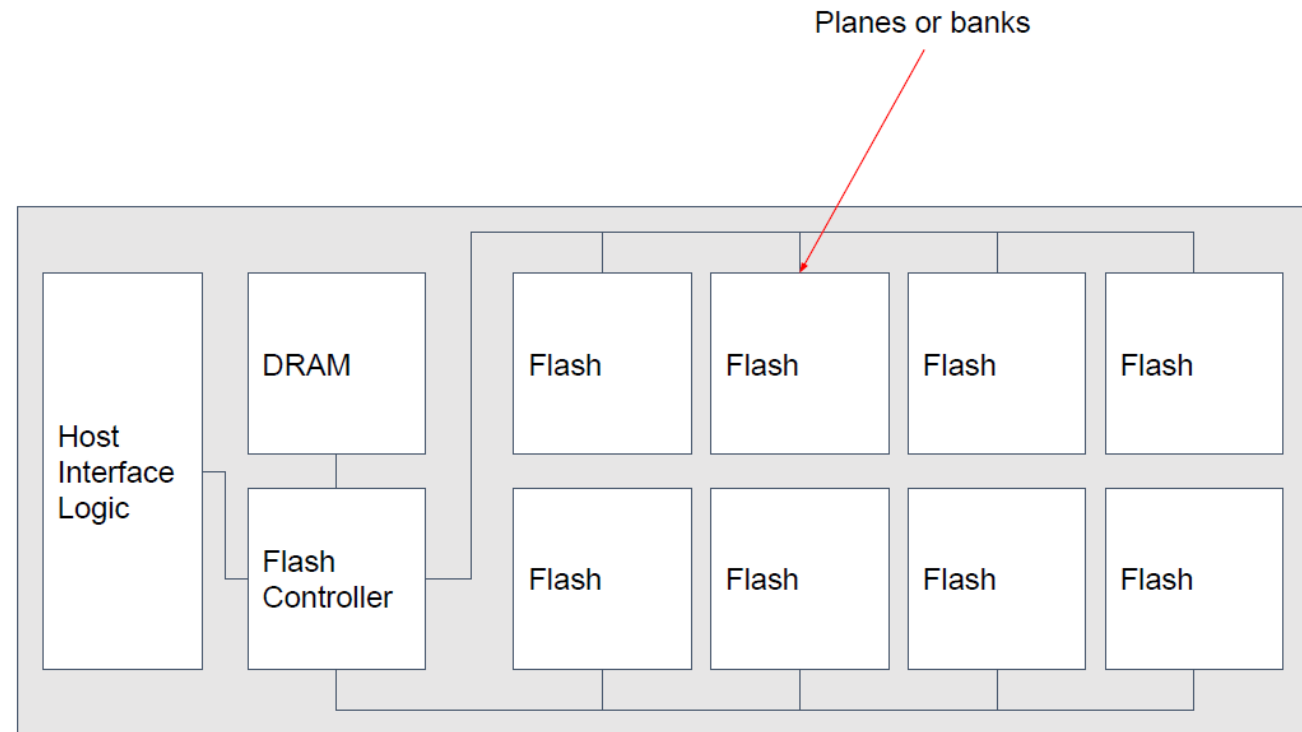


Impact of I/O pattern on magnetic disk

- I/O pattern has a huge impact on I/O performance
 - E.g., 4KB page size
 - Sequential read/write: usually 100 ~ 200+ MB/s
 - Random read/write: 50 ~ 200 IOPS \Leftrightarrow 200 KB ~ 800 KB /s
 - **> 2** orders of magnitude difference in terms of data transfer rate
- Rule of thumb:
 - Random I/O: very slow; avoid reading a lot of data from random location
 - Sequential I/O: better for accessing a lot of data

Flash memory / solid state drive

- NAND Flash is the most common storage media for solid state drives
- No mechanical parts (magnetic disk can have head crash => data corruption/loss)
 - More reliable; less likely to fail due to physical shocks
- Faster than magnetic disk



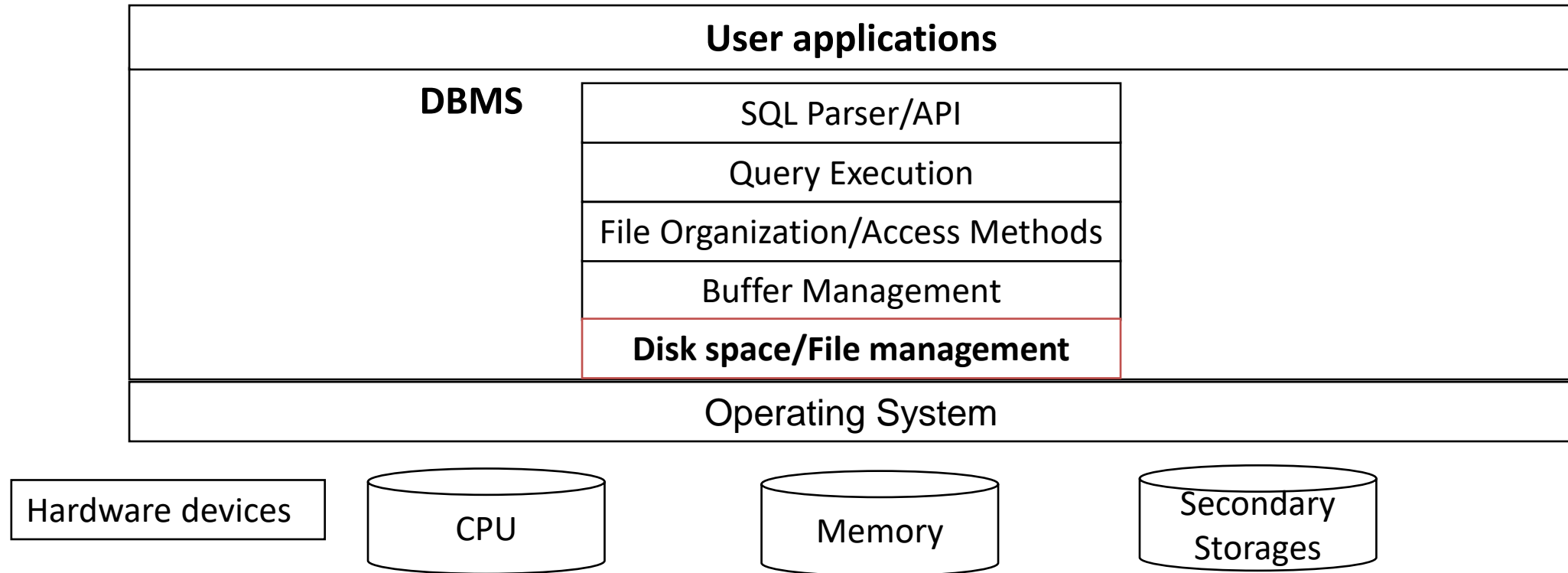
Flash memory / solid state drive

- NAND SSD has asymmetric read/write performance
 - 4KB page, typical SSD internal performance numbers
 - Read latency: 20 to 100 μs ; throughput: > 500 MB/s
 - Write latency: 200 μs ; throughput: > 500 MB/s
 - Erase latency: ~ 2 ms
- Three ops: read/write/erase
 - Read/write works on pages (usually 4KB)
 - Write can only change some bits from 1 to 0 (not the other way around!)
 - Must erase before write a page.
 - Erase works on blocks (e.g., 256 KB)
 - Resets all bits in a block to 1
 - Flash translation layer: indirection of page numbers to physical pages
 - Solves two problems: slow erase and flash wear
- Actual performance also often bound by peripheral bus's bandwidth and IOPS

Flash memory / solid state drive

- NAND SSD has asymmetric read/write performance
 - The performance from DB stand of view?
 - No single answer depending on how you use it
 - I/O queue depth, I/O api, access pattern, page size, peripheral bus type and etc.
 - In a typical case:
 - Sequential I/O is still preferred, although random I/O isn't as bad as in HDD
 - SSDs have much better random I/O performance than magnetic disk
 - 10k - 1M IOPS
 - and higher bandwidth as well
 - up to 7GB/s on PCIe 4.0, ~500MB/s on SATA

Big Picture

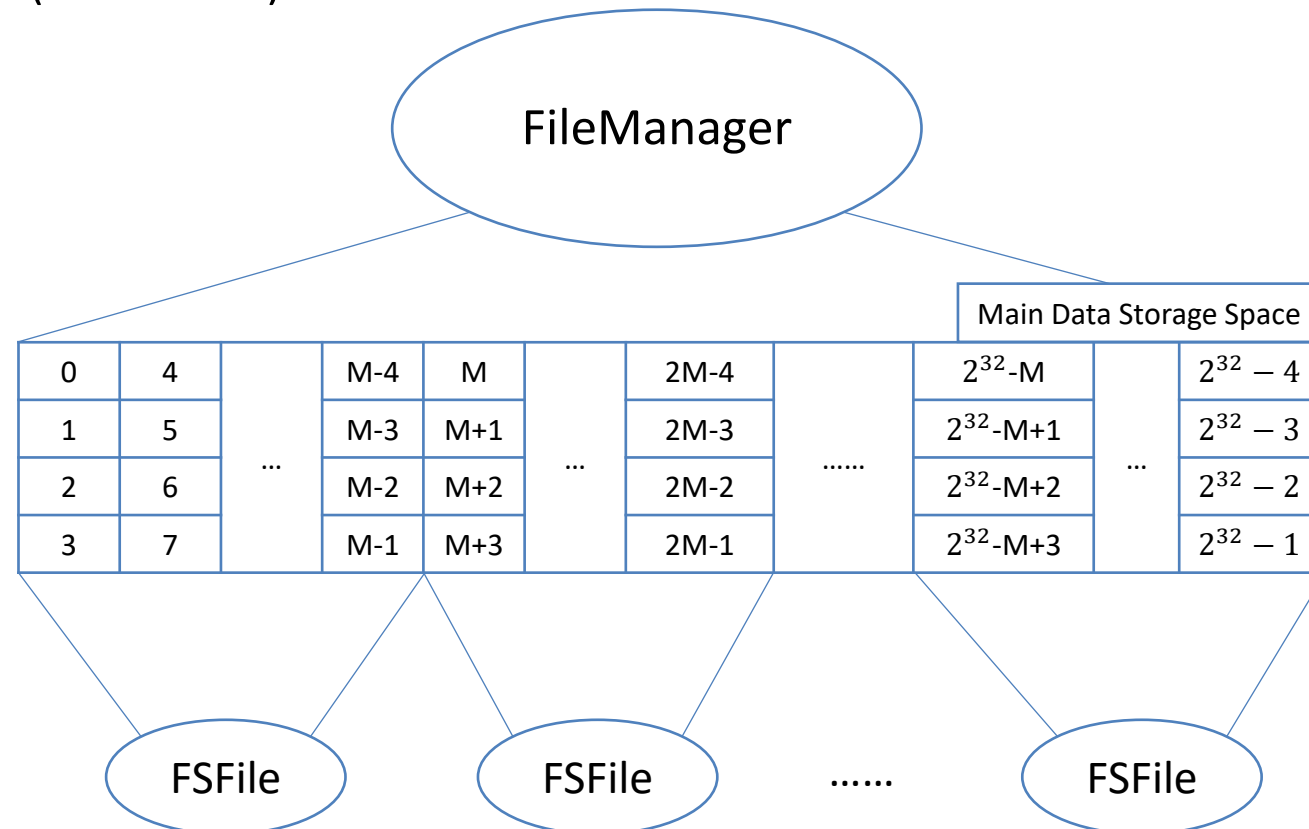


Disk Space Management

- Lowest layer of DBMS software manages space on disk
 - Disk space is usually organized in *pages*
 - which may not necessarily directly be mapped to disk sectors/file system pages!
 - common choices are 4KB, 8KB, 16KB, etc.
 - Using the OS file system or not? Some do and some don't!
 - Even with file system
 - How to organize pages (in one file/multiple files)?
 - How to deal with concurrency/recovery?
 - ...
- Higher levels call upon this layer to:
 - allocate/de-allocate a page
 - read/write a page
- Best if a request for a sequence of pages is satisfied by pages stored sequentially on disk!
 - Responsibility of disk space manager.
 - Higher levels don't know how this is done, or how free space is managed.
 - Though they may assume sequential access for files!
 - Hence, disk space manager should do a decent job.

Disk Space Management in course project Taco-DB

- A flat main data storage page from page 0 to page $2^{32} - 1$
 - Stored as 64GB files on the local file system;
 - FileManager manages many (virtual) files -- (not FSFile)
 - Each is a double-linked list of pages, allocated in groups of 64 consecutive pages
 - Each file maintains its own free list
 - Concurrency? Recovery? (to be done)



Summary

- This lecture
 - Storage hierarchy and storage devices
 - Disk space management
- Next lecture
 - Buffer management
 - File organization in DBMS